

# Reconstructing Voice Identity from Noninvasive Auditory Cortex

## Recordings

Charly Lamothe<sup>1,2</sup> ✉, Etienne Thoret<sup>1,2,3,4</sup>, Régis Trapeau<sup>1</sup>, Bruno L. Giordano<sup>1</sup>, Julien Sein<sup>1,5</sup>, Sylvain Takerkart<sup>1</sup>, Stéphane Ayache<sup>2</sup>, Thierry Artières<sup>2,6,7</sup> ✉, Pascal Belin<sup>1,7</sup> ✉

<sup>1</sup>La Timone Neuroscience Institute UMR 7289, CNRS, Aix-Marseille University, Marseille, France. <sup>2</sup>Laboratoire d'Informatique et Systèmes UMR 7020, CNRS, Aix-Marseille University, Marseille, France. <sup>3</sup>Perception, Representation, Image, Sound, Music UMR 7061, CNRS, Marseille, France. <sup>4</sup>Institute of Language Communication & the Brain, Marseille. <sup>5</sup>Centre IRM-INT@CERIMED, Marseille, France. <sup>6</sup>École Centrale de Marseille, Marseille, France. <sup>7</sup>These authors jointly supervised this work: Thierry Artières, Pascal Belin. ✉e-mail: [charlylmth@gmail.com](mailto:charlylmth@gmail.com), [thierry.artieres@lis-lab.fr](mailto:thierry.artieres@lis-lab.fr), [pascal.belin@univ-amu.fr](mailto:pascal.belin@univ-amu.fr)

## Abstract

The cerebral processing of voice information is known to engage Temporal Voice Areas (TVAs) that respond preferentially to conspecific vocalizations. But how voice information related to the stable physical characteristics of the speaker such as gender, age or identity is represented by neuronal populations in these areas remains poorly understood. Here we used a deep neural network (DNN) to generate a high-level, small-dimension representational space of voice stimuli—the ‘voice latent space’ (VLS)—and examined its linear relation with cerebral activity via encoding, representational similarity and decoding analyses. We find that the VLS maps onto fMRI measures of cerebral activity in response to tens of thousands of voice stimuli from hundreds of different speaker identities, and better accounts for the representational geometry for speaker identity in the TVAs than in A1. Moreover, the VLS allowed TVA-based reconstructions of voice stimuli that preserved important aspects of speaker gender and identity as assessed by both machine classifiers and human listeners. These results demonstrate that a low-dimensional, DNN-derived space accounts well for cerebral voice representations and provide insights into representational differences between A1 and the TVAs, paving the way to noninvasive brain-computer interface applications.

## Introduction

The human voice carries speech, but is also an “auditory face” that carries much valuable information on the stable physical characteristics of the speaker (hereafter, ‘identity-related’; Belin et al., 2004, 2011). The ability of listeners to extract identity-related information in voice such as gender, age, or unique identity even in brief stimuli plays a crucial role in our social interactions, yet its neural bases remain poorly understood compared to those of speech processing. Studies over the past two decades have clearly established via complementary neuroimaging techniques that the cerebral processing of voice information involves a set of temporal voice areas (TVAs) in secondary auditory cortical regions of the human (fMRI: Belin et al., 2000, von Kriegstein et al., 2004, Pernet et al., 2015; EEG, MEG: Charest et al., 2009, Capilla et al., 2013, Barbero et al., 2021; Electrophysiology: Rupp et al., 2022, Zhang et al., 2021) as well as macaque brain (Petkov et al., 2008; Bodin et al., 2021). The TVAs respond more strongly to sounds of voice – with or without speech (Pernet et al., 2015; Rupp et al., 2022; Trapeau et al., 2023)—and categorize voice apart from other sounds (Bodin et al., 2021) but the nature of the information encoded at these stages of cortical processing, especially with respect to speaker identity-related information, remains largely unknown (Blank et al., 2014; Belin et al., 2018).

In recent years, deep neural networks (DNNs) have emerged as a powerful tool for representing complex visual data, such as images (LeCun, Bengio, & Hinton, 2015) or videos (Liu et al., 2020). In the auditory domain, DNNs have been shown to provide valuable representations—so called feature or latent spaces—for modeling the cerebral processing of sound (brain encoding) (speech: Kell et al., 2018, Millet et al., 2022; semantic content Caucheteux et al., 2022, Caucheteux et King, 2022, Caucheteux et al., 2023, Giordano et al., 2023; music: Güçlü et al., 2016), or reconstructing the stimuli listened by a participant (brain decoding) (Akbari et al., 2019). They have not yet been used to explain cerebral representations of identity-related information, due in part to the focus on speech information (von Kriegstein 2003; Morillon et al., 2022).

Here, we addressed this challenge by training a ‘Variational autoencoder’ (VAE; Kingma et Welling, 2014) DNN to reconstruct voice spectrograms from 182,000 250-ms voice samples from 405 different speaker identities in 8 different languages from the CommonVoice database (Ardila et al., 2020). Brief (250ms) samples were used in order to emphasize speaker identity-related information in voice, already available after a few hundreds of milliseconds (Schweingenger et al., 1997; Lavan, 2023), over linguistic information unfolding over longer time periods. While a quarter of a second is admittedly short compared to standards of e.g. computational speaker identification that typically uses 2-3s samples, this short duration is sufficient to allow near perfect gender classification and performance levels well above chance for speaker discrimination (Fig. 5). This brief duration allowed presenting many more stimuli to our participants in the scanner, while preserving acceptable levels of behavioral and classifier performance.

State-of-the-art studies have largely relied on task-optimized neural networks (i.e., DNN trained using supervised learning to classify a category from the input) to study sensory cortex processes (Yamins et DiCarlo, 2016; Schrimpf et al., 2018). They can reach high accuracies in brain encoding (Khaligh-Razavi and Kriegeskorte, 2014; Schrimpf et al., 2018; Han et al., 2019), however there is increasing evidence that unsupervised learning such as used for the VAE also provides plausible computational models for investigating brain processing (Higgins et al., 2021; Zhuang et al, 2021; Millet et al., 2022; Orhan et al., 2022). Thus, the VAE-derived VLS, exploited within encoding, representational similarity and decoding frameworks, offers a potentially promising tool for investigating the representations of voice stimuli in the secondary auditory cortex (Naselaris et al., 2011). Autoencoders learn to compress stimuli with high dimensionality into a lower-dimensional space that nonetheless allows reconstruction of the original stimuli via an inverse transformation learned by the second part of the network called the decoder. Fig. 1a shows the architecture of the VAE, with its encoder that reduces an input spectrogram to a highly compressed, 128-dimension *voice latent space* (VLS) representation, and its decoder that reconstructs the spectrogram

from this VLS representation. Points in the VLS correspond to voice samples with different identities and phonetic content. A line segment in the VLS contains points corresponding to perceptual interpolations between its two extremities (Fig. 1b; Supplementary Audio 1). VLS coordinates of samples presented to the participants averaged by speaker identity suggests that a major organizational dimension of the latent space is along voice gender (Fig. 1b) (colored by age or language in Supplementary Figure 1).

In order to test whether VLS accounts well for cerebral activity in response to voice stimuli, we scanned three healthy volunteers using fMRI to measure an indirect index of their cerebral activity across 10+ hours of scanning each, in response to ~12,000 of the voice samples, denoted *BrainVoice* in the following, used to train the DNN. The small number of participants does not allow for generalization at the level of the general population as in standard fMRI studies, but allows testing for replicability as in comparable studies involving 10+ hours of scanning per participant (VanRullen & Reddy, 2019). Different stimulus sets were used across participants to provide a stringent test of replicability based on subject-level analyses. Stimuli consisted of randomly spliced 250-ms excerpts of speech samples from the CommonVoice database (Ardila et al., 2020) by 119 speakers in 8 different languages. For assessing generalization performances of decoding models and brain-based reconstruction, six different test stimuli were repeated more often (60 times) for each participant to provide robust estimates of their induced cerebral activity (see Methods). We first modeled these responses to voice using a general linear model (GLM) (Friston et al., 1994) with several nuisance regressors as an initial denoising step (Supplementary Figure 4), then used a second GLM modeling cerebral responses to the different speaker identities (Supplementary Figure 3a), resulting in one voxel activity map per speaker (Supplementary Figure 3b). We independently localized in each participant several regions of interest (ROIs) on which subsequent analyses were focused: the anterior, middle and posterior TVAs in each hemisphere (individually localized via an independent ‘voice localizer scan’ and MNI coordinates provided in Pernet et al., 2015; Supplementary Figure 3c) as well as primary

auditory cortex (A1) (using a probabilistic map in MNI space (Penhune et al., 1996; Supplementary Figure 3d).

We first asked how the VLS could account for the brain responses to speaker identities (encoding) measured in A1 and the TVAs, in comparison with a linear autoencoder's latent space (LIN). This approach was chosen because it has been demonstrated that a linear autoencoder with a  $d$ -dimensional hidden layer projects data in the same subspace as the one spanned by the  $d$  first eigenvectors of a principal component analysis (PCA) (Gallinari & LeCun et al., 1987; Plaut et al., 2018). For this, we used a general linear model (GLM) of fMRI responses to the speaker identities, resulting in one voxel activity map per speaker (Supplementary Figure 3). Then, we computed the average VLS coordinates of the fMRI voice stimuli for each speaker identity, which may be seen as a speaker representation in the VLS (see *Identity-based and stimulus-based representations* section). Next we trained a linear voxel-based encoding model to predict the speaker voxel activity maps from the speaker VLS coordinates. As VAE achieves compression through a series of nonlinear transformations (Wetzel, 2017), we choose to contrast its results with a linear autoencoder's latent space. This method has previously been applied to fMRI-based image reconstructions (Cowen et al., 2014; VanRullen & Reddy, 2019; Mozafari et al., 2020).

The extent to which the VLS allows linearly predicting the fMRI recordings does not provide insight into the representational geometries, i.e., the differences between the patterns of cerebral activity for speaker identity. We addressed this subsequent question by using representational similarity analysis (RSA; Kriegeskorte et al., 2008) in order to test which model better accounts for the representational geometry for voice identities in the auditory cortex. Using RSA as a model comparison framework has been shown relevant to examine the brain-model relationship from complementary angles (Diedrichsen et Kriegeskorte, 2017). We built speaker x speaker representational dissimilarity matrices (RDMs) capturing pairwise differences in cerebral activity or model predictions between all pairs of speakers;

then we examined how well the LIN and VLS-derived RDMs correlated with the cerebral RDMs from A1 and the TVAs.

