



OPEN

DATA DESCRIPTOR

# NEBULA101: an open dataset for the study of language aptitude in behaviour, brain structure and function

Alessandra Rampinini<sup>1,2,6</sup>✉, Irene Balboni<sup>1,2,3,4,5,6</sup>, Olga Kepinska<sup>3,4</sup>, Raphael Berthele<sup>2,5</sup> & Narly Golestani<sup>1,2,3,4</sup>

This paper introduces the “NEBULA101 - Neuro-behavioural Understanding of Language Aptitude” dataset, which comprises behavioural and brain imaging data from 101 healthy adults to examine individual differences in language and cognition. Human language, a multifaceted behaviour, varies significantly among individuals, at different processing levels. Recent advances in cognitive science have embraced an integrated approach, combining behavioural and brain studies to explore these differences comprehensively. The NEBULA101 dataset offers brain structural, diffusion-weighted, task-based and resting-state MRI data, alongside extensive linguistic and non-linguistic behavioural measures to explore the complex interaction of language and cognition in a highly multilingual sample. By sharing this multimodal dataset, we hope to promote research on the neuroscience of language, cognition and multilingualism, enabling the field to deepen its understanding of the multivariate panorama of individual differences and ultimately contributing to open science.

## Background & Summary

**Individual differences in language.** This paper describes the “NEBULA101 - Neuro-behavioural Understanding of Language Aptitude” dataset. The dataset collects behavioural and brain imaging data of 101 healthy adults for the study of individual differences in language and cognition.

Human language is a complex behaviour, and crucial to understanding its workings is the fact that individuals differ in the way they manifest this and other cognitive skills<sup>1</sup>. The science of individual differences has expanded in recent years, thanks to a more integrated, less modular take on cognition: by combining the study of behaviour and the brain in a deep phenotyping approach, mindful of individual differences, researchers can gain a comprehensive understanding of complex cognitive functions, language included<sup>2</sup>.

**Language aptitude.** One aspect of language shown to display large individual differences is language *aptitude*. Language aptitude was originally proposed to explain why some people display remarkable abilities when learning additional languages<sup>3–5</sup>. We use the term *additional* language for any language that is not the (or one of the) individual’s *first* language(s), i.e., language(s) to which they were exposed from birth. The term “additional” thus includes *second languages* (e.g. in the context of migration), *foreign languages* (e.g. in classroom learning), or *third/additional languages* of multilinguals who use more than two.

Initially, researchers viewed language aptitude as a stable trait, comprising phonetic coding, grammatical sensitivity, inductive learning, and rote learning abilities. These skills were seen as componential, in the climate of the first cognitive revolution, where the brain–mind–behaviour interaction was seen as the execution of algorithms operating within cognitive modules<sup>6</sup>: in this view, phonetic coding was at the core of production and

<sup>1</sup>Brain and Language Lab, Department of Psychology, Faculty of Psychology and Education Science, University of Geneva, Geneva, Switzerland. <sup>2</sup>National Centre of Competence in Research Evolving Language, Swiss National Science Foundation, Switzerland. <sup>3</sup>Brain and Language Lab, Vienna Cognitive Science Hub, University of Vienna, Vienna, Austria. <sup>4</sup>Department of Behavioural and Cognitive Biology, Faculty of Life Sciences, University of Vienna, Vienna, Austria. <sup>5</sup>Institute of Multilingualism, University of Fribourg, Fribourg, Switzerland. <sup>6</sup>These authors contributed equally: Alessandra Rampinini, Irene Balboni. ✉e-mail: [alessandra.rampinini@unige.ch](mailto:alessandra.rampinini@unige.ch)

perception of speech sounds, grammatical sensitivity underlay the capacity to identify structure in language (i.e. morphosyntax), inductive learning supported the generalisation of language rules from the input, and rote learning skills the construction of vocabulary via routinary use. Research has since suggested a more parsimonious structure, combining grammatical sensitivity and inductive learning into a global language analytic ability<sup>7</sup>.

Nowadays, with advances in neuroscience and experimental psychology methods, the language aptitude construct overall still describes a set of skills operating *across* the hierarchy of all language components, from lower to higher levels of complexity. What has changed is the way we conceptualise language itself, which has also changed the way we view language aptitude: in the second wave of the cognitive revolution, thanks to usage-based linguistics and neural-network psychology, language is seen as a complex adaptive system, rather than a set of rigid structures and fixed operations<sup>8</sup>, and one of the interacting components of human cognition, rather than an isolated function. Language aptitude is part of this dynamic: it can vary with age<sup>9–11</sup> or with multilingual experience (likely increasing meta-linguistic awareness<sup>12–14</sup>, with patterns yet to be clarified<sup>15,16</sup>), and ultimately, with cognition more generally. In this context, we view the brain as a network of interrelated functions giving rise to complex behaviours, including language: to this end, this dataset comprises extensive general and domain-specific cognition tests, together with brain measures. The underlying concept of an integrated mind, arising from an integrated brain, is at the core of our methodological choice and of our first exploration of these data with graph theoretical methods<sup>17</sup>.

Mindful that any quantitative analysis will always require some degree of operationalisation (i.e. to derive scores from tests, we need to identify components of scientific constructs that such tests might tap into), our position is that this view of language as a dynamic system branching out and connecting to more general mechanisms holds promise for better understanding language itself and the human cognitive system more generally.

**The importance of research on the multilingual brain.** The views expressed above are pivotal to the study of multilingualism and its behavioural and brain dynamics. In today's globalised society multilingualism and multicompetence (i.e. using and knowing multiple languages) are becoming normalised<sup>18</sup>. There is, however, a theoretical issue with defining a sociolinguistic construct such as that of any “-lingualism”, due to its multidimensionality<sup>19</sup>, without incurring in stereotypes (e.g. the definition of *native-speakerism*) or perpetrating problematic practices, such as that of selecting certain profiles<sup>20</sup> or matching “the unmatchable” based on intrinsically tricky parameters (e.g. the number of languages<sup>21</sup>). Nonetheless, we need to better define these concepts, as they constitute important characteristics of our present-day society. Knowing and using multiple languages demands a fundamental cognitive (re)organisation<sup>18</sup>, with several psycho-neurobiological correlates that are somewhat hard to reconcile in a comprehensive view<sup>22</sup>. Therefore, we must seek to imagine language competence and use in newer, more naturalistic, and multidimensional ways to better understand their influence on the brains (and lives) of language users.

To this aim, we believe that multimodal datasets with rich phenotypical information, such as the one presented here, are a step forward in this direction. Language aptitude, viewed as embedded in (and interacting with) cognition might be one of the engines driving multilingualism (in its many forms), and understanding its underpinnings might ultimately influence the way we view, use, and even teach language(s). In this context, Switzerland holds a special status as a country with four official languages and as a destination for expatriates from all over the world, including users of Minority, Indigenous, Non-standard, and Dialect (MIND) varieties. These languages have often been disregarded not only by researchers but by users themselves, who might not even recognise their multilingual status when knowing an additional MIND language<sup>20</sup>. However, these *do* appear in the multicompetent panorama of many of our participants, when asked explicitly (see Fig. 2).

These issues have recently been tackled in work calling for a more diverse view of cognitive science in general<sup>23</sup> and neurolinguistics<sup>24</sup> in particular, and underlining the contribution *non-English* (Romance<sup>25</sup>, non-Indo-European<sup>24</sup> as well as MIND<sup>20,21</sup>) languages to the field. The ongoing discourse on language policies and teaching<sup>26–28</sup>, as well as the thriving field of instructed language learning research, especially in the Swiss context<sup>29</sup>, are just a part of the puzzle. What is still lacking in the psycho- and neurolinguistics of multilingualism is a naturalistic perspective rooted in the brain itself, also likely due to the intrinsic difficulty of quantifying language use in a nonparametric and dynamic way. Language varies at all levels, challenging us to consider the remarkable plasticity of advanced human abilities by harnessing diversity as a tool for advancing cognitive science<sup>30</sup>. Thus, while leveraging French knowledge and fluency in our participants - a condition necessary for performance comparability - this dataset also captures and documents their linguistic diversity and multicompetence. This comprehensive documentation has the potential to facilitate investigations into how such diversity relates to both the phenotypical (behavioural) and endophenotypical (brain structural and functional) characteristics of individuals.

**Aptitude for language(s) and individual differences in other domains of cognition.** The construct of language aptitude was developed in the domain of foreign language learning, as explained, but the idea that aptitude only manifests in foreign languages has now been surpassed, since individual differences can be observed in first language skills too, even if it is harder to pinpoint them and isolate them from experiential factors<sup>31</sup>. Moreover, such individual differences might co-exist more globally with individual differences in other domains of cognition<sup>17,29</sup>, giving rise to “neurocognitive profiles” of language aptitude involving the mnemonic domain, fluid reasoning, auditory abilities, and even musicality<sup>32</sup>. It is therefore relevant to ask ourselves whether language, in this modernised and dynamic view, is part of a positive manifold<sup>33,34</sup> originating from the beneficial interactions between cognitive processes, as we proposed in a recent analysis including these data<sup>17</sup>. Both the stable and malleable (or *plastic*) features of the human cognitive system are fundamental to such interactions.

**Stability and malleability, predisposition and experience.** This dataset presents cross-sectional data. Nonetheless, to understand the nature of the language (aptitude) construct, it is important to consider it as the result of a complex interaction between stable traits and malleable states, the first ascribed to genetic predisposition and the latter to environmental influences (what we generally define as “experience”)<sup>35–37</sup>. While we provide a *snapshot* of individual profiles at a given moment in time, in order to formulate relevant questions on (and interpretations of) these data, it is important to remember that the observed measures arise from both stable traits and from malleable skills, or states<sup>38</sup>. Cognitive activity<sup>39–42</sup>, including language learning and use<sup>43–46</sup>, and brain *plasticity* are intrinsically related. Changes in regional activity, network connectivity, or morphology, arising from underlying molecular, neurobiochemical and other changes<sup>39</sup>, subtend this state of malleability in brain function and structure<sup>47</sup>, both during development and in skill learning. However, brain functional and structural architecture is also highly polygenic (i.e. controlled by the complex interaction of multiple genes, which in turn are expressed in multiple variants across individuals), and in addition, different cortical loci and tissue features (thickness, surface area) are affected to different degrees by genetics<sup>48</sup>. Further, genetic factors also likely influence the degree to which neuroplasticity manifests<sup>49</sup>.

**The open science of language aptitude.** Language aptitude has come a long way since its study was confined to the foreign language classroom, and research restricted to military and government access<sup>50</sup>. The construct now encompasses some of the most promising avenues for research on the human cognitive system more generally: the roles of predisposition and experience, the nature of neuroplasticity and the integrated and multivariate organisation of cognitive domains. Because of its relevance for studying questions on language and cognition more generally (as well as investigating the role of the environment, e.g. the experience of multilingualism), this development in the conceptualisation of language aptitude is of particular interest to the ever-growing world of Findable, Accessible, Interoperable, and Reusable (FAIR) neuroscience data<sup>51</sup>, which has faced a substantial growth in the last 15 years<sup>52</sup>, and possibly as many challenges<sup>53,54</sup>. FAIR principles are accelerating our comprehension of the human brain<sup>55</sup>. Concurrently, the existence of standardised protocols such as the Brain Imaging Data Structure (BIDS)<sup>56</sup> supports this recent strive towards data and code accessibility, ensuring that data is organized and its analysis reproducible, enabling more efficient and effective use of shared datasets and streamlining data preparation, thus allowing scientists to focus more on discovery.

Mindful that the panorama of shared neuroimaging data is vast and that datasets might be scattered around repositories of the commercial and institutional type, focusing on either a specific population (clinical, paediatric) or modality<sup>57</sup>, we searched for openly available data similar to ours (MRI and behavioural data from healthy adult participants focusing on language) on OpenNeuro<sup>58</sup>, one of the most recent, easy to access and growing neuroimaging databases. At the time of writing, a search on <http://openneuro.org> for BIDS-valid MRI raw/derivative datasets with more than 50 healthy adults, including the keyword ‘language’ and accompanied by a published or pre-registered data descriptor, yields the following entries (Table 1), with two of the datasets (ds004215, ds000243) being listed but having no relationship with language (and thus not being included in the below table).

A similar, broad search on Google Scholar for data descriptors associated with openly available language datasets (“open AND mri dataset AND language”) yielded more datasets: the MOUS (Mother Of Unification Studies)<sup>59</sup>, a dataset comprising sMRI, DWI, resting-state and task-based fMRI and MEG in 204 participants assigned to either sentence reading or listening; the Alice dataset<sup>60</sup>, where 75 participants listened to the same chapter of Alice in Wonderland during fMRI or during EEG. A third dataset presents a large quantity of data collected in fewer participants with state-of-the-art technology: this is the high-resolution, 7 T fMRI Forrest Gump database<sup>61</sup>, during which participants watched the ‘Forrest Gump’ movie. This dataset can be used to study naturalistic language processing, even though the study did not primarily focus on language. Then, the search yielded the LanA dataset, a probabilistic language atlas derived from brain data in more than 800 individuals<sup>62</sup>. Finally, the search also yielded two speech production datasets of vocal tract MRI<sup>63,64</sup>.

At the behavioural level there is more heterogeneity in the type and quantity of shared data<sup>57</sup>. The origins of such heterogeneity have been discussed for quite some time. In a Nature Neuroscience perspective article from ten years ago, the challenges linked to standardisation and accessibility of behavioural data were already discussed<sup>65</sup>. Importantly, behaviour was defined as a complex, highly dimensional, dynamic, and interconnected phenomenon without distinct separable scales, and this was discussed as being one of the leading causes for lack of standardisation and scarce FAIR compliance. The authors insisted on its foundational and unifying nature, and called for improved standards: “Behaviour (...) is the principal function of the brain. (...) Copious, quantitative and open behavioural data has the potential (...) to solidify the foundations of other [disciplines], including neuroscience”<sup>65</sup> (p.1455). However, even in the cited datasets, which represent timely and relevant steps towards accessibility of neuroscience data, the phenotypical (behavioural) information being shared beyond demographics is still relatively less prominent than the endophenotypical (neural) data: oftentimes, only data from the behavioural tasks being performed during fMRI are available. Moreover, even when behaviour was tested outside the scanner, if the original project did not focus on *both* phenotypical and endophenotypical data and on a specific topic (such as, in the cited examples, pragmatics<sup>66</sup> or reading<sup>67</sup>), the accompanying behavioural information is scarce. Two notable, recent exceptions in the field of individual differences in language are represented by a behavioural dataset including 33 measures from 112 adult Dutch speakers<sup>68</sup>, and by recent work by Berthele and colleagues in children<sup>69</sup>, both representing an important milestone for shared behavioural data.

Given the lack of standardisation in the way we administer behavioural tasks outside the scanner, compared to fMRI task delivery, it seems daunting to force the structure of behavioural paradigms and log files produced by a plethora of software and online platforms to accommodate information within the structure required by the BIDS standard, a difficulty that we encountered in our work, as well. Some initiatives, such as the Behavior

OpenNeuro Accession number	Descriptor reference	Modality	Number of participants	Complete dataset (all participants-measures)	Associated tasks/measures	Main topic
ds003481 <sup>116</sup>	A Dataset to Study Pragmatic Language and Its Underlying Cognitive Processes <sup>66</sup>	sMRI, fMRI, behaviour	145	No	<ul style="list-style-type: none"> <li>• Comprehension of metaphorical phrases and proverbs</li> <li>• Recognition and emotional categorisation of speech acts</li> <li>• Lexical-semantic processing</li> </ul>	Language pragmatics
ds004765 <sup>117</sup>	Relationship between resting state functional connectivity and reading-related behavioural measures in 69 adults <sup>118</sup>	sMRI, fMRI (resting-state), DWI	69	Yes	<ul style="list-style-type: none"> <li>• Word/nonword reading</li> <li>• Spelling</li> <li>• Lexical decision</li> <li>• Spoonerisms</li> <li>• Rapid automatised naming</li> <li>• Non-word repetition</li> <li>• Vocabulary</li> </ul>	Reading
ds002382 <sup>119</sup>	Age-related differences in auditory cortex activity during spoken word recognition <sup>120</sup>	sMRI, fMRI	61	Yes	<ul style="list-style-type: none"> <li>• Word listening</li> <li>• Word repetition</li> </ul>	Word processing
ds004285 <sup>121</sup>	Listening task*	sMRI, fMRI	78	Yes	<ul style="list-style-type: none"> <li>• Word repetition</li> </ul>	Word processing
ds004073 <sup>122</sup>	Comparing language lateralisation using fMRI and fTCD <sup>123</sup>	fMRI, functional transcranial doppler sonography (fTCD)	51	No (minimally missing information)	<ul style="list-style-type: none"> <li>• Word generation</li> <li>• Sentence generation</li> <li>• Phonological decision</li> <li>• Word comprehension</li> <li>• Sentence comprehension</li> <li>• Syntactic decision</li> </ul>	Language lateralisation
ds001747 <sup>124</sup>	Exploring the Resting State Neural Activity of Monolinguals and Late and Early Bilinguals <sup>125</sup>	sMRI, fMRI (resting-state)	92	Yes	<ul style="list-style-type: none"> <li>• L1 and L2 (where applicable) proficiency</li> <li>• Language background</li> </ul>	Bilingualism
ds001796 <sup>126</sup>	Bilingualism and the brain**	sMRI, fMRI (task-based), fMRI (resting-state), DWI	64	Yes	Flanker task	Bilingualism
ds002345 <sup>127</sup>	Narratives: fMRI data for evaluating models of naturalistic language comprehension <sup>128</sup>	sMRI, fMRI	345	No	Naturalistic story listening	Language perception
ds003643 <sup>129</sup>	Le Petit Prince: A multilingual fMRI corpus using ecological stimuli <sup>130</sup>	sMRI, fMRI	112	Yes	Naturalistic story listening	Language perception

**Table 1.** Available neuroimaging datasets on language with more than 50 healthy participants, validated in BIDS. (\*) Possibly including or revising ds002382 but no information on dataset version is provided. (\*\*) No data descriptor available.

project<sup>68</sup>, are proposing data structures which, they claim, can accommodate phenotypical information better than what BIDS can do. However, it is important to note that any individual differences dataset including both phenotypical and endophenotypical information will have the added strength of multimodality, and to date, BIDS is the only data format that can accommodate *both* in a relatively straightforward way, even if it requires extra (and sometimes *post-hoc*) work to prepare the materials intended to be shared. Finally, we must note that when sharing mixed raw and derivative datasets, behavioural data can be easily included as derivatives, these having a more liberal structure, lifting from the end-user the load of reprocessing and calculating basic scores starting from item-level data. This is the route we chose for this dataset, including raw (minimally processed) phenotypic and behavioural data with varying underlying structures, coming from questionnaires and tasks respectively, crucially accompanied by their derivative scores.

In sum, the NEBULA101 dataset aims to promote the study of individual differences in language to better understand a multivariate cognitive system, via the sharing of a truly multimodal dataset in an adequately sized participant sample<sup>70</sup>. We provide sMRI, DWI, task-based and resting-state fMRI in 101 individuals, together with broad phenotypic and behavioural data on linguistic but also on non-linguistic, domain general and domain specific tasks, including cognitive and perceptuomotor tasks. This includes measures of all language aptitude components from phonetics to syntax, measures of reading and reading mediators (e.g. phonological awareness), domain-general cognitive skills, numerical processing, musicality and musical experience and rich multilingual language experience measures. By providing these data to the public domain, we hope to contribute new discoveries to the over-arching construct of language (aptitude), embracing the components of individual behavioural and neural phenotypes as widely as possible.

**Data types.** Behaviour is the phenotype that can be related to, or unify, genetics, neural architecture, neural activity, body structure, physical limitations, and environmental factors<sup>65</sup>. Given the importance of behavioural data in exploring individual differences at the neural level, this dataset includes 28 scores derived from 8 questionnaires, and 74 behavioural measures derived from 21 tasks. Functional neuroimaging data provide information about the neural correlates of cognitive processes, allowing to elucidate how the brain supports specific behaviours and skills. Here, we provide resting-state and task-based functional imaging (fMRI) data, the latter obtained during a language localiser<sup>71</sup>. Finally, NEBULA101 also includes anatomical T1-weighted and diffusion-weighted (DWI) imaging of the brain, which will allow to study brain anatomy and white-matter structural connectivity, to shed light on the brain structural correlates of linguistic behaviours, aptitudes and experiences.

An overview of all measures provided in the NEBULA101 dataset is provided in Table 2. In the Supplementary Information file, Table S1 contains more details specific to the version of the tests used in this dataset, such as any adaptations, modality of administration and derivate scores.

## Methods

**Participants.** According to the Organisation for Economic Cooperation and Development (OECD), approximately 40% of Swiss individuals aged 25 to 34 possess upper secondary or post-secondary non-tertiary education, around 50% have attained tertiary education, and about 10% have below-secondary education<sup>72,73</sup>. Additionally, Switzerland is recognized for its linguistic diversity, as we have extensively discussed in the above paragraphs. The Federal Population Census 2022 Structural Survey<sup>74</sup> indicates that French is spoken by 22.8% of the population. In the survey, 70.1% and 23.4% of the population declared that they speak other national (German, Italian, and/or Romansch) or non-national languages (participants could declare more than one language, therefore there is a degree of overlap in these data). Considering that any study examining the interaction between diverse cognitive domains and language skills will always be somewhat culture-dependent<sup>75</sup>, we chose our target participants with the aim to test at least 100 healthy, relatively multilingual individuals having French as their first or dominant language, which is the primary language of the Canton where the data collection took place.

We recruited 104 adult participants who matched these criteria from the Geneva area, the surrounding French-speaking cantons of Switzerland and neighbouring France through flyers and online advertisements. Participants provided their consent for disclosing their medical history and filled in an online screening survey. To safeguard coherence and avoid confounding factors within the sample, prospective participants were excluded *a priori* if they had musical or simultaneous interpreting professional qualifications (known to interact with language); vision defects that could not be corrected; body implants incompatible with MRI or known claustrophobia; neurological or psychiatric conditions; traumatic head injuries with loss of consciousness; ongoing pregnancies; and past or present illnesses requiring invasive and/or continued medical treatment (such as cancer, chronic and/or autoimmune diseases). Participants with MRI-conditional implants were evaluated on an individual basis upon providing further documentation. Participants with diagnosed developmental dyslexia, as well as those who reported knowing more than 10 languages with a self-reported proficiency equal to or higher than 4 out of 10 in reading, speaking and listening, are not included in this dataset. Fig. 1 reports the number of declared languages for the final sample.

Once eligibility was established, all participants provided a signed informed consent to all subsequent experimental procedures, including anonymised data reuse for open science. One participant subsequently withdrew their consent to data sharing and two were not able to undergo brain imaging due to claustrophobia. Therefore, the final sample includes 101 individuals ( $M_{\text{age}} = 23.35$  years,  $SD = 4.08$ , 68 F,  $M_{\text{education}} = 15.34$  years,  $SD = 2.35$ ). At study completion, all participants received monetary compensation, an image of their brain and a simplified report on their performance on the behavioural tests. Language background and social status show that our participants were quite representative of the well educated and economically stable Swiss society. This information is shown in Fig. 2, with social status represented as a cumulative measure derived from the education and job category of the participant, as well as those of their family of origin and of their partner, if present, calculated with Barratt's Simplified Measure of Social Status (BSMSS)<sup>76</sup>.

**Data collection.** *Data collection.* All interactions and documents provided to participants were in French. Data were collected by six individuals who were either first-language French users, or who had learned French as their second/additional language at an advanced level. All procedures were approved by the Geneva Cantonal Ethical Commission (CCER Protocol N. 2021-01004) prior to the beginning of the study. All participants signed a waiver for anonymised data release in the public domain, and no identifying information was retained, to the best of our knowledge.

Data were collected in two online and two in-person sessions, hereafter named Session 1, 2, 3 and 4, always occurring on different days but in the same order (1–4) between July 2022 and June 2023. Session 1 was unsupervised, and required participants to fill out online questionnaires. Session 2 involved online behavioural data collection, and was supervised by an experimenter. During Session 3 we collected more behavioural data, this time in person. Finally, in Session 4 we collected neuroimaging data.

*Behavioural testing.* Using established and published tests in cognitive psychology is crucial for addressing the reproducibility crisis and curbing the proliferation of one-off tests, ensuring that findings are accurate and that they can be compared and replicated by other researchers<sup>77</sup>. Therefore, where possible, we chose existing questionnaires and behavioural tasks, with the exception of an explicit morphosyntax learning test, which we developed and piloted ourselves due to the lack of such a test in the field (see section “Pilot study”). All instructions and tests were delivered in French. Where the French version of a test was not available, it was *adapted* and checked by at least one first-language French user in the team. Questionnaire instructions (Session 1) were presented in written form. Behavioural task instructions (Sessions 2 and 3) were read by a commercial natural reader (<https://www.naturalreaders.com/commercial/read>), using the voice “Renee France” (if not otherwise specified), at speed -1. fMRI task instructions were given by the experimenter via a microphone. All technical alterations and task adaptations were related to 1) language of delivery, or 2) adapting a paper & pencil test for computer-based delivery. All the above information is thoroughly documented for each test in the Supplementary Information file, Table S1. An overview of the data collection structure is given in Fig. 3 and an extensive theoretical explanation of the tests included in sessions 1–3 can be found in our recent exploratory work on the behavioural correlates of language aptitude, which included this sample<sup>17</sup>.

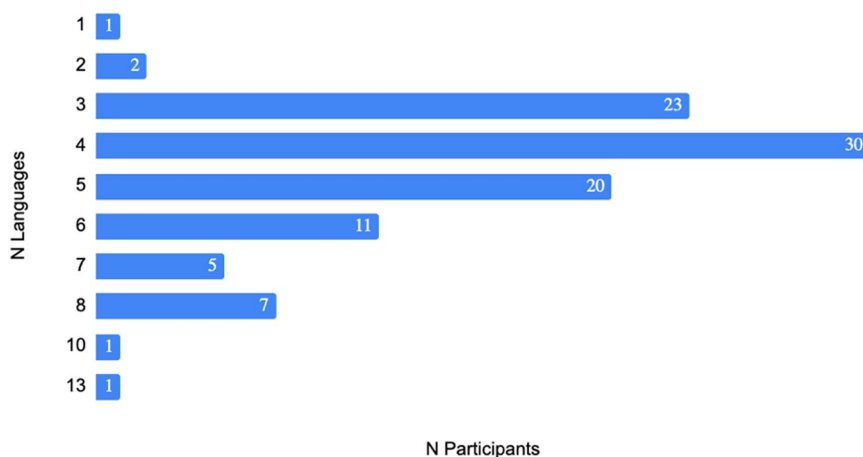
Test	Modality	Construct	Reference
Language Experience and Proficiency Questionnaire (LEAP-Q)	Q	Multilingual language experience	Marian <i>et al.</i> <sup>67</sup>
Code Switching questionnaire	Q	Code switching habits	Rodriguez-Fornells <i>et al.</i> <sup>131</sup>
Motivational Factors Questionnaire (MFQ)	Q	Motivation and attitude towards foreign languages (FL)	Ryan <sup>132</sup> Thompson & Lee <sup>133</sup>
Adult Reading History Questionnaire (AHRQ)	Q	Reading history	Lefly & Pennington <sup>134</sup>
Internal Representations Questionnaire (IRQ)	Q	Modes of internal reasoning	Roebuck & Lupyan <sup>135</sup>
Music Use and Background Questionnaire (MUSEBAQ)	Q	Music training, capacity, preferences, and motivations	Chin <i>et al.</i> <sup>136</sup>
Barratt's Simplified Measure of Socioeconomic Status (BSMSS)	Q	Socioeconomic status	Barratt <sup>76</sup> Rakesh & Whittle <sup>137</sup>
Edinburgh Handedness Inventory (EHI)	Q	Handedness	Oldfield <sup>138</sup> Nedjar <i>et al.</i> <sup>139</sup>
Artgram	T	Language analytic abilities / Morphosyntax	Developed in-house
Modern Language Aptitude Test 5 (MLAT5)	T	Rote learning	Stansfield <sup>140</sup>
Farsi uvular Production Task	T	Foreign sound production	Golestani & Pallier <sup>141</sup>
Hindi Dental Retroflex Contrast	T	Phonological categorisation/discrimination	Golestani <i>et al.</i> <sup>142</sup>
Brocanto	T	Language analytic abilities / Pattern recognition	Kepinska <i>et al.</i> <sup>143</sup> Opitz and Friederici <sup>144</sup>
Raven's Advanced Progressive Matrices (APM)	T	Non-verbal intelligence	Raven <sup>145</sup>
Corsi block	T	Visuospatial memory	Corsi <sup>146</sup> Arce & McMullen <sup>147</sup>
Digit Span	T	Auditory working memory	Wechsler <sup>148</sup> Ryan <i>et al.</i> <sup>149</sup> Conway <i>et al.</i> <sup>150</sup>
Revised Tempo Test	T	Arithmetic abilities	Bellon <i>et al.</i> <sup>151</sup>
Advanced Measures of Music Audiation (AMMA)	T	Music audiation, musicality, musical aptitude	Gordon <sup>152</sup>
Attention Network Test - Interaction (ANT-I)	T	Attention networks: executive control, alerting, orienting	Callejas <i>et al.</i> <sup>153</sup>
California Verbal Learning Task (CVLT)	T	Verbal working memory	Deweer <i>et al.</i> <sup>154</sup>
Finger tapping Test	T	Fine motor skills	Strauss <i>et al.</i> <sup>155</sup> Ashendorf <i>et al.</i> <sup>156</sup>
Purdue Pegboard Test	T	Fine motor skills	Tiffin & Asher <sup>157</sup>
Rapid Automatisated Naming (RAN)	T	Naming automatization	Frederickson <i>et al.</i> <sup>158</sup>
Phoneme suppression	T	Phonological awareness	Rutten <i>et al.</i> <sup>159</sup>
Text Reading	T	Reading skills	"Le Pollueur", Gola-Asmussen <i>et al.</i> <sup>160</sup> "L'Alouette", Lefavrais, 1967 <sup>161</sup>
Word and Pseudoword Reading	T	Reading skills	EVALEC: Sprenger-Charolles <i>et al.</i> <sup>162</sup> ECLA16+: Gola-Asmussen <i>et al.</i> <sup>160</sup>
Spelling task	T	Spelling skills	Gola-Asmussen <i>et al.</i> <sup>160</sup>
Spoonerisms	T	Phonological awareness	Szenkovitz & Ramus, <sup>163</sup>
Non-word repetition	T	Phonological working memory	Majerus, & Van der Linden <sup>164</sup>
Structural MRI (T1-weighted MPRAGE)	N	Brain structural anatomy	Work-in-progress sequence © Siemens Healthineers
Diffusion-Weighted Imaging (DWI)	N	Diffusion gradients	Fedeli <i>et al.</i> <sup>165</sup>
Language Network Functional Localiser (fMRI)	N	Functional activation for language	Malik-Moraleda <i>et al.</i> <sup>71</sup>
Resting-State Functional MRI (fMRI)	N	Resting-state functional activation	Developed in-house
Field maps	N	Intensity of the magnetic field across space	Developed in-house

**Table 2.** Overview of all tasks and modalities. Q = questionnaires; T = cognitive tasks; N = neuroimaging. More details available in Table S1.

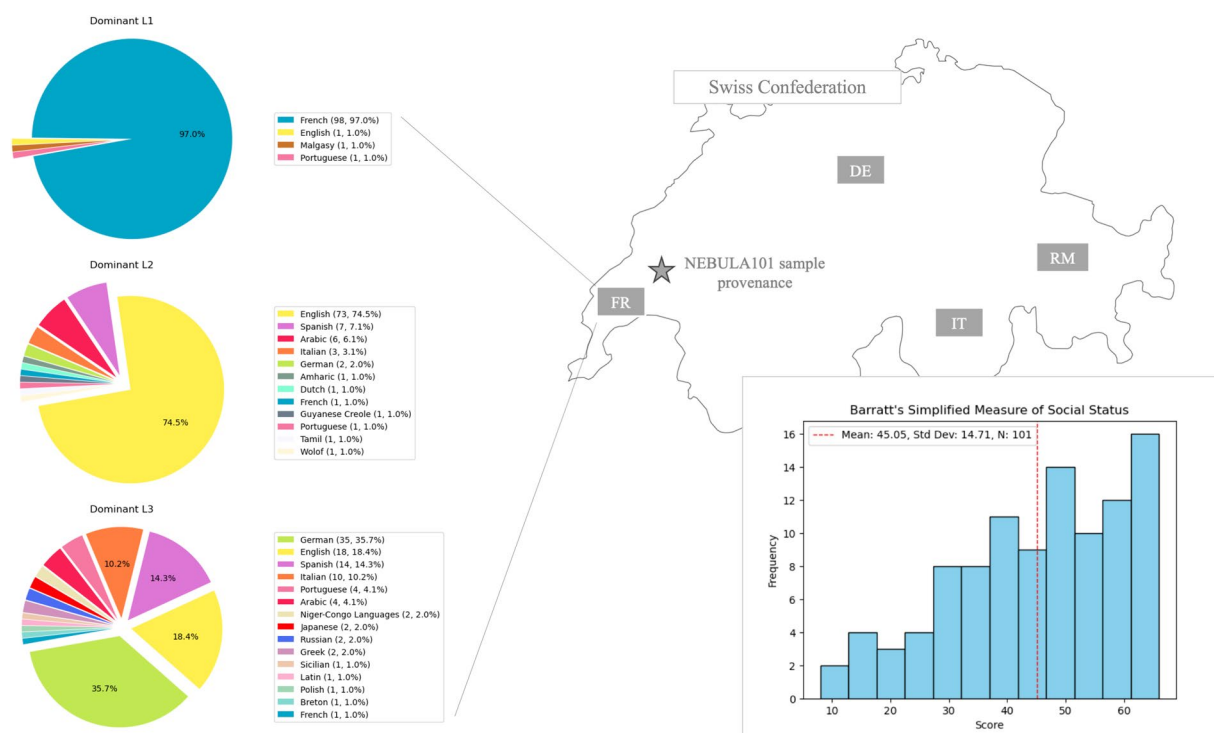
**Pilot study.** Before data collection, an explicit grammar learning task called ArtGram was developed and pilot-tested. ArtGram, designed for adults, extends the PLAB4 task used for language aptitude testing in children and adolescents, for which we could not find an equivalent in adults<sup>78,79</sup>. The test involves learning an artificial, declensional language lexicon with inflected sample sentences, followed by a self-paced multiple-choice, speeded translation task with novel sentences, as described in Table 2 and more extensively in Rampinini *et al.*, (2024)<sup>17</sup>. The pilot study aimed to assess: (1) the task's feasibility, (2) the reliability of the online platform, (3) redundancy with an implicit grammar learning task (Brocanto). Twenty first-language French speakers (11 F,  $M_{\text{age}} = 27.65$ ,  $SD = 8.9$ ) without language or reading disorders participated via video conference. Results showed an insignificant correlation of  $r(18) = 0.44$  between the two grammar tasks.