A strong test of the adequacy of models of brain activity, and a long-standing goal in computational neurosciences, is the reconstruction of a stimulus presented to a participant from the evoked brain responses. While reconstruction of visual stimuli (images, videos) from cerebral activity has been performed by a number of groups (VanRullen et Reddy, 2019; Mozafari et al., 2020; Le et al., 2021; Gaziv et al., 2022; Chen et al., 2023), validating the DNN-derived representational spaces, comparable work in the auditory domain is scarce, almost exclusively concentrated on linguistic information (Santoro et al., 2017). Akbari et al. used a DNN to reconstruct speech stimuli based on ECoG recording of auditory cortex activity, an invasive method compared to techniques like fMRI. They obtained good phonetic recognition rate, but chance-level gender categorization performance from reconstructed spectrograms, and no evaluation of speaker identity discrimination.

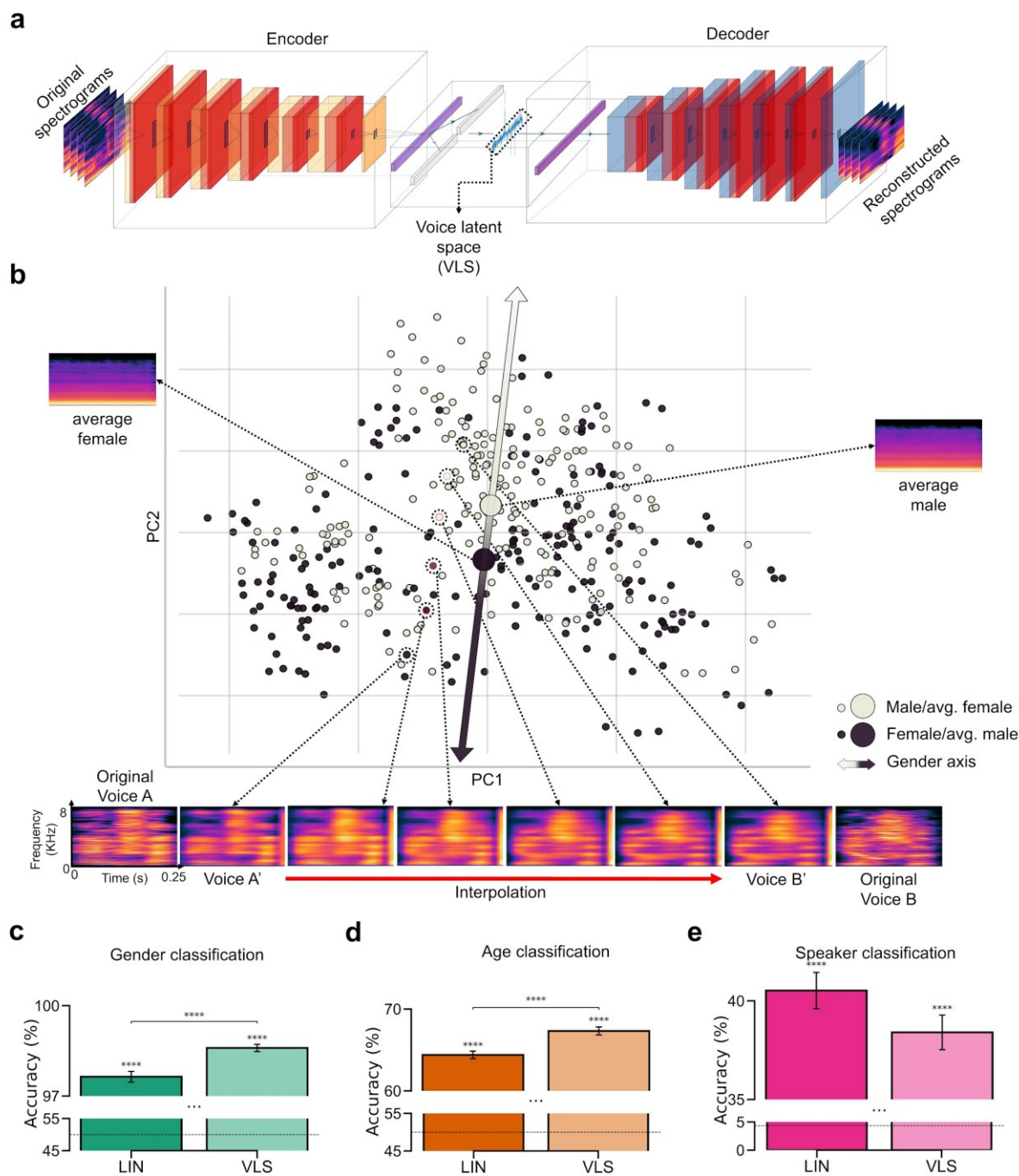
Here we built on the linear relationship uncovered in our encoding analysis between the VLS and the fMRI recordings to invert it and try and predict VLS coordinates from the recorded fMRI data; then, using the decoder, we reconstructed the spectrograms of stimuli presented to the participants (Wu et al., 2006; Naselaris et al., 2011). The voice identity information available in the reconstructed stimuli was finally assessed using both machine learning classifiers and behavioral tasks by human listeners (Fig. 4).

## Results

**Voice Information in the Voice Latent Space (VLS).** In order to probe the informational content of the VLS, linear classifiers were trained to categorize the voice stimuli from 405 speakers by gender (2 classes), age (2 classes) or identity (119 classes, cf Methods) based on VLS coordinates, or their LIN features as control (Fig. 1c,d,e; we aggregated the stimuli from the 3 participants; for each model computed the latent space of each stimulus and averaged the latent spaces by speaker identity, leading to 405 128-dimensional vectors. We

164 then trained linear classifiers using a 5-fold cross validation scheme, see *Characterization of*  
165 *the autoencoder latent space*). The mean of the distribution of accuracies obtained for 100  
166 random classifier initializations (as to account for variance; Bouthillier et al., 2018) was  
167 significantly above chance level (all  $p$ s <  $1e-10$ ) for all classifications (LIN: gender (mean  
168 accuracy  $\pm$  s.d.) =  $97.64 \pm 1.77\%$ ,  $t(99)=266.94$ ; age:  $64.39 \pm 4.54\%$ ,  $t(99)=31.53$ ; identity:  
169  $40.52 \pm 9.14\%$ ,  $t(99)=39.37$ ; VLS: gender:  $98.59 \pm 1.19\%$ ,  $t(99)=406.47$ ; age:  $67.31 \pm 4.86\%$ ,  
170  $t(99)=35.41$ ; identity:  $38.40 \pm 8.75\%$ ,  $t(99)=38.73$ ). We then evaluated the difference in  
171 performance at preserving identity-related information between the VLS and LIN via one-way  
172 ANOVAs. Results showed a significant effect of Feature (LIN/VLS) in categories (all  $F$ s(1,  
173 198) > 225.15, all  $p$ s < .0001) but not in identity. Post-hoc paired t-tests showed that the VLS  
174 was better than the LIN at encoding information related to voice identity, as evidenced by a  
175 significant difference in means for gender ( $t(99)=-6.11$ ,  $p<.0001$ ), age ( $t(99)=-6.10$ ,  $p<.0001$ )  
176 but not for identity classifications ( $t(99)=1.71$ ).





**Fig. 1 | DNN-derived Voice Latent Space (VLS). a, Variational autoencoder (VAE)**

**Architecture.** Two networks learned complementary tasks. An encoder was trained using

182K voice samples to compress their spectrogram into a 128-dimension representation, the

voice latent space (VLS) while a decoder learned the reverse mapping. The network was

trained end-to-end by minimizing the difference between the original and reconstructed

spectrograms. **b, Distribution of the 405 speaker identities along the first 2 principal**



components of the VLS coordinates from all sounds, averaged by speaker identity. Each disk represents a speaker identity colored by gender. PC2 largely maps onto voice gender (ANOVAs on the first two components: PC1:  $F(1, 405)=0.10$ ,  $p=.74$ ; PC2:  $F(1, 405)=11.00$ ,  $p<.001$ ). Large disks represent the average of all male (black) or female (gray) speaker coordinates, with their associated reconstructed spectrograms (note the flat fundamental frequency ( $f_0$ ) and formant frequencies contours caused by averaging). The bottom of spectrograms illustrate an interpolation between stimuli of two different speaker identities: spectrograms at the extremes correspond to two original stimuli (A, B) and their VLS-reconstructed spectrograms (A', B'). Intermediary spectrograms were reconstructed from linearly interpolated coordinates between those two points in the VLS (red line) (cf. Supplementary Audio 1). **c,d e**, Performance of linear classifiers at categorizing speaker gender (chance level: 50%), age (young/adult, chance level: 50%) or identity (119 identities, chance level: 0.84%) based on VLS or LIN coordinates. Error bars indicate standard error of the mean (s.e.m) across 100 random classifier initializations. All  $p<1e-10$ . The horizontal black dashed lines indicate chance levels. \*\*\*\*:  $p<0.0001$ .

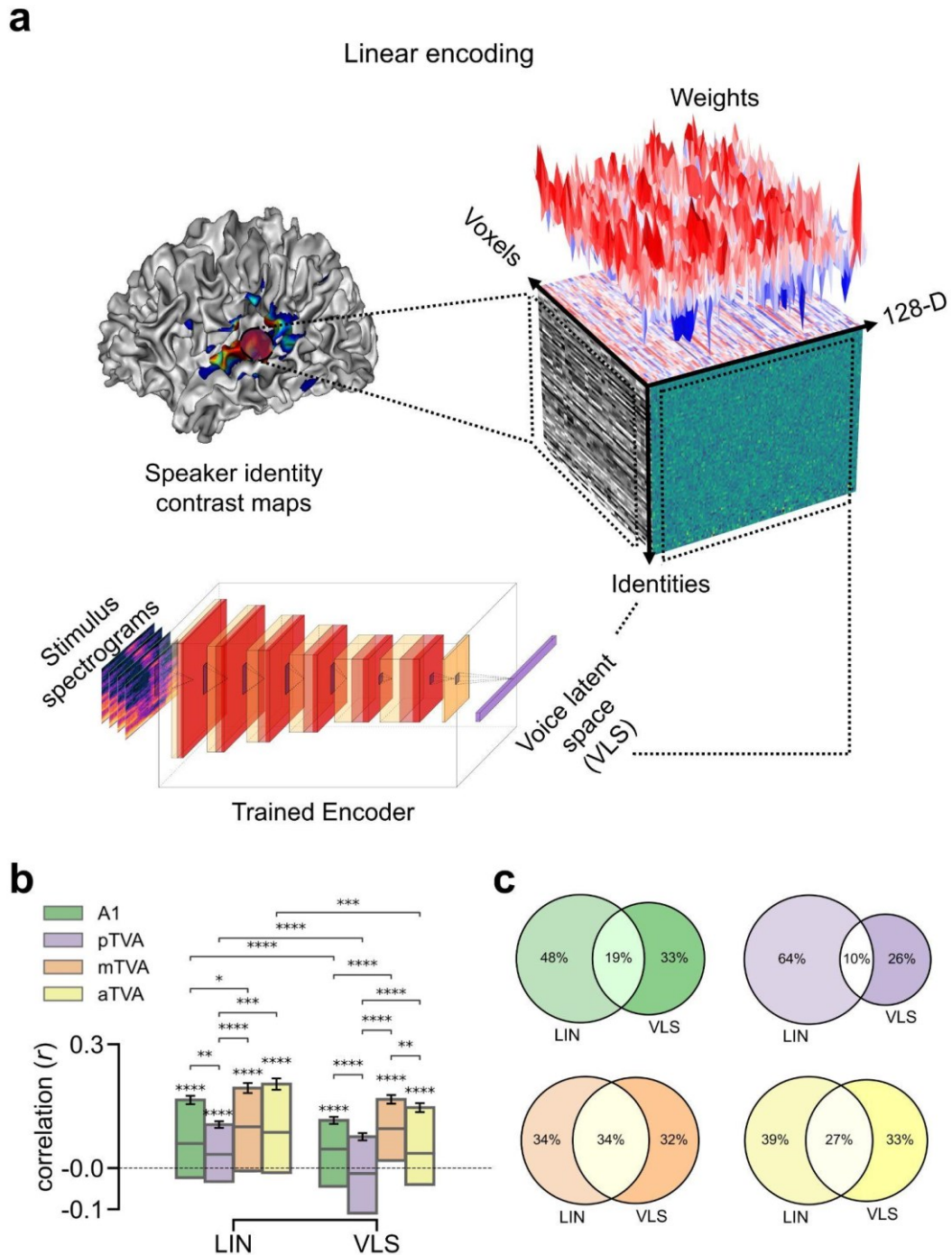
Thus, despite its low number of dimensions (each input spectrogram has  $401 \times 21 = 8421$  parameters and is summarized in the VLS by a mere 128 dimensions), the VLS appears to meaningfully represent the different sources of voice information perceptually available in the vocal stimuli. This representational space therefore constitutes a relevant candidate for linearly modeling voice stimulus representations by the brain.

**Brain Encoding** We used a linear voxel-based encoding model to test whether VLS linearly maps onto cerebral responses to speaker identities measured with fMRI in the different ROIs. A regularized linear regression model (cf. Methods) was trained on a subset of the data (5-fold cross validation scheme) to predict the voxel maps for each speaker identity. For each fold, the trained model was tested on the held-out speaker identities (Fig. 2a). For each ROI, the performance of the model was assessed using the Pearson correlation score between the true and the predicted responses of each voxel (Schrimpf et al., 2021). Similar

predictions were tested with features derived from LIN (cf. Methods). Fig. 2b shows, for each of the ROIs, the distribution of correlation coefficients obtained for the 2 sets of features across voxels, hemispheres and participants.

One-sample t-tests showed that the means of Fisher z-transformed coefficients for both LIN features and VLS were significantly higher than zero (LIN: A1  $t(197)=7.25$ ,  $p<.0001$ , pTVA  $t(175)=4.49$ ,  $p<.0001$ , mTVA  $t(164)=9.12$ ,  $p<.0001$  and aTVA  $t(147)=6.81$ ,  $p<.0001$ ; VLS: A1  $t(197)=4.76$ ,  $p<.0001$ , mTVA  $t(164)=10.12$ ,  $p<.0001$  and aTVA  $t(147)=5.52$ ,  $p<.0001$  but not pTVA  $t(175)=-1.60$ ) (Supplementary Tables 2-3).

A mixed ANOVA performed on the Fisher z-transformed coefficients with Feature (VLS, LIN) and ROI (A1, pTVA, mTVA, aTVA) as factors showed a significant effect of Feature ( $F(3, 683)=56.65$ ,  $p<.0001$ ), a significant effect of ROI ( $F(3, 683)=18.50$ ,  $p<.0001$ ), and a moderate interaction Feature x ROI ( $F(3, 683)=5.25$ ,  $p<.01$ ). Post-hoc comparisons revealed that the mean of correlation coefficients was higher for LIN than for VLS in A1 ( $t(197)=4.02$ ,  $p<.0001$ ), pTVA ( $t(175)=6.64$ ,  $p<.0001$ ), aTVA ( $t(147)=3.78$ ,  $p<.001$ ) but not in mTVA ( $t(164)=0.58$ ) (Supplementary Table 4); and that the voxel patterns are better predicted in mTVA than in A1 for both models (LIN:  $t(361)=2.36$ ,  $p<.05$ ; VLS:  $t(361)=4.91$ ,  $p<.0001$ ) (Supplementary Table 5). However, we found by inspecting the distribution of model-voxel correlations that both models account for different parts of the voice identities responses, and differently across ROIs (Fig. 2c).



**Fig. 2 | Predicting brain activity from the VLS. a, Linear prediction of brain activity from VLS for ~135 speaker identities in the different ROIs. We first fit a GLM to predict the BOLD responses to each voice speaker identity. Then, using the trained encoder, we computed the average VLS coordinates of the voice stimuli presented to the participants based on speaker identity. Finally, we trained a linear voxel-based encoding model to predict the speaker voxel activity maps from the speaker VLS coordinates. The cube illustrates the linear relationship**

between the fMRI responses to speaker identity and the VLS coordinates. The left face of the cube represents the activity of the voxels for each speaker identity, with each line corresponding to one speaker. The right face displays the VLS coordinates for each speaker identity. The top face of the cube shows the weight vectors of the encoding model. **b**, Encoding results. For each region of interest, the performance of the model was assessed using the Pearson correlation score between the true and the predicted responses of each voxel on the held-out speaker identities. Pearson's correlation coefficients were computed for each voxel on the speakers' axis, then averaged across hemispheres and participants. Similar predictions were tested with the LIN features. Error bars indicate standard error of the mean (s.e.m) across voxels. \* $p < 0.05$ ; \*\* $p < 0.01$ ; \*\*\* $p < 0.001$ ; \*\*\*\* $p < 0.0001$ . **c**, Venn diagrams of the number of voxels in each ROI with the LIN, the VLS or both models. For each ROI and each voxel, we checked whether the test correlation was higher than the median of all participant correlations (intersection circle), and if not which model (LIN or VLS) yielded the highest correlation (left or right circles).

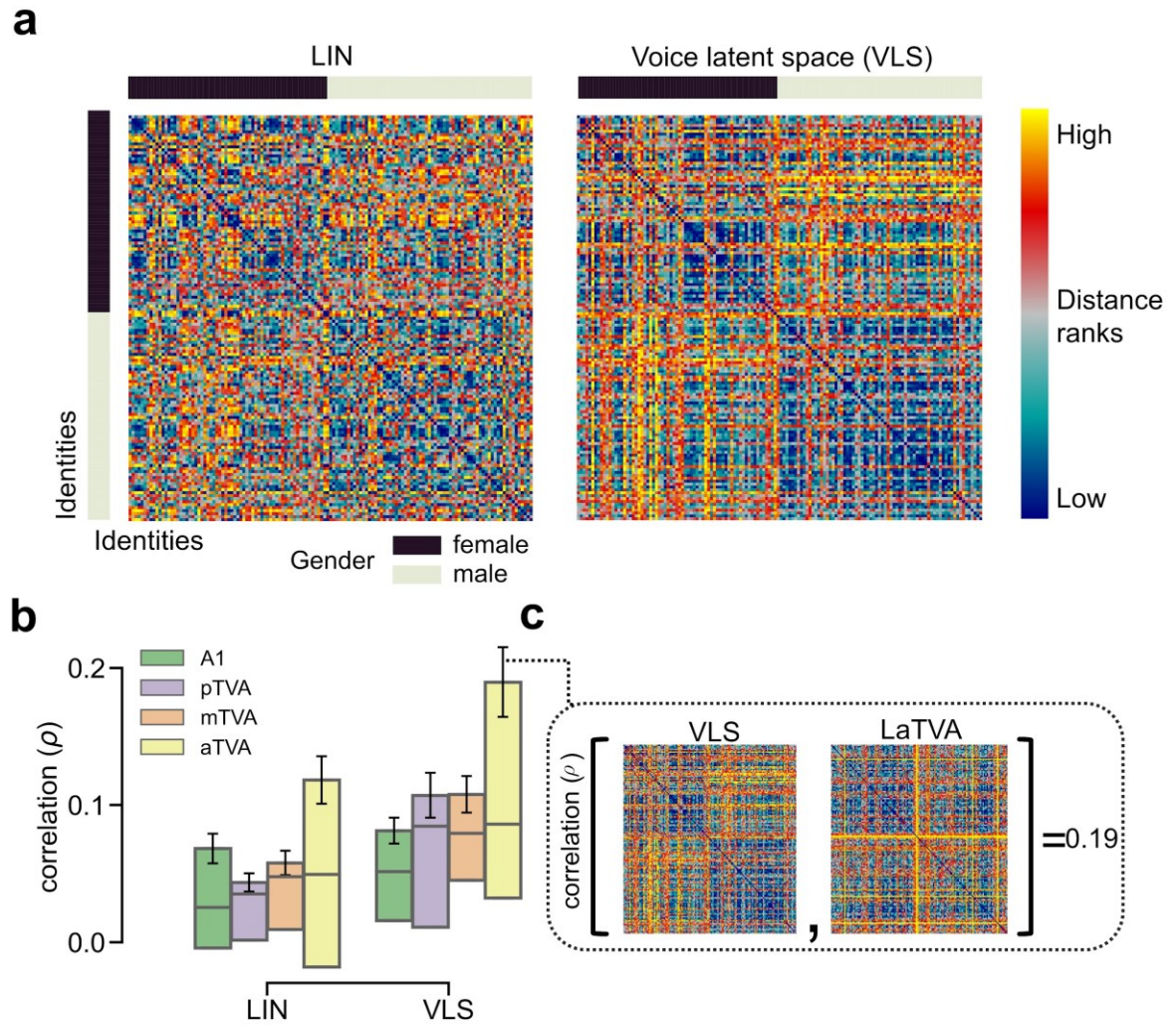
**Representational Similarity Analysis** For RSA, we built speaker x speaker representational dissimilarity matrices (RDMs) capturing for each ROI the dissimilarity in voxel space between each pair of speaker voxel maps ('brain RDMs'; cf. Methods) using Pearson's correlation (Walther et al., 2016). We compared these four bilateral brain RDMs (A1, aTVA, mTVA, pTVA) to two 'model RDMs' capturing speaker pairwise feature differences predicted by LIN and the VLS (Fig. 3a) built using cosine distance (Xing et al., 2015; Bhattacharya et al., 2017; Wang et al., 2018). Fig. 3b shows for each ROI the Spearman correlation coefficients between the brain RDMs and the two model RDMs, for each participant and hemisphere (Kriegeskorte et al., 2008; Fig. 3c for an example of brain-model correlation).