**Session 1: questionnaires.** Participants initially filled out a series of questionnaires online using Qualtrics XM©. The questionnaire sequence was fixed: first came the Edinburgh Handedness Test and the BSMSS questionnaire, completed upon recruitment. Then came the demographics, MFQ, IRQ, AHRQ, MUSEBAQ, Code Switching

## N Participants vs N Languages



**Fig. 1** Frequency distribution of the number of reported languages in the final sample (N = 101).



**Fig. 2** Participants' language background up to their third language (left, languages coded by colour) and Socioeconomic status (bottom right) measured via Barratt's simplified measure of social status<sup>76</sup>.

and LEAPQ (see Table 2). When available from the questionnaire manual, automatic scoring was implemented in Qualtrics via the "Formula Field" and resulted in one column per score in the derivative questionnaire dataset, for each participant.

**Session 2: online behavioural testing.** During session 2, supervised behavioural data collection was conducted via Zoom®. Participants were guided through a demonstration video on how to share both their sound and entire screen, use wired headphones for reliable online measurements<sup>80,81</sup>, and ensure their microphone was functioning properly. Tasks were delivered through the Gorilla web interface<sup>82</sup>. Headphone and microphone tests from the Gorilla open materials section were mandatory before starting the task sequence. Researchers supervised the session, intervening only if technical issues arose. The system prevented participants from logging into the session from mobile phones or tablets, and only allowed the use of Mozilla Firefox® or Google Chrome®.

Participants navigated through tasks autonomously in a predetermined order among 15 possible pseudo-randomizations, one of which was automatically assigned by the system at the start of the testing sequence. Before each task, participants received on-screen written instructions in French, and could not proceed until the natural reader finished delivering the same instructions orally. Fixed-length breaks (3 or 5 minutes) were included after the most intensive tasks to optimize concentration and compliance. Participants could end the break early if ready to continue, with a 60-second timer appearing in the last minute of the break if they had not yet proceeded to the next task.

**Session 3: in-person behavioural testing.** In this session, participants were tested in person for tasks that required closer supervision due to their length (e.g. the ANT-I), or that required manual measurements of reaction times (RT), (e.g. the literacy and literacy mediator tests). Testing was conducted at the Human Neuroscience Platform of “Campus Biotech” in Geneva, in a dedicated, sound-protected room, using the same laptop, mouse, and headset (headphone with microphone) for all participants (see Fig. 3). Tasks were organized and delivered via the Gorilla interface using the previously described pseudo-randomization strategy. The same microphone check was included to ensure the safe recording of tasks requiring vocal responses for later assessment.

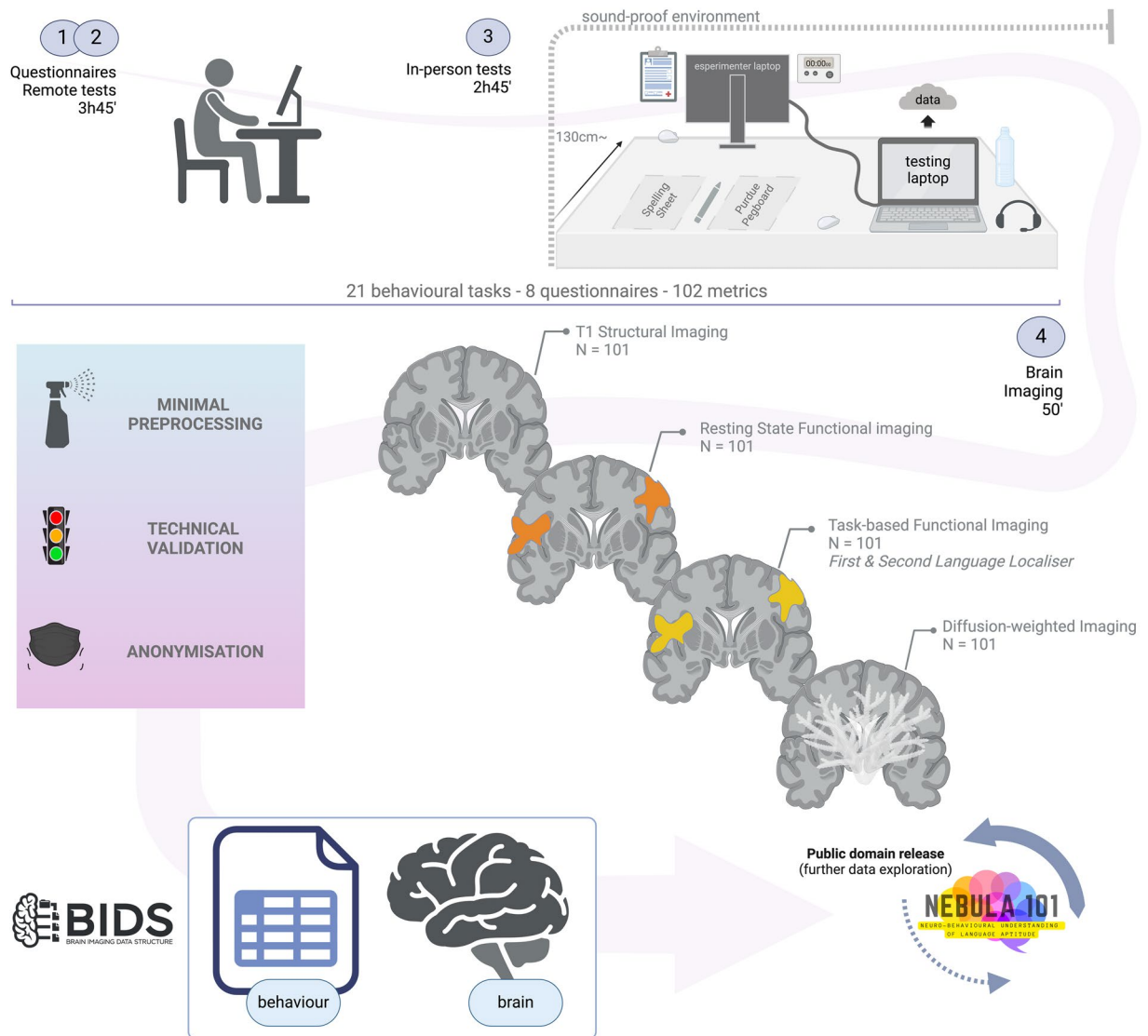
An experimenter closely supervised the session through a connected screen, mouse, and keyboard while facing the participant, intervening only when necessary due to task requirements or technical issues. For tasks requiring a vocal response, the experimenter manually recorded accuracy and/or RT (assessed via a chronometer) in the session booklet. To prevent data loss, these tasks were also audio-recorded, and the responses were later verified for accuracy by a team member with French as their first language. Task events not requiring manual measurement or verification were recorded directly in Gorilla. During this session, participants also performed three extra tasks that were part of another project<sup>17</sup>.

**Session 4: Magnetic resonance imaging.** Following behavioural testing, participants were invited to a brain imaging session on a Siemens 3 T Magnetom-Prisma scanner equipped with a 64-channel head coil, again at Campus Biotech. Prior to scanning, they filled in and signed an MRI safety questionnaire to again verify their eligibility for the procedure. When required, participants were fitted with MRI-compatible goggles to correct their vision. The language localiser task was administered in Matlab r2021b with Psychtoolbox 3, through a computer connected to a screen in the back of the scanner room, that participants could see through a mirror placed on top of the RF head coil. During the whole session, we could observe their eyes through an eye-tracking camera, to check that they were awake. During the imaging session, the field map, resting-state sequence and the language localiser were administered in random order, but the resting-state sequence was always preceded by the T1-weighted anatomical scan to avoid spurious activations due to carrying out an active task just before. After a short break outside the scanner, participants were repositioned and underwent the DWI sequence and its corresponding field map acquisition. During this session, participants also performed three fMRI tasks and one anatomical scan that were part of different projects. An overview of the imaging session parameters is provided in Table S1 of the Supplementary Information file.

**fMRI language localiser.** Participants were instructed to keep their eyes open and look at a black fixation cross on a white background while listening to intact and degraded snippets of the story ‘Alice in Wonderland’, in their first and second most dominant language of choice (L1, L2)<sup>71,83</sup>. This localiser can be used to inspect individual differences in language activation during quasi-naturalistic listening, and their relationship with behavioural measures of language and cognition<sup>84</sup>. The original localiser paradigm is publicly available from the authors’ webpage. As described in Table S1 in the Supplementary Information file, and the Data Records section, we modified this paradigm to include a degraded L2 condition.

**Resting-state fMRI sequence.** Brain connectivity at rest can be linked to individual differences in language<sup>85,86</sup>, reading<sup>87,88</sup>, and other domains of cognition such as executive function<sup>89,90</sup>. To collect information on resting-state brain activity, we asked participants to lie down with their eyes open, instructing them to relax their body and mind as best as they could, while projecting a white fixation cross on a black background.

**BIDS conversion.** This dataset conforms to BIDS v1.9.0 and was validated using the command line version of `bids-validator v1.14.14` (<http://bids-standard.github.io/bids-validator/>). A LINUX Terminal print of the output is provided in the Supplementary Information file (Fig. S1). The dataset has been annotated using the Neurobagel annotation tool (<https://neurobagel.org/>) for enhanced findability<sup>91</sup>, the annotations have been saved in `/neurobagel` at root level, and all JSON sidecars have been validated with the online version of JSONLint. The BIDS data format conversion occurred in several steps, part of which could be planned before data collection (such as naming of the MRI sequences, folder structure, and participant codes), while others were performed *post hoc* to adapt data generated by environments not optimised for BIDS. In general, the procedure aimed at having a BIDS-coherent data structure and file names, (re)organising the content of tabular files, adding sidecar files to accompany data and customising code to work in a BIDS folder structure. These steps were performed on all data collected during the same testing session, and the NEBULA101 data were subsequently imported into the `/nebula101` data space. Nonetheless, for clarity, we provide the specific heuristic for the construction of this dataset after DICOM to NiFTI conversion in `/code/heudiconv/heuristic.py`. Considering this procedure, all code described is specific to the NEBULA101 dataset but is not meant to be rerun, and is given with paths relative to the dataset BIDS root folder, but might reference to folders outside this structure, for example to source data. Outside of the code performing data import or behavioural data cleaning, no further reference is generally made to external/unavailable files. We describe the steps in Table 3.



**Fig. 3** Data acquisition and processing structure (session duration is approximate for behavioural data acquisition due to individual variability in task completion times). [Illustration created with BioRender.com].

### Data Records

The dataset is published on OpenNeuro under a Creative Commons CC0 1.0 (Universal Public Domain Dedication) license with accession number ds005613 at <https://doi.org/10.18112/openneuro.ds005613.v1.0.192>.

In this section we further describe the data structure and its contents. The folder `/nebula101` (Fig. 4) constitutes the root level of the dataset. It contains 101 participant folders with raw data denoted by the code `sub-pp` followed by three digits, the folders `/phenotype`, `/code`, `/derivatives`, `/stimuli`, `/neurobagel`, and the mandatory files required by BIDS (README, participant and dataset description files). The `/phenotype` folder contains tabular data from questionnaires, while `/neurobagel` contains subject-level annotations (harmonized phenotypic properties and imaging metadata) that can be encoded in a knowledge graph (see Technical Validation). As concerns the subject folders, inconsistent numbering is due to non-included participants (see Participants section), the participant who denied consent to share their data, and the two participants who could not undergo brain imaging). We describe the other contents of `/nebula101` in detail here below.

The `/code` folder in Fig. 5 contains subfolders with (1) fMRI paradigm data for the language localiser; (2) code to import, modify or generate the BIDS-compliant files as described in Table 3; (3) the BIDS conversion heuristic; (4) the preprocessing information and code folder; (5) the validation materials folder.

The `/derivatives` folder in Fig. 6 contains the preprocessed behavioural data and their sidecar files, as well as the results of all validation pipelines.

- `cumulative_farsi_rater*` = derivative scores from the Farsi task ratings and their sidecar files.
- `nebula_101_all_questionnaire_scores*` = scores from questionnaires and their sidecar file.

Step	Description	Language and environment	Custom code	Code availability	Code location
1	DICOM to Nifti conversion: conversion was handled in heudiconv via dcm2nii	Python heudiconv	No	No	n/a
2	File structure heuristic: we set up logging, defined constants, and included functions to create keys for file paths and categorise imaging sequences into different modalities based on their metadata	Python heudiconv	Yes	Yes	code/heudiconv/heuristic.py
3	Modifying language localiser log files for BIDS compatibility	Python	Yes	Yes	/code/create_nebula/nebula101_bidsify_mri_logs_aliceloc.py
4	Importing NEBULA101 data from main BIDS data directory, including behaviour and phenotype data	Python	Yes	Yes	/code/create_nebula/create_nebula.py /code/create_nebula/make_beh.py /code/create_nebula/make_phenotype.py
5	Creation of JSON sidecar files for derivative scores	Python	Yes	Yes	/code/create_nebula/create_nebula.py /code/create_nebula/create_score_id_json.py
6	Modifying field map sidecar “intendedFor” field for BIDS compatibility (to remove references to extra scans) and correcting label naming errors in field maps (fid- to acq-)	Python	Yes	Yes	/code/create_nebula/create_nebula.py
7	Importing stimuli from fMRI language localiser to store in BIDS dataset and adjusting their folder structure	Python	Yes	Yes	/code/create_nebula/create_nebula.py
9	Defacing T1w MPRAGE scans	Python PyDeface	Yes*	Yes	code/1_anat/0_pydeface/ nebula101_loop_pydeface.py * The (minimal) customisation involved creating a script to run the code on the NEBULA101 sample
10	Creation of dataset-level sidecar files	Python	Yes	Yes	/code/create_nebula/create_nebula.py

**Table 3.** List of the steps taken to make the dataset BIDS-compliant.

/nebula101/

code	sub-pp011	sub-pp031	sub-pp052	sub-pp088	sub-pp114	sub-pp135	sub-pp162
derivatives	sub-pp012	sub-pp032	sub-pp053	sub-pp091	sub-pp115	sub-pp137	sub-pp163
neurobagel	sub-pp013	sub-pp033	sub-pp054	sub-pp092	sub-pp116	sub-pp145	sub-pp164
phenotype	sub-pp018	sub-pp035	sub-pp067	sub-pp093	sub-pp117	sub-pp147	sub-pp166
stimuli	sub-pp019	sub-pp036	sub-pp068	sub-pp094	sub-pp118	sub-pp149	sub-pp168
sub-pp001	sub-pp020	sub-pp038	sub-pp069	sub-pp095	sub-pp119	sub-pp150	sub-pp169
sub-pp003	sub-pp021	sub-pp039	sub-pp072	sub-pp099	sub-pp120	sub-pp151	sub-pp170
sub-pp004	sub-pp022	sub-pp042	sub-pp073	sub-pp105	sub-pp124	sub-pp152	sub-pp171
sub-pp005	sub-pp023	sub-pp043	sub-pp074	sub-pp106	sub-pp125	sub-pp155	dataset_description.json
sub-pp006	sub-pp024	sub-pp044	sub-pp076	sub-pp107	sub-pp126	sub-pp156	participants.json
sub-pp007	sub-pp025	sub-pp045	sub-pp077	sub-pp108	sub-pp127	sub-pp158	participants.tsv
sub-pp008	sub-pp026	sub-pp046	sub-pp078	sub-pp110	sub-pp128	sub-pp159	README
sub-pp009	sub-pp027	sub-pp048	sub-pp083	sub-pp111	sub-pp129	sub-pp160	
sub-pp010	sub-pp030	sub-pp050	sub-pp085	sub-pp112	sub-pp133	sub-pp161	

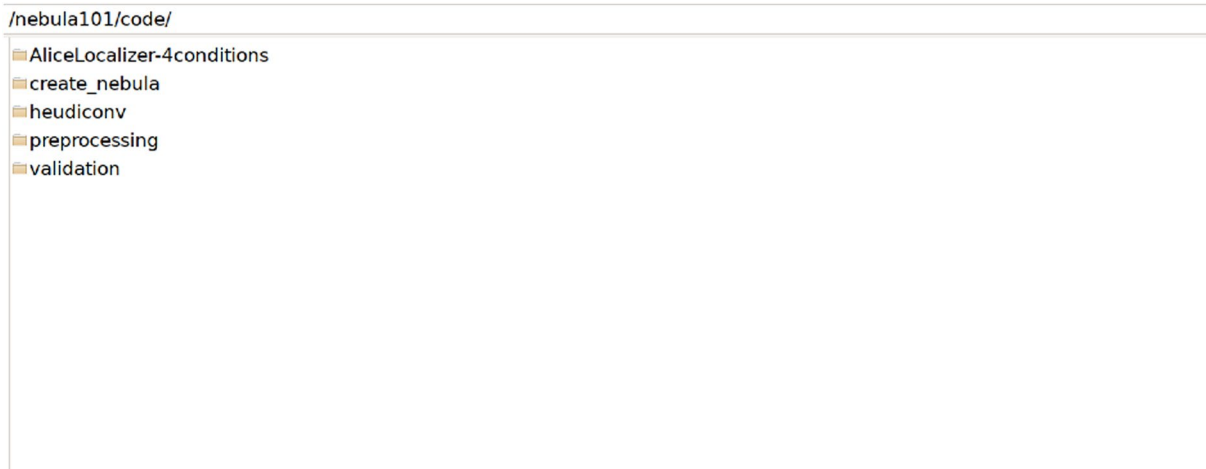
**Fig. 4** Root dataset folder.