These brain-model correlation coefficients were compared to zero using a 'maximum statistics' approach based on random permutations of the model RDMs' rows and columns, (Maris & Oostenveld, 2007; cf. Methods; Fig. 3b). For the LIN model, only one brain-model

264 RDM correlation was significantly different from zero (one-tailed test): in mTVA, right  
 265 hemisphere in S3 ( $p=.0500$ ). For the VLS model, in contrast, 5 significant brain-model RDM  
 266 correlations were observed in all four ROIs: in A1, right hemisphere in S3 ( $p=.0142$ ); pTVA:  
 267 right hemisphere in S3 ( $p=.0160$ ); mTVA: left hemisphere in S3 ( $p=.007$ ); aTVA: left  
 268 hemispheres in S1 ( $p=.0417$ ) and S3 ( $p=.0001$ ) (Supplementary Table 6).

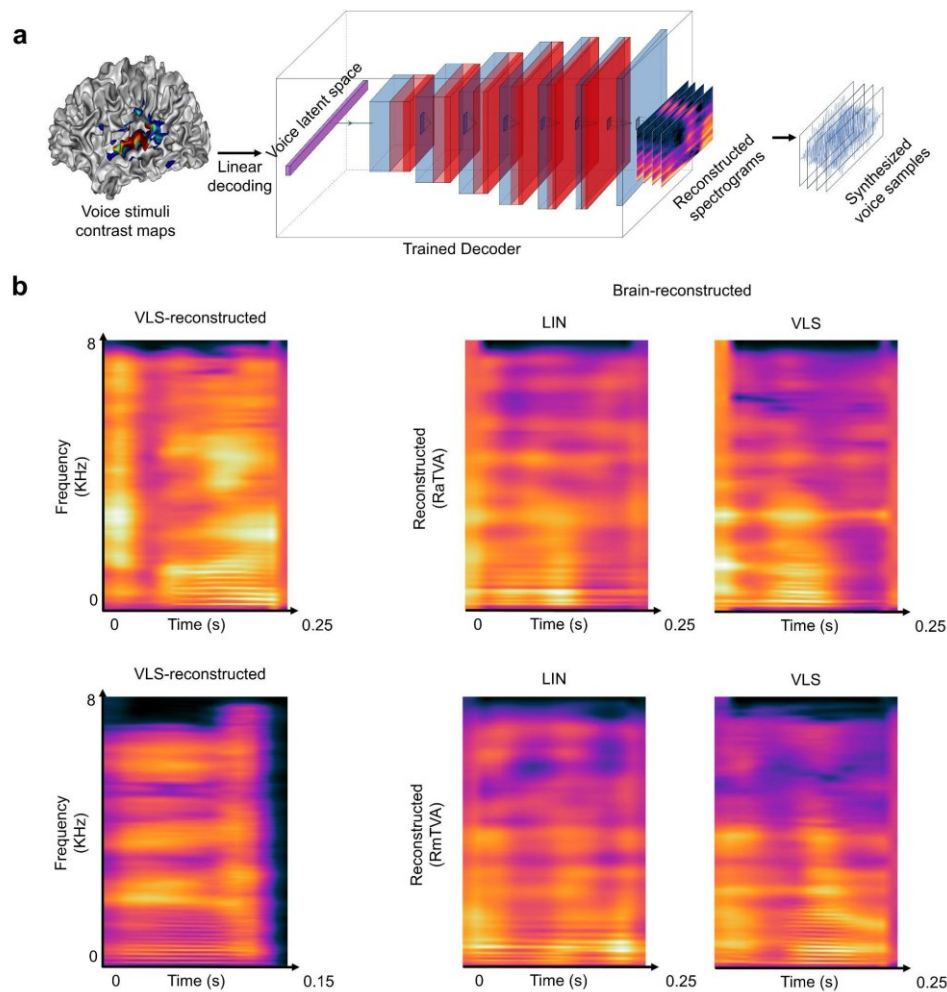
269 A two-way repeated-measures ANOVA with Feature (VLS, LIN) and ROI (A1, pTVA, mTVA,  
 270 aTVA) as factors performed on the Fisher z-transformed correlation coefficients showed a  
 271 tendency towards a significant effect of Feature ( $F(1, 2)=22.53$ ,  $p=.04$ ), and no ROI ( $F(3,$   
 272  $6)=1.79$ ,  $p=.30$ ) or interaction effects ( $F(3, 6)=1.94$ ,  $p=.22$ ). We compared the correlation  
 273 coefficients between the VLS and LIN models within participants and hemispheres using  
 274 one-tailed tests, based on the a priori hypothesis that the VLS models would exhibit greater  
 275 brain-model correlations than the LIN models (cf. Methods). The results revealed two  
 276 significant differences in one of the three participants, both in favor of the VLS model (S3:  
 277 right pTVA,  $p=.0366$ ; left aTVA,  $p=.00175$ ) (Supplementary Table 7).





**Fig. 3 | The VLS better explains representational geometry for voice identities in the TVAs than the linear model.** **a**, Representational dissimilarity matrices (RDMs) of pairwise speaker dissimilarities for ~135 identities (arranged by gender, cf. side bars), according to LIN and VLS. **b**, Spearman correlation coefficients between the brain RDMs for A1 and the 3 TVAs, and the 2 model RDMs. Error bars indicate standard error of the mean (s.e.m) across brain-model correlations. **c**, Example of brain-model RDM correlation in the TVAs. The VLS RDM and the brain RDM yielding one of the highest correlations (LaTVA) are shown in insert.

**Decoding and Reconstruction** We finally inverted the brain-VLS relationship to predict linearly VLS coordinates based on fMRI measurements (Fig. 4a; see ‘Brain decoding’ in Methods), and reconstruct via the trained decoder the spectrograms of 18 Test Stimuli (3 participants x 6 stimuli per participant; see Fig. 4b, and Supplementary Audio 2; audio estimated from spectrogram through phase reconstruction).



**Fig. 4 | Reconstructing voice identity from brain recordings.** *a*, A linear voxel-based decoding model was used to predict the VLS coordinates of 18 Test Stimuli based on fMRI responses to ~12 000 Train stimuli in the different ROIs. To reconstruct the audio stimuli from the brain recordings, the predicted VLS coordinates were then fed to the trained decoder to yield reconstructed spectrograms, synthesized into sound waveforms using the Griffin-Lim phase reconstruction algorithm (Griffin & Lim, 1983). *b*, Reconstructed spectrograms of the stimuli presented to the participants. Left panels show the spectrogram



*of example original stimuli reconstructed from the VLS, and the right panels brain-reconstructed spectrograms via LIN and the VLS (cf Supplementary Audio 2).*

We first assessed the nature of the reconstructed stimuli by using a DNN trained to categorize natural audio events (Howard et al., 2017): all reconstructed versions of the 18 Test Stimuli were categorized as 'speech' (1 class out of 521 - no 'voice' classes). To evaluate the preservation of voice identity information in the reconstructed voices, pre-trained linear classifiers were used to classify the speaker gender (2 classes), age (2 classes), and identity (17 classes) of the 18 reconstructed Test Stimuli. The mean of the accuracy distribution obtained across random classifier initializations (20 per ROI) used on the stimuli reconstructed from the induced brain activity was significantly above chance level for gender (LIN: pTVA (mean accuracy  $\pm$  s.d.):  $72.08 \pm 5.48$ ,  $t(39)=25.15$ ; VLS: A1:  $61.11 \pm 2.15$ ,  $t(39)=32.25$ ; pTVA:  $63.89 \pm 2.78$ ,  $t(39)=31.22$ ), age (LIN: pTVA:  $54.58 \pm 4.14$ ,  $t(39)=6.90$ ; aTVA:  $63.96 \pm 12.55$ ,  $t(39)=6.94$ ; VLS: pTVA:  $65.00 \pm 7.26$ ,  $t(39)=12.89$ ; aTVA:  $60.42 \pm 5.19$ ,  $t(39)=12.54$ ) and identity (LIN: A1:  $9.20 \pm 9.23$ ,  $t(39)=2.24$ ; pTVA:  $9.48 \pm 4.90$ ,  $t(39)=4.59$ ; aTVA:  $9.41 \pm 6.28$ ,  $t(39)=3.51$ ; VLS: pTVA:  $16.18 \pm 7.05$ ,  $t(39)=9.11$ ; aTVA:  $8.23 \pm 4.70$ ,  $t(39)=3.12$ ) (Fig. 5a-c; Supplementary Tables 8-10).