- nebula\_101\_all\_task\_scores\* = scores from tasks and their sidecar file.
- nebula\_101\_leapq\_annotation\_iso\_glottolog\* = mapping of language names to ISO and Glottolog codes, and its sidecar file.
- nebula\_101\_leapq\_data\* = LEAPQ data and their sidecar file.
- nebula\_101\_leapq\_langname\_order\* = LEAPQ languages in order of acquisition and their sidecar file.

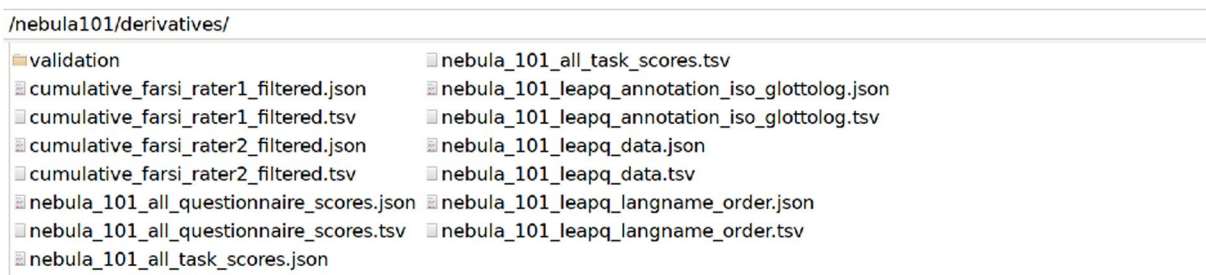
The /stimuli folder contains a subfolder where the language localiser materials are stored and can be referenced by the BIDS validator.

The participant folders, named sub-pp\*/ses-01/, contain /anat, /dwi, /fmap, /func and /beh folders. These host the raw imaging and behavioural data of each participant, the latter having been minimally preprocessed to remove metadata and information unrelated to the scoring, as described in Technical Validation. An example with sub-pp001/ses-01/ is shown in Figs. 7–11.

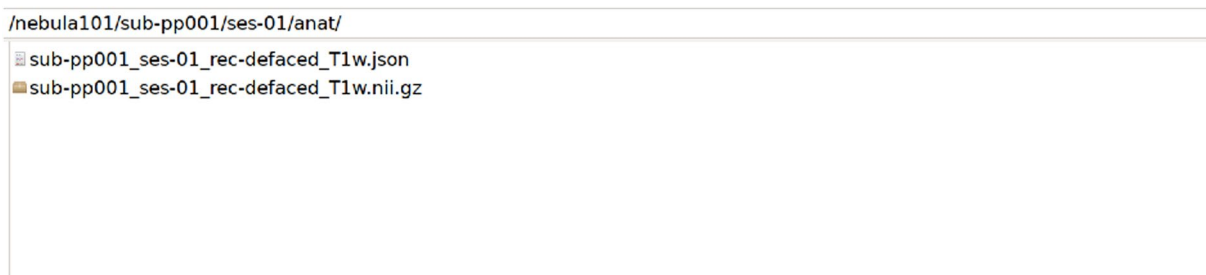
The /code/AliceLocalizer-4conditions/ folder contains our own version of the materials provided in the public domain by the creators of the task<sup>71,83</sup>, which we modified to suit our needs, as described in the fMRI language localiser description (see: Data collection). This version of the code can read the stimuli from the structure of the BIDS folder. It will still require manual intervention to create the events.tsv files from the Matlab log files, as Matlab currently provides a different tabular structure. We remind the reader that we



**Fig. 5** Code folder.



**Fig. 6** Derivatives folder.



**Fig. 7** Defaced T1 MPRAGE data and their sidecar file.

altered the structure of the log files to create simplified BIDS events with no redundant or unneeded information in `/code/create_nebula/bidsify_mri_logs_aliceloc.py`.

The `/code/preprocessing` folder includes numerically ordered subfolders containing code to reprocess the anatomical, field map, functional and diffusion data (see Technical Validation), and the reference code for the raw-to-derivative conversion of behavioural data. Any steps contained in these subfolders are meant to be run sequentially to recreate the materials needed for technical validation, or if the user decides to process raw brain data using our same pipelines.

- `0_beh` contains behavioural data preprocessing code.
- `1_anat` contains the brain extraction code.
- `2_fmap` contains code for field map preparation.
- `3_func` contains code for minimal fMRI preprocessing, to get the displacement parameters used in the technical validation.
- `4_dwi` contains the preprocessing code for the diffusion data necessary to obtain technical validation reports.

```
/nebula101/sub-pp001/ses-01/dwi/
```

```
sub-pp001_ses-01_dwi.bval
sub-pp001_ses-01_dwi.bvec
sub-pp001_ses-01_dwi.json
sub-pp001_ses-01_dwi.nii.gz
```

**Fig. 8** DWI data: beta values, vectors, and sequence data with sidecar file.

```
/nebula101/sub-pp001/ses-01/fmap/
```

```
sub-pp001_ses-01_acq-alpha_magnitude1.json sub-pp001_ses-01_acq-beta_phasediff.json
sub-pp001_ses-01_acq-alpha_magnitude1.nii.gz sub-pp001_ses-01_acq-beta_phasediff.nii.gz
sub-pp001_ses-01_acq-alpha_magnitude2.json sub-pp001_ses-01_acq-gamma_magnitude1.json
sub-pp001_ses-01_acq-alpha_magnitude2.nii.gz sub-pp001_ses-01_acq-gamma_magnitude1.nii.gz
sub-pp001_ses-01_acq-alpha_phasediff.json sub-pp001_ses-01_acq-gamma_magnitude2.json
sub-pp001_ses-01_acq-alpha_phasediff.nii.gz sub-pp001_ses-01_acq-gamma_magnitude2.nii.gz
sub-pp001_ses-01_acq-beta_magnitude1.json sub-pp001_ses-01_acq-gamma_phasediff.json
sub-pp001_ses-01_acq-beta_magnitude1.nii.gz sub-pp001_ses-01_acq-gamma_phasediff.nii.gz
sub-pp001_ses-01_acq-beta_magnitude2.json
sub-pp001_ses-01_acq-beta_magnitude2.nii.gz
```

**Fig. 9** Field maps: phase and magnitude images with their sidecars. This participant has an extra field map (gamma), as explained in the Usage Notes.

```
/nebula101/sub-pp001/ses-01/func/
```

```
sub-pp001_ses-01_task-aliceloc_run-001_bold.json sub-pp001_ses-01_task-aliceloc_run-002_events.tsv
sub-pp001_ses-01_task-aliceloc_run-001_bold.nii.gz sub-pp001_ses-01_task-aliceloc_run-003_bold.json
sub-pp001_ses-01_task-aliceloc_run-001_events.json sub-pp001_ses-01_task-aliceloc_run-003_bold.nii.gz
sub-pp001_ses-01_task-aliceloc_run-001_events.tsv sub-pp001_ses-01_task-aliceloc_run-003_events.json
sub-pp001_ses-01_task-aliceloc_run-002_bold.json sub-pp001_ses-01_task-aliceloc_run-003_events.tsv
sub-pp001_ses-01_task-aliceloc_run-002_bold.nii.gz sub-pp001_ses-01_task-rest_run-001_bold.json
sub-pp001_ses-01_task-aliceloc_run-002_events.json sub-pp001_ses-01_task-rest_run-001_bold.nii.gz
```

**Fig. 10** Functional MRI data with their sidecar files and event files. *Rest* refers to resting-state fMRI, *aliceloc* refers to the language localiser.

```
/nebula101/sub-pp001/ses-01/beh/
```

```
sub-pp001_ses-01_task-amma_beh.json sub-pp001_ses-01_task-cvlt_beh.tsv sub-pp001_ses-01_task-purdue_beh.json
sub-pp001_ses-01_task-amma_beh.tsv sub-pp001_ses-01_task-cvltlong_beh.json sub-pp001_ses-01_task-purdue_beh.tsv
sub-pp001_ses-01_task-ant_beh.json sub-pp001_ses-01_task-cvltlong_beh.tsv sub-pp001_ses-01_task-ran_beh.json
sub-pp001_ses-01_task-ant_beh.tsv sub-pp001_ses-01_task-digitback_beh.json sub-pp001_ses-01_task-ran_beh.tsv
sub-pp001_ses-01_task-apm_beh.json sub-pp001_ses-01_task-digitback_beh.tsv sub-pp001_ses-01_task-rttsub_beh.json
sub-pp001_ses-01_task-apm_beh.tsv sub-pp001_ses-01_task-digitfor_beh.json sub-pp001_ses-01_task-rttsub_beh.tsv
sub-pp001_ses-01_task-artgram_beh.json sub-pp001_ses-01_task-digitfor_beh.tsv sub-pp001_ses-01_task-rttsum_beh.json
sub-pp001_ses-01_task-artgram_beh.tsv sub-pp001_ses-01_task-farsi_beh.json sub-pp001_ses-01_task-rttsum_beh.tsv
sub-pp001_ses-01_task-brocanto_beh.json sub-pp001_ses-01_task-farsi_beh.tsv sub-pp001_ses-01_task-spelling_beh.json
sub-pp001_ses-01_task-brocanto_beh.tsv sub-pp001_ses-01_task-fingertapping_beh.json sub-pp001_ses-01_task-spelling_beh.tsv
sub-pp001_ses-01_task-corsiback_beh.json sub-pp001_ses-01_task-fingertapping_beh.tsv sub-pp001_ses-01_task-textreading_beh.json
sub-pp001_ses-01_task-corsiback_beh.tsv sub-pp001_ses-01_task-hindi_beh.json sub-pp001_ses-01_task-textreading_beh.tsv
sub-pp001_ses-01_task-corsifor_beh.json sub-pp001_ses-01_task-hindi_beh.tsv sub-pp001_ses-01_task-wordreading_beh.json
sub-pp001_ses-01_task-corsifor_beh.tsv sub-pp001_ses-01_task-mlat5_beh.json sub-pp001_ses-01_task-wordreading_beh.tsv
sub-pp001_ses-01_task-cvlt_beh.json sub-pp001_ses-01_task-mlat5_beh.tsv
```

**Fig. 11** Raw behavioural data files with their sidecar files.

The `/code/validation` folder contains the following subfolders (see Technical Validation for details on the operations performed):

- `/anat`: segmentation and sample homogeneity plotting code for the anatomical scans (auto-generated by Matlab): `/cat12/cat_stat_homogeneity.m` and `/cat12/cat_stat_homogeneity_autoplot.m`
- `/beh` contains the following items:
  - The code to run reliability analysis on behavioural data: `calculate_cronbach_alpha.py`
  - Three additional Python notebooks for generating the correlation data presented here and additional tables and matrices contained in the Supplementary Information file (Tables S2, S3 and Fig. S2):
    - `correlate_matrix.py`
    - `calculate_descriptives.py`
    - `farsi_inter_rater.py`
- `/data_checks`: various Python scripts to create data lists of behavioural and imaging data, and missing data heatmaps, as explained in Technical Validation.
- `/dwi`: individual participant pdf QC reports, the code for generating them, and the group report folder.
- `/func`: two Python scripts, one for importing displacement .rms data from FSL, and the other for violin plots of average absolute and relative displacement during fMRI (resting-state and task-based).
  - `1_copy_motion_params.py`
  - `2_mean_abs_rel_disp_violin.py`

For each of the described folders in `/code/validation/`, there is a mirror `/derivatives/validation/` folder where the results of the validation pipelines are stored, and specifically:

- `/derivatives/validation/anat/cat12/` contains the single-subject CAT12 reports in PDF form.
- `/derivatives/validation/beh/` contains the item-level reliability data used to calculate Cronbach alpha and all correlation measures.
- `/derivatives/validation/data_checks/` contains data presence checks generated by the respective code for brain imaging, phenotype and behaviour.
- `/derivatives/validation/dwi/fsl/` contains the reports from subject- and group-level quality assessment of DWI data.
- `/derivatives/validation/func/fsl/aliceloc/` contains the reports from subject- and group-level quality assessment of task fMRI data.
- `/derivatives/validation/func/fsl/rest/` contains the reports from subject- and group-level quality assessment of resting-state fMRI data.

All code is heavily commented for user-friendliness.

## Technical Validation

In compliance with the BIDS indications for mixed raw and derivative datasets, to improve user experience and reduce redundancy, we chose to include preprocessed behavioural data as derivative tables. We also provide quality control (QC) reports for all imaging modalities, together with the code for reproducing them. Below we specify the details of technical validation for each data type.

**Questionnaires.** Questionnaire source data were retrieved from Qualtrics XM® and preprocessed in Python. Each questionnaire was cleaned with its own script and then merged with the others via another script. The overall process aimed at data cleaning, as scoring was generally performed *ante-hoc* via the Qualtrics graphical user interface.

- 1) Log file data collection: for each questionnaire, from the source tabular file containing all participants, we eliminated irrelevant metadata. In rare cases, when duplicates were identified we eliminated them at source level by keeping the first non-empty instance of the questionnaire, blind to the contents, to avoid selection bias.
- 2) Data cleaning and handling of missing data: in all questionnaires, string responses were recoded as numbers. Missing data (NaN) were handled as follows: as a general rule, within-questionnaire, all responses were forced (i.e. progression was halted if a response was missing or on-screen warnings appeared to highlight unfilled responses). When information was indeed unavailable for a participant, we set up options to input specific strings, signalling to us that the information was missing. When missing information was nonetheless identified in a questionnaire with *optional* responses, but without the string signalling unavailability, we evaluated on a case-by-case basis whether it was possible to obtain it from the participant. In such rare cases, the questionnaire was retaken by the participant in session 4, assisted by an experimenter (to avoid data alteration), and only missing responses were filled in. In even rarer cases, when an *entirely* missing questionnaire was identified, participants were invited to complete it through an individual link to that specific questionnaire. Any remaining entries with unavailable information were converted to NaN

during data preprocessing. With the strategies we put in place, all questionnaires have been completed by all participants in this dataset, except for one monolingual participant not having taken the code-switching questionnaire due to only speaking one language.

- 3) LEAPQ cleaning: the LEAPQ, due to its length and structure, required additional, *ad-hoc* handling. In Qualtrics, participants reported the number of languages they knew, including extinct ones and dialects, in order of dominance, that is, how much they used each of them over the others, and in order of acquisition, that is, the temporal order in which they learned their languages. They then provided general information for each language in order of dominance (e.g., names, exposure) and specific information (e.g., learning modality, use context). The resulting data file was restructured by language number and question type, filling gaps with NaN to create a 'staircase-like' structure where participants (rows) and questions (columns) are listed by language number, in ascending order. The resulting file is available in `/derivatives/nebula101_leapq_data.tsv` with a corresponding sidecar file explaining the column contents. The list of languages in order of acquisition was extracted into a separate file to avoid confusion, as it was not always the case that the dominance and acquisition lists coincided. Language order information can be accessed at `/derivatives/nebula_101_leapq_langname_order.tsv`.
- 4) Raw language names as entered by the participants were conformed to ISO 639-3 and Glottolog codes<sup>93</sup> for standardised reference. The conversion map can be found in `/derivatives/nebula_101_leapq_annotation_iso_glottolog.tsv` and its sidecar file.
- 5) Multilingual language experience entropy. The LEAP-Q is a widely used instrument for assessing multilingual language experience, encompassing contexts of use, learning, use choices, history with the language, and native-likeness. Given that composite measures of multilingualism derived from the LEAP-Q would be beneficial for quantitative analyses, we calculated four different continuous 'multilingualism scores' for each participant, to reflect their multilingual experience cumulatively<sup>67</sup>. Three scores were based on their self-reported proficiency in speaking, reading, and comprehension across all reported languages, while the fourth score was based on current exposure to all reported languages. Following previous research<sup>17,94,95</sup>, each participant's multilingualism score combined proficiency or exposure across their different languages using Shannon's entropy equation<sup>96</sup> within the R entropy package<sup>97</sup>.

Item-level questionnaire data are stored in `/nebula101/phenotype/` as individual tabular files, each containing data for one questionnaire, where each subject is represented in a single row as a `participant_id` and their responses in the subsequent columns. Each phenotypic tabular file is accompanied by a JSON sidecar file describing each column:

- `ahrq.tsv` and `ahrq.json` contain data for the Adult Reading History Questionnaire.
- `bsmss.tsv` and `bsmss.json` contain data for the Barratt Simplified Measure of Social Status.
- `code_swt.tsv` and `code_swt.json` contain data for the Code Switching Questionnaire.
- `handedness.tsv` and `handedness.json` contain data for the 10-item French version of the Edinburgh Handedness Inventory.
- `irq.tsv` and `irq.json` contain data for the Internal Reasoning Questionnaire.
- `mfq.tsv` and `mfq.json` contain data for the Motivation Factors Questionnaire.
- `musebaq.tsv` and `musebaq.json` contain data for the Music Experience Use and Engagement Questionnaire.

Given the higher level of processing that the LEAP-Q data went through, its tabular files (responses, language annotations and language order information) reside in `/derivatives`, together with the comprehensive derivate score file of all the questionnaires, `nebula_101_all_questionnaire_scores.tsv`.

**Behavioural tasks.** Behavioural tasks were pre-processed in Python for data cleaning, derivate score calculation and BIDS conversion. The process happened in steps:

- 1) Log file collection: data were downloaded from Gorilla. Given that participants entered 1 out of 15 possible randomised task sequences, automatically assigned to people by Gorilla, the empty log files from the 14 unused randomisations for each participant had to be deleted. Microphone and headphone check files were identified and renamed.
- 2) Log file cleaning: irrelevant metadata were eliminated. Task log files were renamed to reflect the actual task (as Gorilla provides them in encrypted form) and to contain BIDS-compliant strings by adding the relevant labels, then converted to tab-separated values. Within each task's log file, column names were renamed to more easily interpretable strings where necessary. Audio files from voice-recorded tasks were renamed as well, for easier reference.
- 3) Derivate score calculation and tabular file merging: for each task, one or more derivate scores were calculated, based on the task's protocol. See Table 2. and Table S1 for the references and scores that were selected.

The behavioural data preprocessing code is provided for reference in `/nebula101/code/preprocessing/0_beh`, and specifically:

- a. `2DataCleaning_gorilla.py`: this script cleans the source logfiles (not provided in this dataset) from metadata, as explained.
- b. `3Scoring_gorilla.py`: this script scores the files that were generated by Gorilla and cleaned, based

- on an automated pipeline by using the ‘correct’ and ‘incorrect’ columns and/or reaction time information (depending on task-specific scoring protocols).
- c. `4Scoring_manual.py`: this script scores logfiles that were created from tasks requiring human intervention (for example, manual RT measurements of voice responses timed via a chronometer, and/or live pencil scoring by the experimenter, such as the Text or Word and Pseudoword reading tasks, and/or written answers like in the Spelling task). Here, Gorilla was used to run the task within the pseudo-randomisation pipeline but scoring necessitated manual intervention.
  - d. `5Scores-merging.py`: this script merged the scores into a tabular file, from which they were later imported to the `/nebula101` data space via the script `create_nebula.py`.
  - e. `7farsi_raters_assessment.py`: The script calculates cumulative and average scores for the Farsi uvular sound production task, based on assessments made by two independent raters. It produces the derivative file provided in `/nebula101/derivatives/`.

These actions were performed in subsequent steps by Python scripts acting on source data. Users will not be required to rerun these scripts since we provide cleaned tabular files containing the minimally processed raw accuracy and RT scores rather than the source data, as well as, crucially, the derivative scores calculated and ready for use.

For tasks whose data were recorded automatically in Gorilla (when a simple correct/incorrect input was sufficient, and RT could be measured by key press or mouse click), the raw data of each subject will contain their `participant_id`, as minimally mandated by BIDS, accuracy, and RT columns for each trial, if these are sufficient to derivate a score provided in `/derivatives` via the code provided in `/code/preprocessing/0_beh/`. For the manually scored tasks, as described in steps c) and e) above, single-subject raw data resulting from the digitisation of paper-and-pencil materials are provided: such might include block-level data (e.g. each iteration and type of recall trials in the case of the CVLT) or condition-level data (e.g. word types in the case of the Word and Pseudoword reading task, story types in the case of the Text reading task), whose structure will be inherent and specific to the task itself. In the case of the Spelling, Nonword repetition and Spoonerisms tests, scoring was digitised at task level (as per protocol). We therefore included these only in the derivative table. Despite any differences in structure, all behavioural files conform to the minimal BIDS requirements for their data type and are accompanied by extensive JSON descriptions. All materials used to process source data can be made available upon request, as well as source data that can be anonymised.

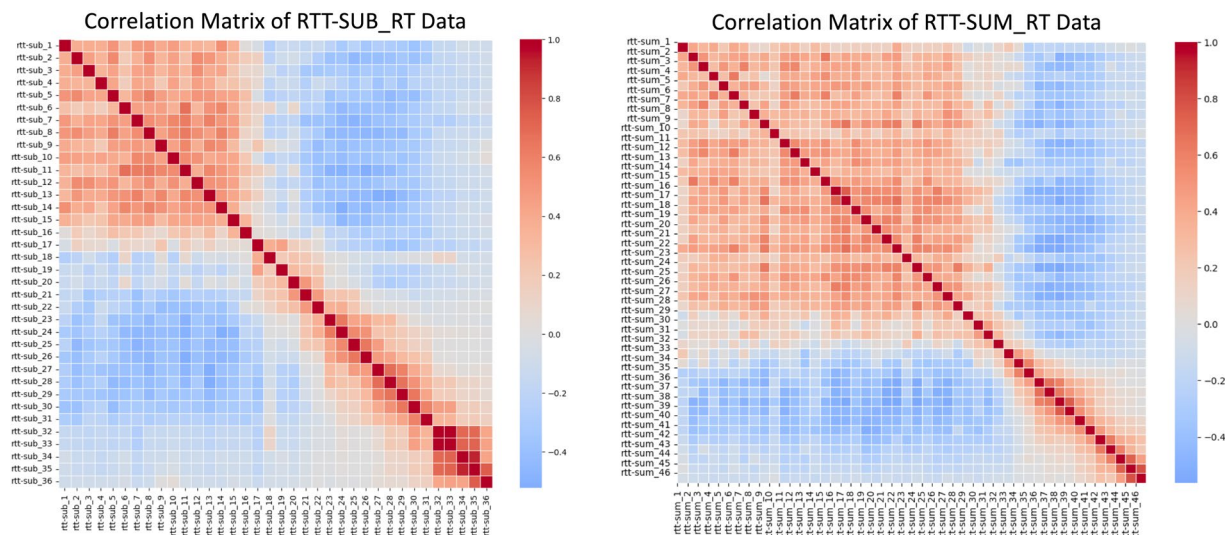
A comprehensive file containing all derivate scores for all behavioural tasks, obtained with steps b) and c) lives in `/derivatives` and is called `nebula_101_all_task_scores.tsv`.

In this location the Farsi task derivate scores obtained with step e) can also be found: these are called `cumulative_farsi_rater*_filtered.tsv` and `cumulative_farsi_rater*_filtered.tsv`. All derivate files are accompanied by JSON sidecars.

*Behavioural data: correlations and reliability.* Behavioural scores, whether resulting from tasks or questionnaires (here we generally refer to them together as “behavioural data”, and we use “scores” or “measures” when referring specifically to their outcomes), derive from item-level data. To validate these measures, we first explored them via correlations, and then ran reliability analyses in the form of internal consistency coefficients, in a conceptually similar approach to our recent exploratory analysis of behavioural data with a larger cohort<sup>17</sup>. Intercorrelation of the scores obtained from behavioural data (tasks and questionnaires) was minimal, and reliability was acceptable to good in most cases. The code for calculating Pearson’s correlations on the z-scored data as well as Cronbach Alpha, along with item-level reliability data and summary output (figures, tables) are available in `/code/validation/beh/` and in the Supplementary Information (Figs. S2 and S3, Tables S2–S5).

To better understand the correlation patterns and reliability of these data, some considerations must be made about correlations and about reliability measured by internal consistency. Regarding correlations, on the one hand, high intercorrelations can be informative about shared variance between tests, but on the other hand, having an excessively intercorrelated dataset can lead to problems linked to multicollinearity<sup>98–101</sup>. On a conceptual level, the separability of latent variables facilitates the interpretation of the overall construct. Our choices, when planning this data collection and selecting the tests, aimed at isolating the best measures to represent each variable that we wanted to test in the context of this exploratory project on language aptitude. All things considered, a relative degree of independence of our measures was not only expected, but desirable. In the Supplementary Information (Table S2), we report all Pearson pairwise linear correlations having a coefficient of at least  $|r(100)| > .5$  and significant at  $p < .05$ . In `/code/validation/beh/` code is provided to generate Table S2 and a matrix (Supplementary Fig. S2) to visualise these data (`correlate_matrix.py`). Keeping in mind that correlations are just a preliminary way to explore data, most significant correlations were between metrics within the same test, and between these and their cumulative score, when present. As concerns correlations across tasks (marked with \* and \*\* for positive and negative relationships respectively, in Table S2), the most evident patterns emerging are the correlations between reading measures and those typically related to (or predictive of) reading skill or deficit (all sub-scores from the RAN, text reading, word and pseudoword reading, spelling, spoonerisms, phoneme suppression and non-word repetition); between the Brocanto and CVLT long-term recognition score; between AMMA (total and tonal scores) and the musical training score from the MUSEBAQ (time, intensity and level of practice reached). No other measures were as highly and significantly correlated.

Regarding reliability, the main advantage of internal consistency measured via Cronbach alpha is that, as long as covariance is inspected and there are minimal to no *nonignorable nonresponses* (MNAR data, i.e. missing not at random), it can be applied quite flexibly, unlike other methods<sup>102</sup>. Minimal reliability requirements are a matter of debate, as the  $\alpha$  value can be affected by the number of items and improved by the independence



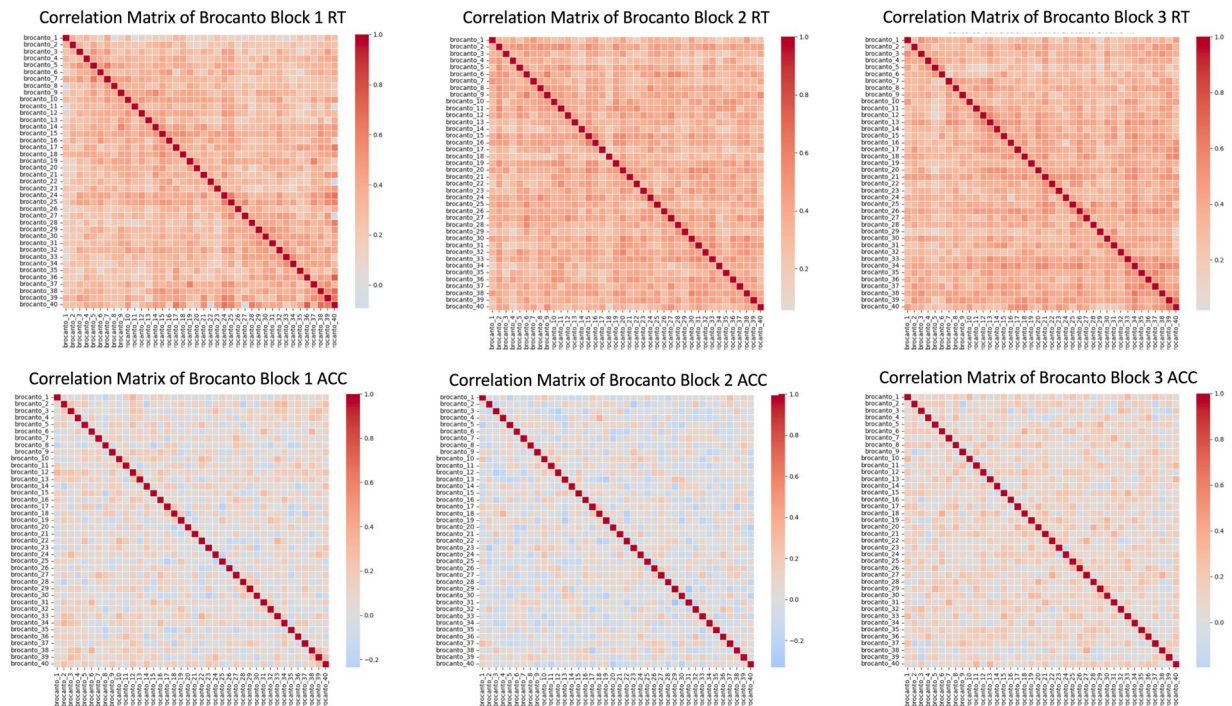
**Fig. 12** Revised tempo test correlation patterns for subtractions (left) and sums (right) RT data. The different number of items across the 2 figures is due to the fact that 1) participants reached a different number of completed operations across the two conditions and 2) no one reached the end of either condition (60 operations). Knowing that for this task, there are data missing for a reason, this can be interpreted as a MNAR pattern.

of the construct these tap into<sup>103–105</sup>. This being said, it is harder to obtain a reliable measurement in tasks that measure *spans* (e.g. Corsi blocks, Digit span) and tasks where participants stop after reaching a maximum of answers with high variability (e.g. Revised tempo test): these will generate a MNAR pattern, where each row ends with a certain number of NaN because the span has been reached, and thus will be less suited to internal consistency measures, providing hardly interpretable coefficients<sup>106</sup>. This can explain the somewhat unreliable Revised tempo test RT, reflected in a pattern where negative or null covariances outweigh positive ones in a MNAR fashion (Fig. 12 and `/code/validation/beh/cronbach_alpha_results.tsv`).

Rapid automatised naming accuracy ( $\alpha = .21$ ) as well as accuracy on the first run of Brocanto ( $\alpha = .25$ ) also had low reliability. We ascribe the former to the difficulty of the task for the participant, and therefore reaction times may provide more reliable measures than accuracy: RAN measures naming latency<sup>107</sup>, whose motor component may show more consistent RTs than accuracy (and this indeed shows in our data, as the two had  $\alpha$  coefficients of .93 and .21, respectively). The lower reliability of the first run of Brocanto is conceptually interesting if compared to the subsequent runs and to the RT data: it shows that upon learning an artificial language completely inductively, responses are more variable at the beginning of the learning process, where people tend to guess more frequently, while the time to make a grammaticality judgment is overall always consistent. We provide the correlation matrices for accuracy and RT across the 3 blocks of Brocanto in Fig. 13. In a few cases, data were recorded at the item level but digitised at the task level (spoonerisms, non-word repetition and phoneme suppression): item-level data are available upon request. To facilitate readers, in the Supplementary Information file we show the internal consistency values for tasks where  $\alpha > .5$ , and tasks with  $\alpha > .6$  are additionally highlighted in bold (Table S3). Complete reliability data are available in `/code/validation/beh/` and can be regenerated with `calculate_cronbach_alpha.py`. Descriptive statistics and plots of the test metrics are available in the Supplementary Information file: Tables S4, S5; Figs. S4–S6).

The Farsi uvular sound production task was rated by two first-language Farsi speakers who heard each recorded utterance from participants and gave it a ‘native-like production score’: therefore, any internal consistency metric would reflect the way the raters scored the task, more than the way the participant performed it (even though these are clearly related). Thus, for this task, we chose to calculate inter-rater reliability between the first and second rater, in a procedure identical to a previous study where this task was used<sup>108</sup>. For the ANT-I task, given its structure (each trial reflecting more than one possible condition from which a score can be derived), we chose to measure the maximal split-half coefficient to assess reliability<sup>109</sup> (the Supplementary Information file contains the code for this procedure and results are reported in Table S3).

**Brain imaging.** *Anatomical imaging anonymisation, QC and brain extraction.* Facial features in anatomical MRI scans violate the principle of anonymity in open data. Therefore, T1-weighted MPRAGE scans were defaced with PyDeface (<https://pypi.org/project/pydeface/>) and their quality was assessed in the Computational Anatomy Toolbox for SPM (hereon, CAT12)<sup>110</sup>. Specifically, we used the CAT12 segmentation and sample homogeneity toolboxes, providing easily interpretable quality measures at the participant and group level. The weighted overall image quality (IQR) and the quartic mean Z-score are the two key indicators of image quality. IQR combines noise and spatial resolution measurements before pre-processing, while the mean quartic Z-score assesses the homogeneity of data after pre-processing, with deviations increasing variance and reducing statistical power. The *product* of IQR and quartic mean Z-scores combines these quality measures, with a low number indicating



**Fig. 13** Top row: correlation matrix for the RT data of Brocanto across 3 blocks. Bottom row: accuracy data for Brocanto across 3 blocks. Plots show that overall, RTs were a more internally consistent metric than accuracy.

high quality. For each participant, we provide a PDF with the CAT12 report, as well as information on sample homogeneity, in `/code/validation/anat/cat12/`. The CAT12 toolbox can be easily run via a graphical user interface (GUI) in Matlab and requires no custom code. We provide the group distribution of the described measures in Fig. 14.