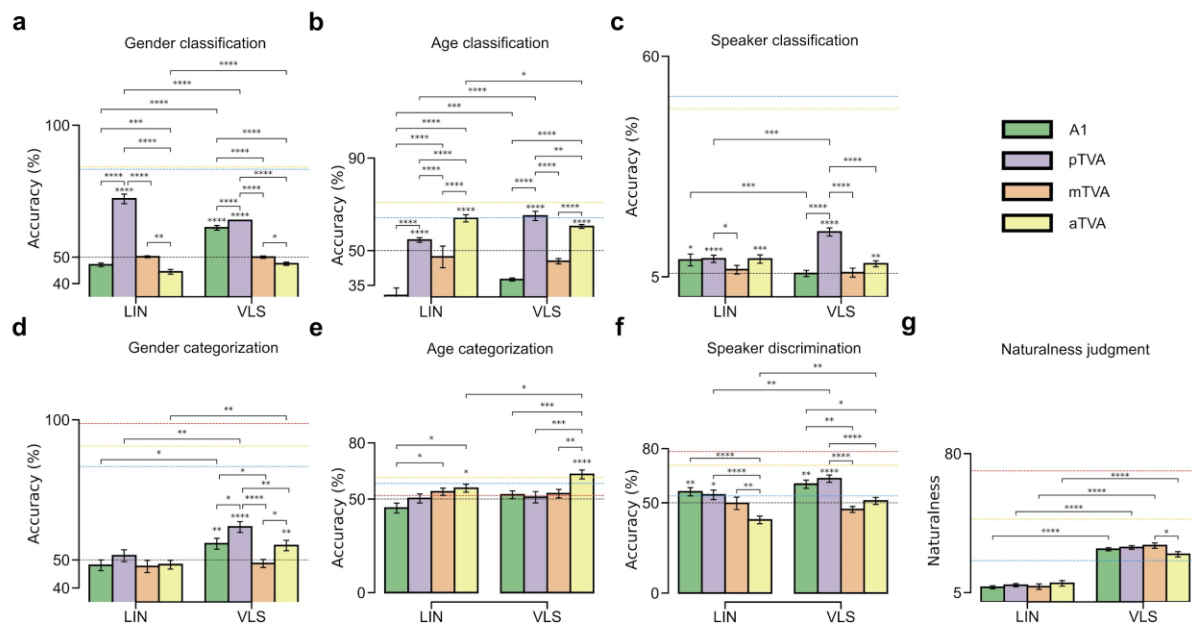
Two-way ANOVAs with Feature (VLS, LIN) and ROI (A1, pTVA, mTVA, aTVA) as factors performed on classification accuracy scores (gender, age, identity) revealed for gender classifications significant effects of Feature  $F(1, 312)=12.82$ ,  $p<.0005$  and of ROI (gender:  $F(3, 312)=245.06$ ,  $p<.0001$ ; age:  $F(3, 312)=64.49$ ,  $p<.0001$ ; identity:  $F(3, 312)=14.49$ ,  $p<.0001$ ), as well as Feature x ROI interactions (gender:  $F(3, 312)=56.74$ ,  $p<.0001$ ; age:  $F(3, 312)=4.31$ ,  $p<.001$ ; identity:  $F(3, 312)=8.82$ ,  $p<.0001$ ). Post-hoc paired t-tests indicated that the VLS was better than LIN in preserving gender, age and identity information in at least one TVA compared with A1 (gender: aTVA:  $t(39)=5.13$ ,  $p<.0001$ ; age: pTVA:  $t(39)=9.78$ ,  $p<.0001$ ; identity: pTVA:  $t(39)=4.01$ ,  $p<.0005$ ) (all tests in Supplementary Table 11). Post-hoc two sample t-tests comparing ROIs revealed significant differences in all classifications, in particular with pTVA outperforming other ROIs in gender (LIN: pTVA vs A1:

t(78)=22.40,  $p<.0001$ ; pTVA vs mTVA: t(78)=10.92,  $p<.0001$ ; pTVA vs aTVA: t(78)=31.47,  $p<.0001$ ; VLS: pTVA vs A1: t(78)=4.94,  $p<.0001$ ; pTVA vs mTVA: t(78)=13.96,  $p<.0001$ ; pTVA vs aTVA: t(78)=22.06,  $p<.0001$ ), age (LIN: pTVA vs A1: t(78)=7.26,  $p<.0001$ ; pTVA vs mTVA: t(78)=10.11,  $p<.0001$ ; VLS: pTVA vs A1: t(78)=5.71,  $p<.0001$ ; pTVA vs mTVA: t(78)=10.11,  $p<.0001$ ; pTVA vs aTVA: t(78)=3.21,  $p<.005$ ) and identity (LIN: pTVA vs mTVA: t(78)=2.27,  $p<.05$ ; VLS: pTVA vs A1: t(78)=6.45,  $p<.0001$ ; pTVA vs mTVA: t(78)=6.62,  $p<.0001$ ; pTVA vs aTVA: t(78)=5.85,  $p<.0001$ ) (Supplementary Table 12).

We further evaluated voice identity information in the reconstructed stimuli by testing human participants ( $n=13$ ) in a series of 4 online experiments assessing the reconstructed stimuli on: (i) naturalness judgment; (ii) gender categorization; (iii) age categorization; and (iv) speaker categorization (cf Methods). The naturalness rating task showed that the VLS-reconstructed stimuli sounded more natural compared to LIN-reconstructed ones, as revealed by a two-way repeated-measures ANOVA (factors: Feature and ROI) with a strong effect of Feature ( $F(1, 12)=53.72$ ,  $p<.0001$ ) and a small ROI x Feature interaction ( $F(3, 36)=5.36$ ,  $p<.005$ ). Post-hoc paired t-tests confirmed the greater naturalness of VLS-reconstructed stimuli in both A1 and the TVAs (all  $ps<.0001$ ) (Fig. 5g).

For the gender task, one-sample t-tests showed that categorization of the reconstructed stimuli was only significantly above chance level for the VLS (A1: (mean accuracy  $\pm$  s.d.)  $55.77\pm 10.84$ ,  $t(25)=2.66$ ,  $p<.01$ ; pTVA:  $61.75\pm 7.11$ ,  $t(25)=8.26$ ,  $p<.0001$ ; aTVA:  $55.13\pm 9.23$ ,  $t(25)=2.78$ ,  $p<.01$ ). Regarding the age and speaker categorizations, results also indicated that both the LIN- and VLS-reconstructed stimuli yielded above-chance performance in the TVAs (age: LIN: aTVA,  $55.77\pm 14.95$ ,  $t(25)=1.93$ ,  $p<.05$ ; VLS: aTVA,  $63.14\pm 11.82$ ,  $t(25)=5.56$ ,  $p<.0001$ ; identity: LIN: pTVA:  $54.38\pm 9.34$ ,  $t(17)=1.93$ ,  $p<.05$ ; VLS: pTVA:  $63.33\pm 6.75$ ,  $t(17)=8.14$ ,  $p<.0001$ ) (Supplementary Tables 13-15). Two-way repeated-measures ANOVAs revealed a significant effect of ROI for all categories (gender:  $F(3, 27)=5.90$ ,  $p<.05$ ; age:  $F(3, 36)=14.25$ ,  $p<.0001$ ; identity:  $F(3, 24)=38.85$ ,  $p<.0001$ ), and a Feature effect for gender ( $F(1, 9)=43.61$ ,  $p<.0001$ ) and identity ( $F(1, 8)=14.07$ ,  $p<.001$ ), but not for age ( $F(1, 12)=4.01$ ,  $p=0.07$ ), as well as a ROI x Feature interaction for identity

discrimination ( $F(3, 24)=3.52$ ,  $p<.05$ ) (Supplementary Tables 16-17 for the model and ROI comparisons).



**Fig. 5 | Behavioural and machine classification of the reconstructed stimuli. a,b,c,** Decoding of voice identity information in brain-reconstructed spectrograms. Performance of linear classifiers at categorizing speaker gender (chance level: 50%), age (chance level: 50%), and identity (17 identities, chance level: 5.88%). Error bars indicate s.e.m across 40 random classifier initializations per ROI (instance of classifiers; 2 hemispheres x 20 seeds). The horizontal black dashed line indicates chance level. The blue and yellow dashed lines indicate ceiling levels for the LIN and the VLS respectively.  $*p < .05$ ;  $**p < .001$ ,  $***p < .001$ ;  $****p < .0001$ . **d,e,f,** Listener performance at categorizing speaker gender (chance level: 50%) and age (chance level: 50%), and at identity discrimination (chance level: 50%) in the brain-reconstructed stimuli. Error bars indicate s.e.m across participant scores. The horizontal black dashed line indicates chance level, while the red, blue and yellow dashed lines indicate the ceiling levels for the original stimuli, the LIN-reconstructed and the VLS-reconstructed, respectively.  $*p < .05$ ;  $**p < .01$ ;  $***p < .001$ ,  $****p < .0001$ . **g,** Perceptual ratings of voice naturalness in the brain-reconstructed stimuli' as assessed by human listeners.  $*p < .05$ ,  $****p < .0001$ .

## Discussion

In this study we examined to what extent the cerebral activity elicited by brief voice stimuli can be explained by machine-learned representational spaces, with a specific focus on identity-related information. We trained a linear model and a DNN model to reconstruct 100,000s of short voice samples from 100+ speakers, providing low-dimensional spaces (LIN and VLS) which we related to fMRI measures of cerebral response to thousands of these stimuli. We find: (i) that 128 dimensions are sufficient to explain a sizeable portion of the brain activity elicited by the voice samples and yield brain-based voice reconstructions that preserve identity-related information; (ii) that the DNN-derived VLS outperforms the LIN space particularly in yielding more brain-like representational spaces and more naturalistic voice reconstructions; (iii) that different ROIs have different degrees of brain-model relationship, with marked differences between A1 and the the a, m, and pTVAs.

Low-dimensional spaces generated by machine learning have been used to approximate cerebral face representations and reconstruct recognizable faces based on fMRI (VanRullen et Reddy, 2019; Dado et al, 2022). In the auditory domain, however, they have mostly been used with a focus on linguistic (speech) information, ignoring identity-related information (but see Akbari et al., 2019). Here we applied them to brief voice stimuli—with minimal linguistic content but already rich identity-related information—and found that as little as 128 dimensions account reasonably well for the complexity of cerebral responses to thousands of these voice samples as measured by fMRI (Fig. 2). LIN and VLS both showed brain-like representational geometries, particularly the VLS in the aTVAs (Fig. 3). They made possible what is to our knowledge the first fMRI-based voice reconstructions to preserve voice-related identity information such as gender, age or even individual identity, as indicated by above-chance categorization or discrimination performance by both machine classifiers (Fig. 5a-c) and human listeners (Fig. 5d-f). Note that LIN and VLS also represent the limited linguistic content of the brief stimuli, as indicated by high language classification performance (Supplementary Figure 4).

Estimation of fMRI responses (encoding) by LIN yielded correlations largely comparable to those by VLS (Fig. 2b) although many voxels were only explained by one or the other space (Fig. 2c). But in the RSA, VLS yielded higher overall correlations with brain RDMs (Fig. 3), suggesting a representational geometry closer to that instantiated in the brain than LIN. Further, VLS-reconstructed stimuli sounded more natural than the LIN-reconstructed ones (Fig. 5g) and yielded both the best speaker discrimination by listeners (Fig. 5f) and speaker classification by machine classifiers (Fig. 5c). Unlike LIN, which was generated via linear transforms, VLS was obtained through a series of nonlinear transformations (Wetzel, 2017). The fact that the VLS outperforms LIN in terms of decoding performance is an indication that nonlinear transformation is required to better account for brain representation of voices (Naselaris et al., 2011; Cowen et al., 2014; Han et al., 2019).

Comparisons between ROIs revealed important differences between A1 and the a, m and pTVAs. For both LIN and VLS, predictions of fMRI signal (encoding) were more accurate for the mTVAs than for A1, and for A1 than for the pTVAs (Fig. 2b). The aTVAs yielded the highest correlations with the models in the RSA (Fig. 3). Stimulus reconstructions (Fig. 4) based on the TVAs also yielded better gender, age and identity classification than those based on A1, with gender and identity best preserved in the pTVA-, and to a lesser extent, in the aTVA-based reconstructions (Fig. 5). These results show that the a and pTVAs not only respond more strongly to vocal sounds than A1, they also better represent identity-related information in voice better than mTVA, which was previously anticipated in some neuroimaging studies (Latinus et al., 2011; Charest et al., 2013; Aglieri et al., 2021).