To prepare T1 MPRAGE anatomical scans to be fed to FSL<sup>111</sup> (see *Functional neuroimaging* section), we provide code for running improved skull-stripping as a loop in the NEBULA101 dataset in `/code/preprocessing/1_anat/1_optiBET/nebula101_run_optiBET.sh`. This code will process T1-weighted MRI images stored in a BIDS directory to perform improved brain extraction with optiBET<sup>84</sup> via the following actions:

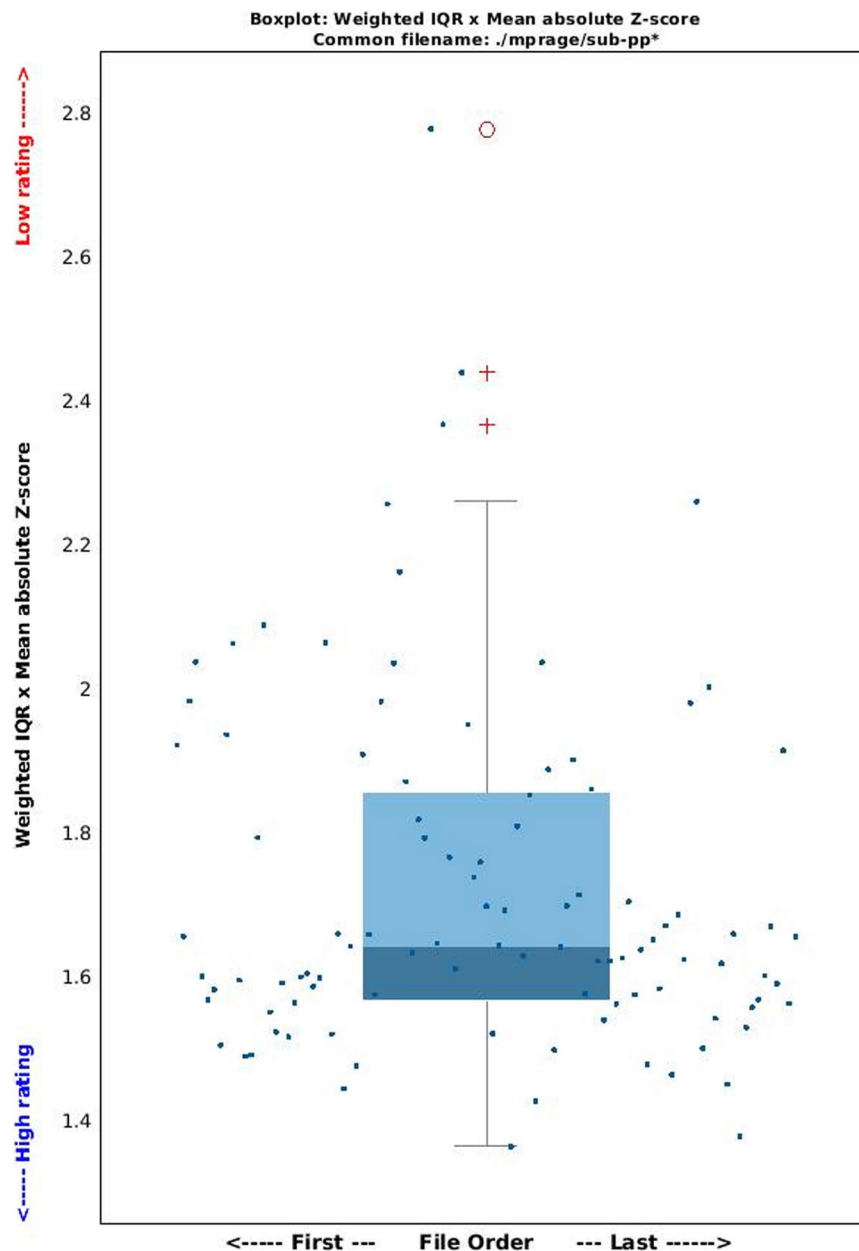
- 1) Iterating through all participant directories in the BIDS directory.
- 2) For each participant, creating an output directory in the derivatives folder if it does not yet exist.
- 3) Copying the T1-weighted images (`sub-pp*_ses-01_T1w.nii.gz`) to the corresponding output directory.
- 4) Checking if each T1-weighted image has already been processed.
- 5) If an image has not been processed, running the `nebula101_optiBET.sh` master script on the image, which performs actual brain extraction.

The script referenced in step 5 is the master optiBET script, which needs to be downloaded from <https://montilab.psych.ucla.edu/fmri-wiki/optibet/> and stored in the location where the loop script is stored.

**Field maps.** We acquired field maps to accompany the DWI and fMRI sequences. Field map preparation consists of a few steps that are partly scanner-specific, and in the case of Siemens Prisma, consisted in recomposing the phase and magnitude images. We provide field map preprocessing code in `/code/preprocessing/2_fmmap/fsl/1_nebula101_run_fmmap_prep_all.sh`

The script will perform the following steps:

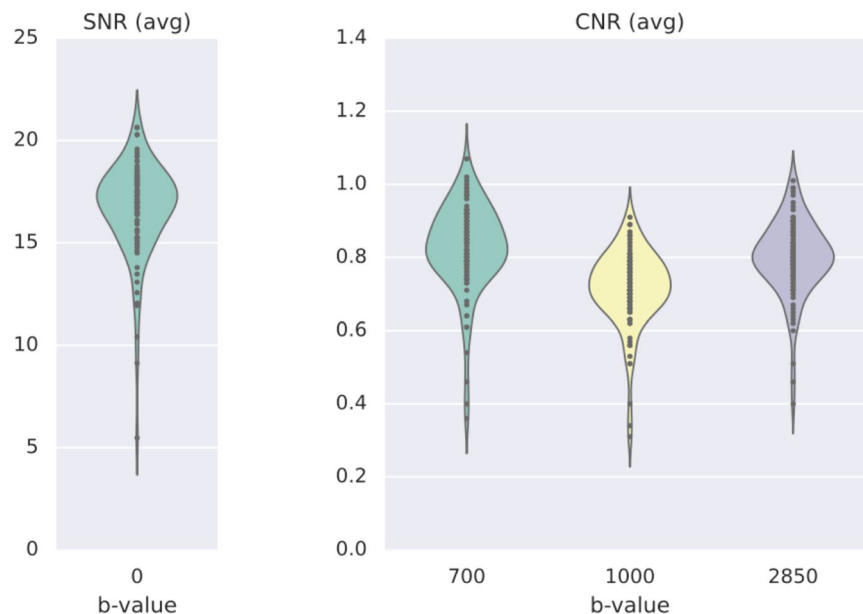
- 1) Transformation Calculation: calculating the transformation matrix from the anatomical T1 image to the magnitude image and applying the calculated transformation to the T1 brain mask to align it with the magnitude image (FLIRT).
- 2) Magnitude Image Masking: using the transformed brain mask to mask the magnitude image, creating an untrimmed brain image (`fslmaths`).
- 3) Brain Mask Trimming (`fslmaths`) consisting of the following steps:
  - a. Binarising the untrimmed mask.
  - b. Smoothing the mask with an 8 mm kernel.
  - c. Thresholding the smoothed mask at 75%.
  - d. Binarising the thresholded mask.



**Fig. 14** CAT12 output. Boxplot showing the weighted IQR by mean absolute z-score scaled by a factor of 4 (i.e. quartic) to emphasize outliers. As the plot shows, all but 1 participant in this distribution lie in the good to optimal range of the quality measure. Of note, even when outliers are detected, these are not necessarily data to discard on an absolute basis, as the score is calculated relatively to the specific sample being analysed.

- e. Trimming the original brain image using the final binary mask.
  - f. Removing intermediate untrimmed files.
- 4) Preparing the final field map (`fsl_prepare_fieldmap`) using the processed magnitude and phase difference images, with a specified echo spacing of 2.46 ms.

**Diffusion-weighted imaging.** To validate our diffusion data, we ran the QUAD (participant) and SQUAD (group) programs within the FSL-FDT<sup>111,112</sup> suite as part of the diffusion pre-processing pipeline. The DWI data and prepared field map of each participant were fed to EDDY to correct for susceptibility, eddy currents, inter- and intra-volume displacement and signal dropout, using the prepared field map via the `-field` flag. Finally, we ran `eddyQC`<sup>113</sup> for quality assessment. We provide a pdf with QC statistics for each NEBULA101 participant and for the group, as well as the code to reproduce this procedure, in `/code/validation/dwi/fsl/dwi_qc_quad_squad.py`. The code is configured to automatically grab the acquisition parameters and the correct number of shells from the DWI sidecars to generate the `acqparam.txt` and `index.txt` files required by



**Fig. 15** Signal-to-noise and contrast-to-noise ratios from the QC-ed DWI data.

EDDY, and will raise flags if the expected numbers do not match with the data. An `eddy_quad_gc_paths.txt` file has also been set up to this aim in the same location, which the user can alter for their needs, should they decide to preprocess DWI data (which we provide in raw form) with the same pipeline (`/code/preprocessing/dwi/fsl/1_dwi_prep.py`). Fig. 15 reports the signal- and contrast-to-noise ratios from the group QC metrics generated using SQUAD. The complete PDF report is available in `code/validation/dwi/fsl/squad/group_qc.pdf`.

**Functional imaging.** As concerns fMRI, we provide QC metrics obtained from preprocessing of the whole dataset in FSL, and the code to reproduce the procedure on the NEBULA101 sample.

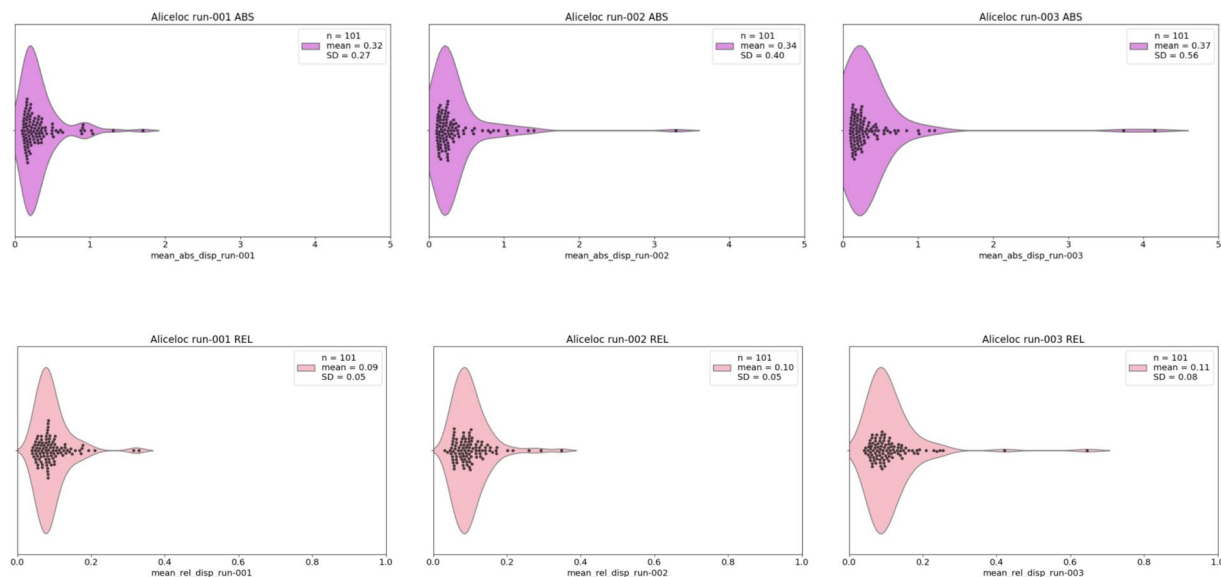
Figs. 16, 17 report violin plots of the absolute and relative displacement by run for the language localiser and for the resting-state fMRI sequence, respectively. Absolute displacement reflects the movement with respect to the reference volume, while relative displacement reports volume-to-volume information. Both can impact SNR but while absolute displacement can be counteracted, relative displacement is usually more harmful, and both must be checked when assessing the quality of fMRI data<sup>14</sup>. Displacement information and plots (Figs. 16, 17) can be obtained by running `code/validation/func/fsl/1_copy_motion_params.py` and `code/validation/func/fsl/2_mean_abs_rel_disp_violin.py`. In principle, if the data are not being reprocessed, `1_copy_motion_params` does not need to be rerun as we provide `.rms` output from FSL, containing absolute and relative displacement data for each participant, in `code/validation/func/fsl/aliceloc/and/rest`.

A folder containing the code for performing the basic preprocessing steps of fMRI for the language localiser and resting-state sequences has been set up in `/code/preprocessing/3_func/fsl/1_preprocl/1_loop_fsl_preprocl.py`. This code performs the first step of fMRI preprocessing with FSL FEAT, consisting of:

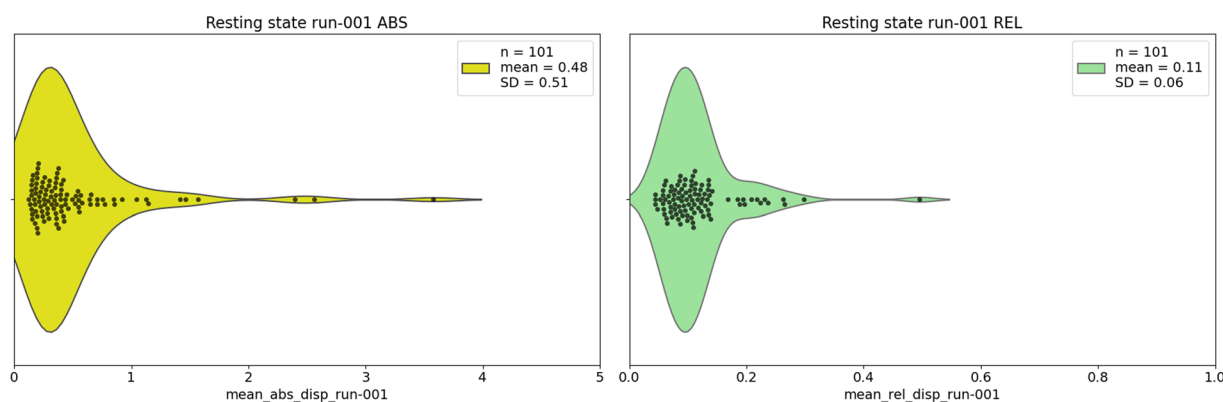
- 1) Motion correction
- 2) 4D mean intensity normalization
- 3) Spatial smoothing (5 mm FWHM)
- 4) B0 unwarping (includes BBR `reg`)
- 5) Registration (BBR, FNIRT)

This code will run using `/code/preprocessing/3_func/fsl/1_preprocl/design_pp.fsf`, a design file required by FSL which has been set up to contain information on the preprocessing loop, should the user wish to preprocess our data with the same pipeline.

**Overall validation of data structure.** This dataset is mostly complete, with minimal missing information. However, to provide the user with information that is traceable without necessarily having this data descriptor at hand, we have created a series of scripts within the `/code/data_checks/` folder aimed at providing report logs on the presence of behavioural and brain imaging data. Such logs are saved to `/derivatives/validation/data_checks/`.



**Fig. 16** Absolute (ABS) and relative (REL) displacement (mm) across runs of the language localiser.



**Fig. 17** Absolute (ABS) and relative (REL) displacement (mm) during the single resting-state sequence run.

- `nebula101_bids_report_beh.py`: this script plots behavioural data presence checks on the raw and derivate dataset and compiles the findings to a report log called `nebula101_beh_report_[timestamp].txt`. It will also create:
  - A tabular file `nebula101_datalist_[timestamp].tsv` listing data as 1 and 0 if present or absent, for easy reference.
  - A heatmap with the above information, called `data_presence_heatmap_[timestamp].png`.
  - A report specific to the LEAPQ data, to have missing and present information at a glance, called `nebula101_leapq_data_report_[timestamp].txt`
- `nebula101_bids_report_mri.py`: this script plots imaging data presence checks and compiles the findings to a report log called `nebula101_mri_report_[timestamp].txt`. This can be used for easy referencing to participants having undergone different conditions (e.g. different field map, see: Usage Notes).

### Usage Notes

**Behavioural data.** This dataset is mostly complete, but a few technical malfunctions during testing have caused some minimal missing data, as described in Table 4. This information and a heatmap for visualising missing data at a glance can easily be regenerated using the scripts in `/code/validation/data_checks/` called `nebula101_bids_report_beh.py` and `nebula101_missing_data_log.py`.

Participant ID	Columns with missing data at task level
sub-pp031	phon_suppr_rt_manual
sub-pp050	brocanto_corr1, brocanto_incorr1, brocanto_rt1, cvlt_long_corr, cvlt_long_incorr, cvlt_long_rt, brocanto
sub-pp116	cvlt_long_rt
sub-pp118	apm_corr, apm_incorr, apm_rt
sub-pp155	non_word_rep_span, non_word_rep_acc

**Table 4.** Missing behavioural data.

**Imaging data.** Here we provide some details for easier navigation of the imaging data. This information can be generated by `/code/validation/data_checks/nebula101_bids_report_mri.py`. JSON sidcar files have already been adjusted to report the below nuances:

- sub-pp161: for this participant, the `beta` field map must be used for resting-state fMRI processing, and the `delta` field map must be used for DWI processing.
- sub-pp001, sub-pp010, sub-pp013, sub-pp023, sub-pp043, sub-pp053: for these participants, the `gamma` field map must be used for resting-state fMRI processing.
- sub-pp032: the voxel dimensions of `/sub-pp032/ses-01/anat/sub-pp032_ses-01_rec-defaced_T1w.nii.gz` slightly differ in the `y` and `z` directions (`pixdim1 = 1.000000`, `pixdim2 = 1.039060`, `pixdim3 = 1.039060`).

When running the BIDS validator, this information is reported as *warnings* and saved in the log, unless the `ignoreWarnings` flag is raised. Warnings will not impede BIDS validation and the dataset will pass it every time. When processing the data, this information will be read automatically by software that can work primarily in BIDS (such as `fMRIPrep`<sup>115</sup> in the case of functional imaging). When using software that lacks this functionality, such as FSL, we recommend processing these participants manually. The `.bidsignore` file contains instructions for skipping the validation of `/neurobagel` as it contains extra files.

### Code availability

All code that can be rerun on these data is provided within the `/code` folder, as described in the Technical Validation section. Raw-to-derivative code is provided for all raw data. Given that we do not provide *source* behavioural data, source-to-raw preprocessing code for tasks and questionnaires is not included, but can be shared upon request.

Received: 2 August 2024; Accepted: 20 December 2024;

Published online: 06 January 2025

### References

- Kidd, E., Donnelly, S. & Christiansen, M. H. Individual Differences in Language Acquisition and Processing. *Trends Cogn Sci* **22**, 154–169 (2018).
- Vogel, E. K., & Edward A. How to exploit diversity for scientific gain: Using individual differences to constrain cognitive theory. *Curr Dir Psych Sci* **17.2**, 171–176 (2008)
- Carroll, J. B. & Sapon, S. M. *Modern Language Aptitude Test. Modern language aptitude test.* (Psychological Corporation, San Antonio, TX, US, 1959).
- Carroll, J. B. Twenty-five years of research on foreign language aptitude. in *Individual differences and universals in language learning aptitude* 867–873 (1981).
- Stansfield, C. W. & Reed, D. J. The Story Behind the Modern Language Aptitude Test: An Interview With John B. Carroll (1916–2003). *Lang Assess Q* **1**, 43–56 (2004).
- Miller, G. A. The cognitive revolution: a historical perspective. *Trends Cogn Sci* **7**, 141–144 (2003).
- Li, S. The Associations between Language Aptitude and Second Language Grammar Acquisition: A Meta-Analytic Review of Five Decades of Research. *Appl Linguist* **36**, 385–408 (2015).
- Harré, R. The Second Cognitive Revolution. in *After cognitivism* (ed. Leidlmaier, K.) 182–187 (Springer, Dodrecht, 2009).
- Robinson, P. Individual differences, aptitude complexes, SLA processes, and aptitude test development. *Second Language Learning and Teaching* **4**, 57–75 (2012).
- Robinson, P. Aptitude and second language acquisition. *Annu Rev Appl Linguist* **25**, 46–73 (2005).
- Wen, Z., Biedroń, A. & Skehan, P. Foreign Language Aptitude Theory: Yesterday, Today and Tomorrow. *Language Teaching* **50** (2017).
- Eisenstein, M. Childhood bilingualism and adult language learning aptitude. *International Review of Applied Psychology* **29**, 159–172 (1980).
- Grigorenko, E. L., Sternberg, R. J. & Ehrman, M. E. A theory-based approach to the measurement of foreign language learning ability: The canal-F theory and test. *Modern Language Journal* **84**, 390–405 (2000).
- Sparks, R. L., Ganschow, L. & Patton, J. Prediction of Performance in First-Year Foreign Language Courses: Connections Between Native and Foreign Language Learning. *J Educ Psychol* **87**, 638–655 (1995).
- Harley, B. & Hart, D. Language aptitude and second language proficiency in classroom learners of different starting ages. *Stud Second Lang Acquis* **19**, 379–400 (1997).
- Sawyer, M. Language Aptitude and Language Experience: Are They Related? *International University of Japan, Departmental Bulletin Paper* (1992).
- Rampinini, A., Balboni, I., Golestani, N. & Berthele, R. A behavioural exploration of language aptitude and experience, cognition and more using Graph Analysis. *Brain Res* **1842**, 149109 (2024).
- Kecskes, I. Dual and multilingual systems. *International Journal of Multilingualism* **7**, 91–109 (2010).
- Ellis, E. M. Defining and investigating monolingualism. *Sociolinguistic Studies*, **2**, 311–330 (2008).

20. Kirk, N. W. MIND your language(s): Recognizing Minority, Indigenous, Non-standard(ized), and Dialect variety usage in monolinguals. *Appl Psycholinguist* **44**, 358–364 (2023).
21. Leivada, E., Rodríguez-Ordóñez, I., Parafita Couto, M. C. & Perpiñán, S. Bilingualism with minority languages: Why searching for unicorn language users does not move us forward. *Appl Psycholinguist* **44**, 384–399 (2023).
22. Higby, E., Kim, J. & Obler, L. K. Multilingualism and the brain. *Annual Review of Applied Linguistics* **33**, 68–101 (2013).
23. Blasi, D. E., Henrich, J., Adamou, E., Kemmerer, D. & Majid, A. Over-reliance on English hinders cognitive science. *Trends in Cognitive Sciences* vol. 26 1153–1170 (2022).
24. Collart, A. & Collart, A. A decade of language processing research: Which place for linguistic diversity? *Glossa Psycholinguistics* **3**, (2024).
25. Bambini, V. & Canal, P. Neurolinguistic research on the Romance languages *osf.io* Preprint at <https://doi.org/10.31219/osf.io/c9yxn> (2021).
26. Berthele, R. Introduction: What's Special About Multilingualism? *Language Learning* **71**, 5–11 (2021).
27. Berthele, R. The Extraordinary Ordinary: Re-engineering Multilingualism as a Natural Category. *Lang Learn* **71**, 80–120 (2021).
28. Berthele, R. The selective celebration of linguistic diversity: evidence from the Swiss language policy discourse. *J Multiling Multicult Dev* **42**, 125–136 (2021).
29. Udry, I. & Berthele, R. The smart, the motivated and the self-confident: The role of language aptitude, cognition, and affective variables in early instructed foreign language learning. in *Individual Differences in Early Instructed Language Learning: The Role of Language Aptitude, Cognition, and Motivation* (eds. Berthele, R. & Udry, I.) 71–90 (Language Science Press, Berlin, 2021).
30. Evans, N. & Levinson, S. C. The myth of language universals: Language diversity and its importance for cognitive science. *Behavioral and Brain Sciences* **32**, 429–448 (2009).
31. Dąbrowska, E. Experience, aptitude and individual differences in native language ultimate attainment. *Cognition* **178**, 222–235 (2018).
32. Turker, S., Seither-Preisler, A. & Reiterer, S. M. Examining Individual Differences in Language Learning: A Neurocognitive Model of Language Aptitude. *Neurobiology of Language* **2**, 389–415 (2021).
33. Van Der Maas, H. L. *et al.* A Dynamical Model of General Intelligence: The Positive Manifold of Intelligence by Mutualism. *Psychol Rev* **113**(4), 842–61 (2006).
34. Pluck, G. & Cerone, A. A Demonstration of The Positive Manifold of Cognitive Test Inter-correlations, and how it Relates to General Intelligence, Modularity, and Lexical Knowledge. *UC Merced Proceedings of the Annual Meeting of the Cognitive Science Society* (2021).
35. Kong, X. Z. *et al.* Gene Expression Correlates of the Cortical Network Underlying Sentence Processing. *Neurobiology of Language* **1**, 77–103 (2020).
36. Amelink, J. S. *et al.* Imaging genetics of language network functional connectivity reveals links with language-related abilities, dyslexia and handedness. *Commun Biol* **7**, 1209 (2024).
37. Nayak, S. *et al.* The Musical Abilities, Pleiotropy, Language, and Environment (MAPLE) Framework for Understanding Musicality-Language Links Across the Lifespan. *Neurobiology of Language* **3**, 615–664 (2022).
38. Steyer, R., Mayer, A., Geiser, C. & Cole, D. A. A Theory of States and Traits-Revised. *Ann Rev Psych* **11**, 71–98 (2014).
39. Fuchs, E. & Flügge, G. Adult Neuroplasticity: More Than 40 Years of Research. *Neural Plast* **2014**, 541870 (2014).
40. Dityatev, A. & Schachner, M. Extracellular matrix molecules and synaptic plasticity. *Nature Reviews Neuroscience* **2003** 4:6 **4**, 456–468 (2003).
41. Ramirez-Amaya, V. Molecular Mechanisms of Synaptic Plasticity Underlying Long-Term Memory Formation. *Neural Plasticity and Memory: From Genes to Brain Imaging* 47–66 (2007).
42. Grafman, J. Conceptualizing functional neuroplasticity. *J Commun Disord* **33**, 345–356 (2000).
43. Wang, R. *et al.* Functional and structural neuroplasticity associated with second language proficiency: An MRI study of Chinese-English bilinguals. *J Neurolinguistics* **56**, 100940 (2020).
44. DeLuca, V., Rothman, J., Bialystok, E. & Pliatsikas, C. Redefining bilingualism as a spectrum of experiences that differentially affects brain structure and function. *Proc Natl Acad Sci USA* **116**, 7565–7574 (2019).
45. Pliatsikas, C. Multilingualism and Brain Plasticity. *The Handbook of the Neuroscience of Multilingualism* 230–251 (2019).
46. Li, P., Legault, J. & Litcofsky, K. A. Neuroplasticity as a function of second language learning: Anatomical changes in the human brain. *Cortex* **58**, 301–324 (2014).
47. DeFelipe, J. Brain plasticity and mental processes: Cajal again. *Nat Rev Neurosci* **7**, 811–817 (2006).
48. Grasby, K. L. *et al.* The genetic architecture of the human cerebral cortex. *Science* **367**, 6484 (2020).
49. Navarri, X. *et al.* A biologically informed polygenic score of neuronal plasticity moderates the association between cognitive aptitudes and cortical thickness in adolescents. *Dev Cogn Neurosci* **60**, 101232 (2023).
50. Petersen, C. R. & Al-Haik, A. R. The Development of the Defense Language Aptitude Battery (DLAB). *Educ Psychol Meas* **36**, 369–380 (1976).
51. Wilkinson, M. D. *et al.* The FAIR Guiding Principles for scientific data management and stewardship. *Sci Data* **2016** **3**, 1–9 (2016).
52. Poldrack, R. A. & Gorgolewski, K. J. Making big data open: Data sharing in neuroimaging. *Nat Neurosci* **17**, 1510–1517 (2014).
53. Poldrack, R. A. & Poline, J. B. The publication and reproducibility challenges of shared data. *Trends Cog Sci* **19**, 59–61 (2015).
54. Botvinik-Nezer, R. *et al.* Variability in the analysis of a single neuroimaging dataset by many teams. *Nature* **582**, 84–88 (2020).
55. Milham, M. P. *et al.* Assessment of the impact of shared brain imaging data on the scientific literature. *Nat Comm* **9** (2018).
56. Gorgolewski, K. J. *et al.* The brain imaging data structure, a format for organizing and describing outputs of neuroimaging experiments. *Sci Data* **3**, (2016).
57. Eickhoff, S., Nichols, T. E., Van Horn, J. D. & Turner, J. A. Sharing the wealth: Neuroimaging data repositories. *NeuroImage* **124**, 1065–1068 (2016).
58. Markiewicz, C. J. *et al.* The openneuro resource for sharing of neuroscience data. *Elife* **10**, (2021).
59. Schoffelen, J. M. *et al.* A 204-subject multimodal neuroimaging dataset to study language processing. *Sci Data*, **6**(1), 17 (2019)
60. Bhattasali, S., Brennan, J., Luh, W. M., Franzluebbers, B. & Hale, J. The Alice Datasets: fMRI & EEG observations of natural language comprehension. In *Proceedings of the Twelfth Language Resources and Evaluation Conference* (2020).
61. Hanke, M. *et al.* A high-resolution 7-Tesla fMRI dataset from complex natural stimulation with an audio movie. *Sci Data* **1**, (2014).
62. Lipkin, B. *et al.* Probabilistic atlas for the language network based on precision fMRI data from > 800 individuals. *Sci Data* **9**, 529 (2022).
63. Isaieva, K. *et al.* Multimodal dataset of real-time 2D and static 3D MRI of healthy French speakers. *Sci Data* **8**(1), 258 (2021).
64. Lim, Y. *et al.* A multispeaker dataset of raw and reconstructed speech production real-time MRI video and 3D volumetric images. *Sci Data* **8**(1), 187 (2021).
65. Gomez-Marin, A., Paton, J. J., Kampff, A. R., Costa, R. M. & Mainen, Z. F. Big behavioral data: Psychology, ethology and the foundations of neuroscience. *Nat Neurosci* **17**, 1455–1462 (2014).
66. Rasgado-Toledo, J. *et al.* A Dataset to Study Pragmatic Language and Its Underlying Cognitive Processes. *Front Hum Neurosci* **15**, 666210 (2021).
67. Marian, V., Blumenfeld, H. K. & Kaushanskaya, M. The Language Experience and Proficiency Questionnaire (LEAP-Q): Assessing language profiles in bilinguals and multilinguals. *Journal of Speech, Language, and Hearing Research* **50**(4), 940–967 (2007).

68. Hintz, F., Dijkhuis, M., van 't Hoff, V., McQueen, J. M. & Meyer, A. S. A behavioural dataset for studying individual differences in language skills. *Sci Data* **7**, (2020).
69. Berthele, R. & Udry, I. *Individual Differences in Early Instructed Language Learning: The Role of Language Aptitude, Cognition, and Motivation*. (Language science press, 2021).
70. Turner, B. O., Paul, E. J., Miller, M. B., & Barbey, A. K. Small sample sizes reduce the replicability of task-based fMRI studies. *Communications biology* **1**, 62 (2018).
71. Malik-Moraleda, S. *et al.* An investigation across 45 languages and 12 language families reveals a universal language network. *Nat Neurosci* **25**, 1014–1019 (2022).
72. OECD. *Education at a Glance 2023: OECD Indicators*. <https://doi.org/10.1787/e13bef63-en> (OECD, 2023).
73. OECD. *Education at a Glance 2023 Sources, Methodologies and Technical Notes*. (Organisation for Economic Co-operation and Development, Paris, 2023).
74. Federal Statistical Office Section Demography and Migration. Swiss Federal Population Census Structural Survey - *Language* (2022).
75. Milfont, T. L. & Klein, R. A. Replication and Reproducibility in Cross-Cultural Psychology. *Cultural Psychology* (2018).
76. Barratt, W. The Barratt simplified measure of social status (BSMSS): Measuring SES. *Unpublished manuscript, Indiana State University* (2006).
77. Elson, M., Hussey, I., Alsalti, T. & Arslan, R. C. Psychological measures aren't toothbrushes. *Communications Psychology* **2023** 1:1 1, 1–4 (2023).
78. Pimsleur, P., Reed, D. J. & Stansfield, C. W. *Pimsleur Language Aptitude Battery: PLAB: Manual*. (Second Language Testing Foundation, 2004).
79. Pimsleur, P. *Pimsleur Language Aptitude Battery (Form S)*. (Harcourt, Brace and world, Incorporated, 1966).
80. Woods, K. J. P., Siegel, M. H., Traer, J. & McDermott, J. H. Headphone screening to facilitate web-based auditory experiments. *Atten Percept Psychophys* **79**, 2064–2072 (2017).
81. Zhao, S., Brown, C. A., Holt, L. L. & Dick, F. Robust and Efficient Online Auditory Psychophysics. *Trends Hear* **26**, 23312165221118790 (2022).
82. Anwyll-Irvine, A. L., Massonnié, J., Flitton, A., Kirkham, N. & Evershed, J. K. Gorilla in our midst: An online behavioral experiment builder. *Behav Res Methods* **52**, 388–407 (2020).
83. Malik-Moraleda, S. *et al.* The universal language network: A cross-linguistic investigation spanning 45 languages and 12 language families. *bioRxiv* 2021.07.28.454040 <https://doi.org/10.1101/2021.07.28.454040> (2022).
84. Balboni, I., Rampinini, A., Kepinska, O., Berthele, R. & Golestani, N. Brain activation for language and its relationship to cognitive and linguistic measures: a multimodal exploration. in *Society for the Neurobiology of Language Annual Meeting* (Marseille, France, 2023).
85. Liu, X. & Yang, L. Individual differences in the language task-evoked and resting-state functional networks. *Front Hum Neurosci* **17**, 1283069 (2023).
86. Deng, Z., Chandrasekaran, B., Wang, S. & Wong, P. C. M. Resting-state low-frequency fluctuations reflect individual differences in spoken language learning. *Cortex* **76**, 63–78 (2016).
87. Achal, S., Hoefl, F. & Bray, S. Individual Differences in Adult Reading Are Associated with Left Temporo-parietal to Dorsal Striatal Functional Connectivity. *Cerebral Cortex* **26**(10), 4069–4081 (2016).
88. Zhang, G. *et al.* Individual differences in first-pass fixation duration in reading are related to resting-state functional connectivity. *Brain Lang* **213**, (2021).
89. Reineberg, A. E., Gustavson, D. E., Benca, C., Banich, M. T. & Friedman, N. P. The relationship between resting state network connectivity and individual differences in executive functions. *Front Psychol* **9**, 361864 (2018).
90. Reineberg, A. E., Andrews-Hanna, J. R., Depue, B. E., Friedman, N. P. & Banich, M. T. Resting-state networks predict individual differences in common and specific aspects of executive function. *Neuroimage* **104**, 69–78 (2015).
91. Queder, N. *et al.* NIDM-Terms: community-based terminology management for improved neuroimaging dataset descriptions and query. *Front Neuroinform* **17**, (2023).
92. Rampinini, A., Balboni, I., Kepinska, O., Bertele, R. & Golestani, N. NEBULA101 NeuroBehavioural Understanding of Language Aptitude. *OpenNeuro dataset* <https://doi.org/10.18112/openneuro.ds005613.v1.0.1>.
93. Nordhoff, S. & Hammarström, H. Glottolog/Langdoc: Defining dialects, languages, and language families as collections of resources. in *First International Workshop on Linked Science 2011-In conjunction with the International Semantic Web Conference (ISWC 2011)* (2011).
94. Gullifer, J. W. & Titone, D. Characterizing the social diversity of bilingualism using language entropy. *Bilingualism: Language and Cognition* **23**, 283–294 (2020).
95. Kepinska, O. *et al.* Language combinations of multilinguals are reflected in their first-language knowledge and processing. *Sci Rep* **13**, 1947 (2023).
96. Shannon, C. E. A mathematical theory of communication. *The Bell system technical journal* **27**, 379–423 (1948).
97. Hausser, J., Strimmer, K. & Strimmer, M. K. Package 'entropy'. *R Foundation for Statistical Computing: Vienna, Austria* (2012).
98. Allen, M. P. The problem of multicollinearity. *Understanding regression analysis* 176–180 (1997).
99. Chan, J. Y. L. *et al.* Mitigating the multicollinearity problem and its machine learning approach: a review. *Mathematics* **10**, 1283 (2022).
100. Slinker, B. K. & Glantz, S. A. Multiple regression for physiological data analysis: the problem of multicollinearity. *American Journal of Physiology-Regulatory, Integrative and Comparative Physiology* **249**, R1–R12 (1985).
101. Paul, R. K. Multicollinearity: Causes, effects and remedies. *IASRI, New Delhi* **1**, 58–65 (2006).
102. Green, S. B. & Yang, Y. Evaluation of Dimensionality in the Assessment of Internal Consistency Reliability: Coefficient Alpha and Omega Coefficients. *Educational Measurement: Issues and Practice* **34**, 14–20 (2015).
103. Vaske, J. J., Beaman, J. & Sponarski, C. C. Rethinking Internal Consistency in Cronbach's Alpha. *Leis Sci* **39**, 163–173 (2017).
104. Drennan, J. Quantitative health research: issues and methods. in (ed. Curtis, E.) (McGraw-Hill Education UK, 2013).
105. Watson, R. Issues and debates in validity and reliability. in *Quantitative Health Research: Issues And Methods* (eds. Curtis, E. & Drennan, J.) (McGraw-Hill Education UK, 2013).
106. Cronbach, L. J. & Hartmann, W. A note on negative reliabilities. *Educ Psychol Meas* **14**, 342–346 (1954).
107. Wile, T. L. & Borowsky, R. What does rapid automatized naming measure? A new RAN task compared to naming and lexical decision. in *Brain and Language* **90**, 47–62 (2004).
108. Golestani, N. & Pallier, C. Anatomical Correlates of Foreign Speech Sound Production. *Cerebral Cortex* **17**, 929–934 (2006).
109. Thompson, B. L., Green, S. B. & Yang, Y. Assessment of the maximal split-half coefficient to estimate reliability. *Educ Psychol Meas* **70**, 232–251 (2010).
110. Gaser, C. *et al.* CAT: a computational anatomy toolbox for the analysis of structural MRI data. *Gigascience* **13** (2024)
111. Jenkinson, M., Beckmann, C. F., Behrens, T. E. J., Woolrich, M. W. & Smith, S. M. FSL. *Neuroimage* **62**, 782–790 (2012).
112. Jbabdi, S., Sotiropoulos, S. N., Savio, A. M., Graña, M. & Behrens, T. E. J. Model-based analysis of multishell diffusion MR data for tractography: how to get over fitting problems. *Magn Reson Med* **68**, 1846–1855 (2012).
113. Bastiani, M. *et al.* Automated quality control for within and between studies diffusion MRI data using a non-parametric framework for movement and distortion correction. *Neuroimage* **184**, 801–812 (2019).