Overall, we show that a DNN-derived representational space provides an interesting approximation of the cerebral representations of brief voice stimuli that can preserve identity-related information. We find remarkable that such results could be obtained to explain sound representations despite the poor temporal resolution of fMRI. Future work combining more complex architectures to time-resolved measures of cerebral activity such as magneto-

encephalography (Défossez et al., 2023) or ECoG (Pasley et al., 2012) will likely yield better models of the cerebral representations of voice information.

## Methods

### Experimental procedure overview

Three participants attended 13 MRI sessions each. The first session was dedicated to acquire high-resolution structural data, as well as to identify the voice-selective areas of each participant using a ‘voice localizer’ based on different stimuli than those in the same experiment (Pernet et al., 2015; see below).

The next 12 sessions began with the acquisition of two fast structural scans for inter-session realignment purposes, followed by six functional runs, during which the main stimulus set of the experiment was presented. Each functional run lasted approximately 12 minutes.

Participants 1 and 2 attended all scanning sessions (72 functional runs in total); due to technical issues, Participant 3 only performed 24 runs.

Participants were instructed to stay in the scanner while listening to the stimuli. To maintain participants’ awareness during functional scanning, they were asked to press an MRI-compatible button each time they heard the same stimulus two times in a row, a rare event occurring 3% of the time (correct button hits (median accuracy  $\pm$  s.d.): S1=96.67 $\pm$ 7.10, S2=100.00 $\pm$ 0.89, S3=95.00 $\pm$ 3.68).

Scanning sessions were spaced by at least two days to avoid possible auditory fatigue due to the exposure to scanner noise. To ensure that participants’ hearing abilities did not vary across scanning sessions, hearing thresholds were measured before each session using a standard audiometric procedure (Martin and Champlin, 2000; ISO 2004) and compared with the thresholds obtained prior the first session.

## Participants

This study was part of the project 'Réseaux du Langage' and was promoted by the National Center for Scientific Research (CNRS). It has been given favorable approval by the local ethics committee (Comité de Protection des Personnes Sud-Méditerranée) on the date of 13th February 2019. The National Agency for Medicines (ANSM) has been informed of this study, which is registered under the number 2017-A03614-49. Three native French human speakers were scanned (all females; 26-33 years old). Participants gave written informed consent and received a compensation of 40€ per hour for their participation. All were right-handed and no one had hearing disorder or neurological disease. All participants had normal hearing thresholds of 15 dB HL, for octave frequencies between 0.125 and 8 kHz.

## Stimuli

The auditory stimuli were divided into two sequences. One 'voice localizer' sequence to identify the voice-selective areas of each participant (Pernet et al., 2015) and a main voice stimuli.

*Voice localizer stimuli.* The voice localizer stimuli consisted of 96 complex sounds of 500ms grouped in four categories of human voice, macaque vocalizations, marmoset vocalizations, and complex non-vocal sounds (more details in Bodin et al., 2021).

*Main voice stimuli.* The main stimulus set consisted of brief human voice sounds sampled from the Common Voice dataset (Ardila et al., 2020). Stimuli were organized into four main category levels: language (English, French, Spanish, Deutch, Polish, Portuguese, Russian, Chinese), gender (female/male), age (young/adult; young: teenagers and twenties; adult: thirties to sixties included) and identity (S1: 135 identities; S2: 142 identities; S3: 128 identities; ~44 samples per identity). Throughout the manuscript, the term 'gender' rather than 'sex' was utilized in reference to the demographic information obtained from the participants of the Common Voice dataset (Ardila et al., 2020), as it was the terminology employed in the survey ('male/female/other'). Stimulus sets were different for each



participant and the number of stimuli per set also varied slightly (number of unique stimuli: Participant 1, N=6150; Participant 2, N=6148; Participant 3, N=5123). For each participant, six stimuli were selected randomly among the sounds having a high energy (as measured with the amplitude envelope) from their stimulus set and were repeated extensively (60 times), to improve the performance of the brain decoding (VanRullen et Reddy, 2019; Horikawa & Kamitani, 2017; Chang et al., 2019); these will be called the “repeated” stimuli hereafter, the remaining stimuli were presented twice. The third participant attended 5 BrainVoice sessions instead of 12, one BrainVoice session corresponding to 1030 stimuli (1024 unique stimuli and 6 ‘test’ stimuli). Specifically, 5270 stimuli were presented to the third participant instead of ~12,000 for the two others. Among these 5270 stimuli, 5120 unique stimuli were presented once, as for the two other participants, 6 ‘test’ stimuli were presented 25 times (150 trials). The stimuli were balanced within each run according to language, gender, age, and identity, as to avoid any potential adaptation effect. In addition, identity was balanced across sessions.

All stimuli of the main set were resampled at 24414 Hz and adjusted in duration (250 ms). For each stimulus, a fade-in and a fade-out were applied with a 15 ms cosine ramp to their onset and offset, and were normalized by dividing the root mean square amplitude. During fMRI sessions, stimulus presentations were controlled using custom Matlab scripts (Mathworks, Natick, MA, USA) interfaced with an RM1 Mobile Processor (Tucker-David Technologies, Alachua, USA). The auditory stimuli were delivered pseudo-randomly through MRI-compatible earphones (S14, SensiMetrics, USA) at a comfortable sound pressure level that allowed for clear and intelligible listening.

## **Computational models**

We used two computational models to learn representational space for voice signals, Linear Autoencoder (LIN) and Deep Variational Autoencoder (VAE; Kingma et Welling., 2014). Both are encoder-decoder models that are learnt to reproduce at their output their input while going through a low dimensional representation space usually called latent space (that we

will call *voice latent space* since they are learnt on voice data). The autoencoders were trained on a dataset of 182K sounds from the Common Voice dataset (Ardila et al., 2020), balanced in gender, language and identity to reduce the bias in the synthesis (Gutierrez et al., 2021). Both models operate on sounds which were represented as spectrograms that we describe below. These representations were tested in all the encoding/decoding and RSA analyses.

## **Spectrograms**

We used amplitude spectrograms as input of the models that we describe below. Short term Fourier transforms of the waveform were computed using a sliding window of length 50 ms with a hop size of 12.5 ms (hence an overlap of 37.5 ms) and applying a Hamming window of size 800 samples before computing the Fourier transform of each slice. Only the magnitude of the spectrogram was kept and the phase of the complex representation was removed. At the end, a 250 ms sound is represented by a 21×401 matrix with 21 time steps and 401 frequency bins.

We used a custom code based on *numpy.fft* package (Harris et al., 2020). The size and the overlap between the sliding windows of the spectrogram were chosen to conform with the uncertainty principle between time and frequency resolution. The main constraint was to find a trade-off between accurate phase reconstruction with the Griffin & Lim algorithm (1983) and a reasonable size of the spectrogram.

We standardized each of the 401 frequency bands separately, by centering all the data corresponding to each frequency band at every time step in all spectrograms, which involved removing their mean, and dividing by their standard deviation. This separate standardization of frequency bands resulted in a smaller reconstruction error compared to standardizing across all the bands.

## Deep neural network

We designed a deep variational autoencoder (VAE; Kingma et Welling, 2014) of 15 layers with an intermediate hidden representation of 128 neurons that we refer to as the *voice latent space* (VLS). In an autoencoder model, the two sub-network components, the *Encoder* and the *Decoder*, are jointly learned on complementary tasks (Fig. 1a). The Encoder network (noted *Enc* hereafter; 7 layers) learns to map an input,  $s$  (a spectrogram of a sound), onto a (128-dimensional) *voice latent space* representation ( $z$ ; in blue in the middle of Fig. 1a), while the Decoder (noted *Dec* hereafter; 7 layers) aims at reconstructing the spectrogram  $s$  from  $z$ . The learning objective of the full model is to make the output spectrogram  $Dec(Enc(s))$  as close as possible to the original one  $s$ . This reconstruction objective is defined as the L2 loss,  $||Dec(Enc(s)) - s||^2$ . The parameters of the Encoder and of the Decoder are jointly learned using gradient descent to optimize the average L2 loss computed on the training set  $\sum_{s \in Training\ Set} ||Dec(Enc(s)) - s||^2$ . We trained this DNN on the Common Voice dataset (Ardila et al., 2020) according to VAE learning procedure (as explained in Kingma et Welling., 2019) until convergence (network architecture and particularities of the training procedure are provided in Supplementary Table 1), using the PyTorch python package (Paszke et al., 2019).

## Linear autoencoder

We trained a linear autoencoder on the same dataset (described above) to serve as a linear baseline. Both the *Encoder* and the *Decoder* networks consisted of a single fully-connected layer, without any activation functions. Similar to the VAE, the latent space obtained from the *Encoder* was a 128-dimensional vector. The parameters of both the *Encoder* and of the *Decoder* were jointly learned using gradient descent to optimize the average L2 loss computed on the training set.

## Neuroimaging data acquisition

Participants were scanned using a 3 Tesla Prisma scanner (Siemens Healthcare, Erlangen, Germany) equipped with a 64-channel receiver head-coil. Their movements were monitored during the acquisition using the software FIRMM (Dosenbach et al., 2017). The whole-head high-resolution structural scan acquired during the first session was a T1-weighted multi-echo MPRAGE (MEMPRAGE) (TR = 2.5 s, TE = 2.53, 4.28, 6.07, 7.86 ms, TI=1000 ms flip angle: 8°, matrix size = 208 × 300 × 320; resolution 0.8 × 0.8 × 0.8 mm<sup>3</sup>, acquisition time: 8min22s). Lower resolution scans acquired during all other sessions were T1-weighted MPRAGE scans (TR = 2.3 s, TE = 2.88 ms, TI=900ms, flip angle: 9°, matrix size = 192 × 240 × 256; resolution 1 × 1 × 1 mm<sup>3</sup>, sparse sampling with 2.8 times undersampling and compressed sensing reconstruction, acquisition time: 2min37). Functional imaging was performed using an EPI sequence (multiband factor = 5, TR = 462 ms, TE = 31.2 ms, flip angle: 45°, matrix size = 84 × 84 × 35, resolution 2.5 × 2.5 × 2.5 mm<sup>3</sup>). Functional slices were oriented parallel to the lateral sulci with a z-axis coverage of 87.5 mm, allowing it to fully cover both the TVAs (Pernet et al., 2015) and the FVAs (Aglieri et al., 2018). The physiological signals (heart rate and respiration) were measured with the external sensors of Siemens.

## Pre-processing of neuroimaging data and general linear modeling

Tissue segmentation and brain extraction was performed on the structural scans using the default segmentation procedure of SPM 12 (Ashburner et al., 2012). The preprocessing of the BOLD responses involved correcting motion, registering inter-runs, detrending and smoothing the data. Each functional volume was realigned to a reference volume taken from a steady period in the session that was spatially the closest to the average of all sessions. Transformation matrices between anatomical and functional data were computed using boundary-based registration (FSL; Smith et al., 2004). The data were respectively detrended and smoothed using the *nilearn* functions *clean\_img* and *smooth\_img* (kernel size of 3mm)

(Abraham et al., 2014), resulting in the matrix  $Y \in R^{S \times V}$ , with  $S$  the number of scans and  $V$  the number of voxels.

A first general linear model (GLM) was fit to regress out the noise by predicting  $Y$  from a “denoised” design matrix, composed of  $R = 38$  regressors of nuisance (Supplementary Figure 4). These regressors of nuisance, also called covariates of no interest, included: 6 head motion parameters (3 variable for the translations, 3 variables for the rotations); 18 ‘RETROICOR’ regressors (Glover et al., 2000) using the *TAPAS PhysIO* package (Kasper et al., 2017) (with the hyperparameters set as specified in Snoek et al.) were computed from the physiological signals; 13 regressors modeling slow artifactual trends (sines and cosines, cut frequency of the high-pass filter = 0.01 Hz); and a confound-mean predictor. The design matrix was convolved with an hemodynamic response function (HRF) with a peak at 6s and an undershoot at 16s (Glover et al., 1999), we note the convolved design matrix as  $X_d \in R^{S \times R}$ . The “denoise” GLM’s parameters  $\beta_d \in R^{R \times V}$  were optimized to minimize the amplitude of the residual  $\beta_d = \underset{\beta \in R^{R \times V}}{\operatorname{argmin}} ||Y - X_d \beta||^2$ . We used a lag-1 autoregressive model (ar(1)) to model the temporal structure of the noise (Friston et al., 2002). The *denoised* BOLD signal  $Y_d$  was then obtained from the original one according to  $Y_d = Y - (X_d \beta_d) \in R^{S \times V}$ .

A second “stimulus” GLM model was used to predict the denoised BOLD responses for each stimulus using a design matrix  $X_s \in R^{S \times (N_s + 1)}$  (which was convolved with an hemodynamic response function, HRF as above) and a parameters matrix  $\beta_s \in R^{(N_s + 1) \times V}$  where  $N_s$  stands for the number of stimuli. The last row (resp. column) of  $\beta_s$  (resp.  $X_s$ ) stands for a silence condition. Again,  $\beta_s$  was learned to minimize the residual  $\beta_s = \underset{\beta \in R^{(N_s + 1) \times V}}{\operatorname{argmin}} ||Y_d - X_s \beta||^2$ . Once learned, each of the first  $N_s$  line of  $\beta_s$  was corrected by subtracting the  $(N_s + 1)^{th}$  line, yielding the contrast maps for stimuli  $\tilde{\beta}_s \in R^{N_s \times V}$ . We note hereafter  $\tilde{\beta}_s[i, :] \in R^V$  the contrast map for a given stimulus, it is the  $i^{th}$  line of  $\tilde{\beta}_s$ .

A third “identity” GLM was fit to predict the BOLD responses of each voice speaker identity, using a design matrix  $\beta_i \in R^{(N_i+1) \times V}$  and a design matrix  $X_i \in R^{S \times (N_i+1)}$  (which was again convolved with an hemodynamic response function, HRF) where  $N_s$  stands for the number of unique speakers. Again the last row/column in  $\beta_i$  and  $X_i$  stands for the silent condition.  $\beta_i$  is learned to minimize the residual  $\beta_i = \operatorname{argmin}_{\beta \in R^{(N_i+1) \times V}} \|Y_d - X_i \beta\|^2$  (Supplementary Figure 3a). Again, the final speaker contrast maps were obtained by contrasting (i.e., subtracting) the regression coefficients in a row of  $\beta_i$  with the silence condition (last row; Supplementary Figure 3a), yielding  $\tilde{\beta}_i \in R^{N_s \times V}$ . Here the  $j^{\text{th}}$  row of  $\tilde{\beta}_i$ ,  $\tilde{\beta}_i[j, :] \in R^V$ , represents the amplitude of the BOLD response of the contrast map for speaker  $j$  (i.e. to all the stimuli from this speaker).

A fourth “localizer” GLM model was used to predict the denoised BOLD responses of each sound category from the *Voice localizer stimuli* presented above. The procedure was similar as described for the two previous GLM models. Once the GLM was learned, we contrasted the human voice category with the other sound categories in order to localize for each participant the posterior Temporal Voice Area (pTVA), medial Temporal Voice Area (mTVA) and anterior Temporal Voice Area (aTVA) in each hemisphere. The center of each TVA corresponded to the local maximum of the voice > non voice t-map whose coordinates were the closest to the TVAs reported in (Pernet et al., 2015). The analyses were carried on for each region of interest (ROI) of each hemisphere.

Additionally, we defined for each participant the primary auditory cortex (A1) as the maximum value of the probabilistic map (non-linearly registered to each participant functional space) of Heschl’s gyri provided with the MNI152 template (Penhune et al., 1996), intersected with the sound vs silence contrast map.

## Identity-based and stimulus-based representations

We performed analyses either at the stimulus level, e.g. predicting the neural activity of a participant listening to a given *stimulus* ( $\tilde{\beta}_s$ ’s lines) from the *voice latent space* representation

of this stimuli, or at the speaker identity level, e.g. predicting the average neural activity in response to stimuli of a given speaker *identity* ( $\tilde{\beta}_i$ 's lines) from this speaker's *voice latent space* representation. The identity-based analyses were used for the characterization of the *voice latent space* (Fig. 1), the brain encoding (Fig. 2), and the representational similarity analysis (Fig. 3), while the stimulus-based analyses were used for the brain decoding analyses (Fig. 4, 5).

We conducted stimulus-based analyses to examine the relationship between stimulus contrast maps in neural activity ( $\tilde{\beta}_s$ ) and the encodings of individual stimulus spectrograms computed by the encoder of an autoencoder model (either linear or deep variational autoencoder) on the computational side. We will note  $z_s^{lin} \in R^{N_s \times 128}$  encodings of stimuli by the LIN model and  $z_s^{vae} \in R^{N_s \times 128}$  the encodings of stimuli computed by the VAE model. The encoding of the  $k^{\text{th}}$  stimuli by one of these models is the  $k^{\text{th}}$  row of the corresponding matrix and it is noted as  $z_s^{model}[k, :]$ .

For identity-based analyses we studied relationships between identity contrast maps in  $\tilde{\beta}_i$  on the neural activity side, and an encoding of speaker identity in the VLS implemented by an autoencoder model (LIN or VAE) on the computational side, e.g. we note  $z_i^{vae}[j]$  the representation of speaker  $j$  as computed by the *vae* model. We chose to define a speaker identity-based representation as the average of a set of sample-based representations for stimuli from this speaker, e.g.  $z_i^{model}[j] = 1/|S_j| \sum_{k \in S_j} z_s^{model}[k, :]$  where  $S_j$  stands for the set of stimuli by speaker  $j$  and *model* stands for *vae* or *lin*. Averaging in the *voice latent space* is expected to be much more powerful and relevant than averaging in the input space spectrograms (VanRullen & Reddy, 2019).

## Characterization of the autoencoder latent space

We characterized the organization of the *voice latent space* (VLS) and of the features computed by the linear autoencoder (LIN) by measuring through classification experiments



the presence of information about speaker's gender, age, and identity in the representations learned by these models.

We first computed the speaker's identity *voice latent space* representations for each of the 405 speakers in the main voice dataset (135+142+128 see *Stimuli* section) as explained above.

Next we used these speakers' *voice latent space* representation to investigate if the gender, age, identity were encoded in the VLS. To do so we divided the data in separate train and test sets and learned classifiers to predict gender, age, or identity from the train set. The balanced (to avoid the small effects associated with unbalanced folds) accuracy of the classifiers were then evaluated on the test set. The higher the performance on the test set the more we are confident that the information is encoded in the VLS. More specifically for each task (gender, age, identity), we trained a Logistic Regression classifier (linear regularized logistic regression; L2 penalty, tol=0.0001, fit\_intercept=True, intercept\_scaling=1, max\_iter=100) using the scikit-learn python package (Pedregosa et al., 2018).

In order to statistically evaluate the significance of the results and to avoid a potential overfitting, the classifications were repeated 20 times with 20 different initializations (*seed*) and the metrics were then averaged for each voice category (gender, age). More specifically, we repeated the following experiment 20 times with 20 different random seeds. For each seed, we performed 5 train-test splits with 80% of the data in the training and 20% in the test set. For each split we used 5-fold cross validation on the training set to select the optimal value for the regularization hyperparameter C (searching between 10 values logarithmically spaced on the interval [-3, +3]). We then computed the generalization performance on the test set of the model trained on the full training set with the best hyperparameter value. Reported results were then averaged over 20 experiments. Note that data were systematically normalized with a scaler fitted on the training set. We used a robust

scaling strategy for these experiments (removing the median, then scaling to the quantile range; 25<sup>th</sup> quantile and 75<sup>th</sup> quantile) which occurs to be more relevant with a small training set.

To investigate how speaker identity information is encoded in the latent space representations of speakers' voices, we computed speaker identity *voice latent space* representations by averaging 20 stimulus-based representations, in order to obtain a limited amount of data per identity that could be distributed across training and test datasets.

We first tested whether the mean of the distribution of accuracy scores obtained for 20 seeds was significantly above chance level using one-sample t-tests. We then evaluated the difference in classification accuracy between the VLS and LIN via one-way ANOVAs (dependent variable: test balance accuracy; between factor: Feature), for each category (speaker gender, age, identity). We performed post-hoc planned paired t-tests between the models to test the significance of the VLS-LIN difference.

## Brain encoding

We performed encoding experiments on identity-based representations for each of the three participants (Fig. 2). For each participant we explored the ability to learn a regularized linear regression that predicts a speaker-based neural activity, e.g. the  $j^{\text{th}}$  speaker's contrast map  $\tilde{\beta}_i[j] \in R^V$ , from this speaker's voice latent space representation, that we note  $z_i^{\text{model}}[j] \in R^{128}$  (Fig. 2a). We carried out these regression analyses for each ROI (A1, pTVA, mTVA, aTVA) in each hemisphere and participant, independently.