114. Friston, K. J., Williams, S., Howard, R., Frackowiak, R. S. J. & Turner, R. Movement-Related effects in fMRI time-series. *Magn Reson Med* **35**, 346–355 (1996).
115. Esteban, O. *et al.* fMRIPrep: a robust preprocessing pipeline for functional MRI. *Nat Methods* **16**, 111–116 (2019).
116. Rasgado-Toledo, J. *et al.* Pragmatic Language. *OpenNeuro dataset* <https://doi.org/10.18112/openneuro.ds003481.v1.0.3> (2021).
117. Bathelt, J., Taylor, J. & Rastle, K. Language fMRI. *OpenNeuro dataset* <https://doi.org/10.18112/openneuro.ds004765.v1.0.0> (2023).
118. Bathelt, J., Rastle, K. & Taylor, J. S. H. Relationship between resting state functional connectivity and reading-related behavioural measures in 69 adults. *Neurobiology of Language*, 1–19 (2024).
119. Rogers, C. S. *et al.* Age-related differences in auditory cortex activity during spoken word recognition. *OpenNeuro dataset* <https://doi.org/10.18112/openneuro.ds002382.v1.0.1> (2022).
120. Rogers, C. S. *et al.* Age-Related Differences in Auditory Cortex Activity During Spoken Word Recognition. *Neurobiology of Language* **1**, 452–473 (2020).
121. Rogers, C. S., Jones, M. S., McConkey, S. & Peelle, J. E. Listening task. *OpenNeuro dataset* <https://doi.org/10.18112/openneuro.ds004285.v1.0.0> (2022).
122. Woodhead, Z. *et al.* Comparing language lateralisation using fMRI and fTCD. *OpenNeuro dataset* <https://doi.org/10.18112/openneuro.ds004073.v1.0.1> (2023).
123. Bishop, D. V. M., Woodhead, Z. V. J. & Watkins, K. E. Approaches to Measuring Language Lateralisation: An Exploratory Study Comparing Two fMRI Methods and Functional Transcranial Doppler Ultrasound. *Neurobiology of Language* **5**, 409–431 (2024).
124. Gold, C. E., Howell, A. L., Burdis, J., Kirwan, C. B. & Thompson, G. L. Exploring the Resting State Neural Activity of Monolinguals and Late and Early Bilinguals. *OpenNeuro dataset* <https://doi.org/10.18112/openneuro.ds001747.v1.1.0> (2023).
125. Gold, C. Exploring the Resting State Neural Activity of Monolinguals and Late and Early Bilinguals. (Brigham Young University, 2018).
126. DeLuca, V. & Plitsikas, C. Bilingualism and the brain. *OpenNeuro dataset* <https://doi.org/10.18112/openneuro.ds001796.v1.7.0> (2022).
127. Nastase, S. A. *et al.* Narratives. *OpenNeuro dataset* <https://doi.org/10.18112/openneuro.ds002345.v1.1.4> (2020).
128. Nastase, S. A. *et al.* The “Narratives” fMRI dataset for evaluating models of naturalistic language comprehension. *Sci Data* **8**, (2021).
129. Li, J., Hale, J. & Pallier, C. Le Petit Prince: A multilingual fMRI corpus using ecological stimuli. *OpenNeuro dataset* <https://doi.org/10.18112/openneuro.ds003643.v2.0.5> (2024).
130. Li, J. *et al.* Le Petit Prince multilingual naturalistic fMRI corpus. *Sci Data* **9**, 1–15 (2022).
131. Rodriguez-Fornells, A., Kramer, U., Lorenzo-Seva, U., Festman, J. & Münte, T. Self-Assessment of Individual Differences in Language Switching. *Front Psychol* **2**, (2012).
132. Ryan, S. The ideal L2 selves of Japanese learners of English. (University of Nottingham, 2008).
133. Thompson, A. S. & Lee, J. The Motivational Factors Questionnaire in the Korean EFL context: predicting group membership according to English proficiency and multilingual status. *Language Learning Journal* **46**, 398–414 (2018).
134. Lefly, D. L. & Pennington, B. F. Reliability and validity of the adult reading history questionnaire. *J Learn Disabil* **33**, 286–296 (2000).
135. Roebuck, H. & Lupyan, G. The Internal Representations Questionnaire: Measuring modes of thinking. *Behav Res Methods* **52**, 2053–2070 (2020).
136. Chin, T. C., Coutinho, E., Scherer, K. R. & Rickard, N. S. MUSEBAQ: A modular tool for music research to assess musicianship, musical capacity, music preferences, and motivations for music use. *Music Percept* **35**, 376–399 (2018).
137. Rakesh, D. & Whittle, S. Socioeconomic status and the developing brain—A systematic review of neuroimaging findings in youth. *Neurosci Biobehav Rev* **130**, 379–407 (2021).
138. Oldfield, R. C. The assessment and analysis of handedness: The Edinburgh inventory. *Neuropsychologia* **9**, 97–113 (1971).
139. Nedjar, T., Touari, M., Mesbah, M., Lellouch, J. & Dellatolas, G. La préférence manuelle dans une population d'étudiants algériens francophones et comparaison avec la population française. *Annee Psychol* **89**, 239–253 (1989).
140. Stansfield, C. W. Language Aptitude Reconsidered. *ERIC Digest* (1989).
141. Golestani, N. & Pallier, C. Anatomical correlates of foreign speech sound production. *Cerebral Cortex* **17**, 929–934 (2007).
142. Golestani, N., Paus, T. & Zatorre, R. J. Anatomical Correlates of Learning Novel Speech Sounds. *Neuron* **35**, 997–1010 (2002).
143. Kepinska, O., de Rover, M., Caspers, J. & Schiller, N. O. On neural correlates of individual differences in novel grammar learning: An fMRI study. *Neuropsychologia* **98**, 156–168 (2017).
144. Opitz, B. & Friederici, A. D. Interactions of the hippocampal system and the prefrontal cortex in learning language-like rules. *Neuroimage* **19**, 1730–1737 (2003).
145. Raven, J. C. Matrices Progressives de Raven Avancées - Abregées. *Pearson Clinical & Talent Assessment* (1998).
146. Corsi, P. M. Human memory and the medial temporal region of the brain. (Mc Gill University, Montréal, Canada, 1972).
147. Arce, T. & McMullen, K. The Corsi Block-Tapping Test: Evaluating methodological practices with an eye towards modern digital frameworks. *Computers in Human Behavior Reports* **4**, 100099 (2021).
148. Wechsler, D. Wechsler Adult Intelligence Scale—Fourth Edition (WAIS-4). *Pearson Clinical & Talent Assessment* (2008).
149. Ryan, J. J., Townsend, J. M. & Kreiner, D. S. Comparison of Oral, Written, and Pointing Responses to WAIS-IV Digit Span. *Appl Neuropsychol Adult* **21**, 94–97 (2014).
150. Conway, A. R. A., Kane, M. J. & Al, C. E. T. Working memory span tasks: A methodological review and user's guide. *Psychon Bull Rev* **12**, 769–786 (2005).
151. Bellon, E., van Bergen, E. & Dowker, A. D. Is Parental Mathematics Anxiety Associated with Young Children's Arithmetical Performance? *Educ Sci* **12**, 812 (2022).
152. Gordon, E. E. Music Aptitude and Related Tests An Introduction. *GIA Publications Inc* (1989).
153. Callejas, A., Lupiáñez, J., Funes, M. J. & Tudela, P. Modulations among the alerting, orienting and executive control networks. *Exp Brain Res* **167**, 27–37 (2005).
154. Deweer, B., Poitrenaud, J., Kalafat, M. & der Linden, M. CVLT - Test d'apprentissage et de mémoire verbale. *Pearson Clinical & Talent Assessment* (2008).
155. Strauss, E., Sherman, E. M. S. & Spreen, O. *A Compendium of Neuropsychological Tests: Administration, Norms, and Commentary, 3rd Ed.* (Oxford University Press, New York, NY, US, 2006).
156. Ashendorf, L., Horwitz, J. E. & Gavett, B. E. Abbreviating the Finger Tapping Test. *Archives of Clinical Neuropsychology* **30**, 99–104 (2015).
157. Tiffin, J. & Asher, E. J. The Purdue Pegboard: norms and studies of reliability and validity. *Journal of applied psychology* **32**, 234 (1948).
158. Frederickson, N., Frith, U. & Reason, R. *Phonological Assessment Battery (PhAB): Manual and Test Materials.* (NFER-Nelson, Windsor, 1997).
159. Rutten, S., Santoro, R., Hervais-Adelman, A., Formisano, E. & Golestani, N. Cortical encoding of speech enhances task-relevant acoustic information. *Nat Hum Behav* **3**, 974–987 (2019).
160. Gola-Asmussen, C., Lequette, C., Pouget G., Rouyer, C. & Zorman, M. ECLA-16+. Evaluation des compétences en lecture chez l'adulte de plus de 16 ans. *Université de Provence Aix-Marseille I-Cognisciences LSE Université Pierre Mendès, Grenoble* (2010).
161. Lefavrais, P. Test de l'Alouette: Test d'analyse de la lecture et de la dyslexie. *Paris: Editions du Centre de Psychologie Appliquée* (1967).

162. Sprenger-Charolles, L., Colé, P., Béchennec, D. & Kipffer-Piquard, A. French normative data on reading and related skills from EVALEC, a new computerized battery of tests. *European Review of Applied Psychology* **55**, 157–186 (2005).
163. Szenkovitz, G. & Ramus, F. Exploring dyslexics' phonological deficit I: Lexical vs sub-lexical and input vs output processes. *Dyslexia* **11**, 253–268 (2005).
164. Majerus, S., Van der Linden, M., Mulder, L., Meulemans, T. & Peters, F. Verbal short-term memory reflects the sublexical organization of the phonological language network: Evidence from an incidental phonotactic learning paradigm. *J Mem Lang* **51**, 297–306 (2004).
165. Fedeli, D., Del Maschio, N., Sulpizio, S., Rothman, J. & Abutalebi, J. The bilingual structural connectome: Dual-language experiential factors modulate distinct cerebral networks. *Brain Lang* **220**, 104978 (2021).

## Acknowledgements

This work was supported by the Swiss National Science Foundation [Grant #100014\_182381], and by the NCCR Evolving Language, Swiss National Science Foundation [Agreement #51NF40\_180888]. In addition to our funding sources, we would like to express our sincere gratitude to several collaborators who contributed their valuable resources and support throughout this large and intensive project. We gratefully acknowledge the following people: Roberto Martuzzi, Loan Mattera and Nathalie Philippe from the Human Neuroscience Platform at Campus Biotech for their invaluable and continued assistance with MRI setup and data collection; we thank our student assistants Vanessa Gottofrey, Melody Cascioli, and Aspasia Sfakaki for their contribution to data collection, and Rixa Gruhnert for helping with the LEAPQ annotations. We are thankful to Michael Dayan for his help with setting up the BIDS heuristic; to Priscilla Borges, Jutta Mueller and Hettie Roebuck for discussions on the Modes of Internal Reasoning Questionnaire use and its scoring; to Mark Eckert and Davide Fedeli for their helpful hints on technical validation of anatomical and diffusion-weighted MRI; to Gwendoline Mahé for providing us with the French version of the Adult Reading History Questionnaire; to Peter Schneider for consultations on the musical experience and musicality tests; to Evelina Fedorenko and her team for discussing with us the modifications to the Alice in Wonderland Localiser. We also thank Isabelle Udry for translating the Motivation Factor Questionnaire to French; Neiloufar Family and our colleague Sevil Maghsadagh for rating the Farsi task. Finally, we would like to acknowledge the collaborators who gave us access to and information on materials they have developed themselves: Tobias Kober and Siemens Healthineers© for the WIP-Compressed Sensing MPRAGE sequence, which crucially shortened the neuroimaging session without compromising on data quality; Will Barratt for sharing his Socioeconomic Status Measure; Franck Ramus and Steve Majerus for sharing their literacy and phonological awareness tasks; Charles Stansfield for granting permission to use the MLAT5; Elsjé Van Bergen for providing the Tempo Test Revised. Lastly, our heartfelt gratitude goes to two anonymous reviewers whose suggestions greatly improved this dataset, and to the 101 participants who took part in our study, for their commitment and generosity with their time.

## Author contributions

CRedit (Contributor Roles Taxonomy) is provided, as follows: Conceptualization: A.R., I.B., O.K., R.B., N.G.; Investigation: A.R., I.B.; Methodology and Formal Analysis: A.R., I.B., O.K.; Data Curation: A.R., I.B.; Technical Validation: A.R.; Writing - Original Draft A.R.; Writing - Review & Editing: A.R., I.B., O.K., R.B., N.G.; Visualization A.R.; Resources: R.B., N.G.; Supervision: R.B., N.G.; Funding acquisition: R.B., N.G.

## Competing interests

The authors declare no competing interests.

## Additional information

**Supplementary information** The online version contains supplementary material available at <https://doi.org/10.1038/s41597-024-04357-y>.

**Correspondence** and requests for materials should be addressed to A.R.

**Reprints and permissions information** is available at [www.nature.com/reprints](http://www.nature.com/reprints).

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2025