The regression model parameters  $\hat{W}_{\text{encod}} \in R^{128 \times V}$  were learned according to:

$$\hat{W}_{\text{encod}} = \underset{W_{\text{encod}} \in R^{128 \times V}}{\operatorname{argmin}} \sum_{j=1..N_i} (z_i^{\text{model}}[j] \times W_{\text{encod}} - \tilde{\beta}_i[j])^2 + \lambda \|W_{\text{encod}}\|^2$$

where  $\lambda$  is a hyperparameter tuning the optimal tradeoff between the data fit and the penalization terms above. We used the ridge regression with built-in cross-validation as implemented as *RidgeCV* in the scikit-learn library (Pedregosa et al., 2018).

The statistical significance of each result was assessed with the following procedure. We repeated the following experiment 20 times with different random seeds. Each time, we performed 5 train-test splits with 80% of the data in the training and 20% in the test set. For each split we used RidgeCV (relying on leave-one-out) on the training set to select the optimal value for the hyperparameter  $\lambda$  (searching between 10 values logarithmically spaced on the interval  $[10^{-1}; 10^8]$ ). Following standard practice in machine learning, we then computed the generalization performance on the test set of the model trained on the full training set with the best hyperparameter value. Reported results are then averaged over 20 experiments. Note that here again with small training sets data were systematically normalized in each experiment using robust scaling.

Evaluation relied on the ‘brain score’ procedure (Schrimpf et al., 2018) which evaluates the performance of the ridge regression with a Pearson’s correlation score. Correlations between measured neural activities  $\tilde{\beta}_i$  and predicted ones  $\widehat{z_t^{model}} * W_{encod}$  were computed for each voxel and averaged over repeated experiments (folds and seeds) yielding one correlation value for every voxel and for every setting. The significance of the results was assessed with one-sample t-tests for the Fisher z-transformed correlation scores (3 x participants x 2 hemispheres x V voxels). For each region of interest, the scores are reported across participants and hemispheres (Fig. 2b). The exact same procedure was followed for the LIN modeling.

In order to determine which of the two feature spaces (VLS, LIN) and which of the two ROI (A1, TVAs) yielded the best prediction of neural activity, we compared the means of distributions of correlations coefficients using a mixed ANOVA performed on the Fisher z-

transformed coefficients (dependent variable: correlation; between factor: ROI; repeated measurements: Feature; between-participant identifier: voxel).

For each ROI, we then used t-tests to perform post-hoc contrasts for the VLS-LIN difference in brain encoding performance (comparison tests in Fig. 2b; Supplementary Table 4). We finally conducted two-sample t-tests between the brain encoding model's scores trained to predict A1 and those trained to predict temporal voice areas to test the significance of the A1-TVAs difference (Supplementary Table 5).

The statistical tests were all performed using the *pingouin* python package (Vallat., 2018).

## Representational similarity analysis

The RSA analyses were carried out using the package *rsatoolbox* (Schütt et al., 2021; <https://github.com/rsagroup/rsatoolbox>). For each participant, region of interest and hemisphere, we computed the cerebral Representational Dissimilarity Matrix (RDM) using the Pearson's correlation between the speaker identity-specific response patterns of the GLM estimates  $\tilde{\beta}_i$  (Walther et al., 2016) (Fig. 3a). The model RDMs were built using cosine distance (Xing et al., 2015; Bhattacharya et al., 2017; Wang et al., 2018), capturing speaker pairwise feature differences predicted by the computational models LIN and the VLS (Fig. 3a). For greater comparability with the rest of the analyses described here, the GLM estimates and the computational models' features were first normalized using robust scaling. We computed the Spearman correlations coefficients between the brain RDMs for each ROI, and the two model's RDMs (Fig. 3b). We assessed the significance of these brain-model correlation coefficients within a permutation-based 'maximum statistics' framework for multiple comparison correction (one-tailed inference; N permutations = 10,000 for each test; permutation of rows and columns of distance matrices, see Giordano et al., 2023 and Maris & Oostenveld, 2007; see Fig. 3b). We evaluated the VLS-LIN difference using a two-way repeated-measures ANOVA on the Fisher z-transformed Spearman correlation coefficients (dependent variable: correlation; within factors: ROI and Feature; participant identifier:

participant hemisphere pair). The same permutation framework was also used to assess the significance of the difference between the RSA correlation for the VLS and LIN models.

## Brain decoding

Brain decoding was investigated at the stimulus level. The stimuli's voice latent space representations  $z_s^{model} \in R^{N \times 128}$  and voice samples' contrast maps  $\tilde{\beta}_s \in R^{N \times V}$  were divided into train and test splits, normalized across voice samples using robust scaling, then fit to the training set. For every participant and each ROI, we trained a  $L_2$ -regularized linear model  $W \in R^{V \times 128}$  model to predict the voice samples' latent vectors from the voice samples' contrast maps (Fig. 4a). The hyperparameter selection and optimization was done similarly as in the Brain encoding scheme. Training was performed on non repeated stimuli (see Stimuli section). We then used the trained models to predict for each participant the 6 repeated stimuli that were the most presented. Waveforms were estimated starting from the reconstructed spectrograms using the Griffin-Lim phase reconstruction algorithm (Griffin & Lim, 1983).

We then used classifier analyses to assess the presence of voice information (gender, age, speaker identity) in the reconstructed latent representations (i.e., the latent representation predicted from the brain activity of a participant listening to a specific stimulus) (Fig. 5a, b, c). To this purpose, we first trained linear classifiers to categorize the training voice stimuli (participant 1, N = 6144; participant 2, N = 6142; participant 3, N = 5117; total, N = 17403) by gender (2 classes), age (2 classes) or identity (17 classes) based on VLS coordinates. Secondly, we used the previously trained classifiers to predict the identity information based on the VLS derived from the brain responses of the 18 Test voice stimuli (3 participants x 6 stimuli). We first tested using one-sample t-tests that the mean of the distribution of accuracy scores obtained across random classifier initializations of classifiers (2 hemispheres x 20 seeds = 40) was significantly above chance level, for each category, ROI and model. We then evaluated the difference in performance at preserving identity-related information depending on the model or ROI via two-way ANOVAs (dependent variable: accuracy;

between factors: Feature and ROI). We performed post-hoc planned paired t-tests between each model pair to test the significance of the VLS-LIN difference. Two-sample t-tests were finally used to test the significance of the A1-TVAs difference.

## **Listening tests**

We recruited 13 participants through the online platform Prolific ([www.prolific.co](http://www.prolific.co)) for a series of online behavioral experiments. All participants reported having normal hearing. The purpose of these experiments was to evaluate how well voice identity information and naturalness are preserved in fMRI-based reconstructed voice excerpts. In the main session, participants carried out 4 tasks, in the following order: ‘speaker discrimination’ (~120 min), ‘perceived naturalness’ (~30 min), ‘gender categorization’ (~30 min), ‘age categorization’ (~30 min). The experiment lasted 3 hours and 35 minutes, and each participant was paid £48.

Prior to the main experiment session, participants carried out a short loudness-change detection task to ensure that they wore headphones, and that they were attentive and properly set up for the main experiment (Woods et al., 2017). On each of 12 trials, participants heard 3 tones and were asked to identify which tone was the least loud by clicking on one of 3 response buttons: ‘First’, ‘Second’, or ‘Third’. Participants were admitted to the main experiment only if they achieved perfect performance in this task. We additionally refined the participant pool by excluding those who performed badly on the original stimuli.

The next three tasks were each carried out on the same set of 342 experimental stimuli, each presented on a different trial: 18 original stimuli, 36 stimuli reconstructed directly from the LIN and the VLS models, and 18 stimuli x 2 models x 4 regions of interest x 2 hemispheres= 288 brain-reconstructed stimuli.

In the ‘perceived naturalness’ task, participants were asked to rate how natural the voice sounded on a scale ranging from ‘Not at all natural’ to ‘Highly natural’ (i.e., similar to a real recording), and were instructed to use the full range of the scale.

During the ‘gender categorization’ task, participants categorized the gender by clicking on a ‘Female’ or ‘Male’ button.

Finally, in the ‘age categorization’ task, participants categorized the age of the speaker by clicking on a ‘Younger’ or ‘Older’ button.

In the ‘speaker discrimination’ task, participants carried out 684 trials (342 experimental stimuli x 2) with short breaks in between. On each trial, they were presented with 2 short sound stimuli, one after the other, and participants had to indicate whether they were from the same speaker or not.

To evaluate the performance of the participants, we firstly conducted one-sample t-tests to examine whether the mean accuracy score calculated from their responses was significantly higher than the chance level for each model and ROI. Next, we used two-way repeated-measures ANOVAs to assess the variation in participants’ performances in identifying identity-related information (dependent variable: accuracy; between-participant factors: Feature and ROI). To determine the statistical significance of the VLS-LIN difference, we carried out post-hoc planned paired t-tests between each model pair. Finally, we employed two-sample t-tests to evaluate the statistical significance of the A1-TVAs difference.

## **Data and code availability**

All data and codes will be made publicly available upon the article publication.



## References

- Abraham, A., Pedregosa, F., Eickenberg, M., Gervais, P., Mueller, A., Kossaifi, J., Gramfort, A., Thirion, B., & Varoquaux, G. (2014). Machine learning for neuroimaging with scikit-learn. *Frontiers in Neuroinformatics*, 8. <https://www.frontiersin.org/articles/10.3389/fninf.2014.00014>
- Akbari, H., Khalighinejad, B., Herrero, J. L., Mehta, A. D., & Mesgarani, N. (2019). Towards reconstructing intelligible speech from the human auditory cortex. *Scientific Reports*, 9(1), 874. <https://doi.org/10.1038/s41598-018-37359-z>
- Ardila, R., Branson, M., Davis, K., Henretty, M., Kohler, M., Meyer, J., Morais, R., Saunders, L., Tyers, F. M., & Weber, G. (2020). *Common Voice: A Massively-Multilingual Speech Corpus* (arXiv:1912.06670). arXiv. <https://doi.org/10.48550/arXiv.1912.06670>
- Ashburner, J. (2012). SPM: A history. *NeuroImage*, 62(2), 791–800. <https://doi.org/10.1016/j.neuroimage.2011.10.025>
- Barbero, F. M., Calce, R. P., Talwar, S., Rossion, B., & Collignon, O. (2021). Fast Periodic Auditory Stimulation Reveals a Robust Categorical Response to Voices in the Human Brain. *ENeuro*, 8(3), ENEURO.0471-20.2021. <https://doi.org/10.1523/ENeuro.0471-20.2021>
- Belin, P., Bestelmeyer, P. E. G., Latinus, M., & Watson, R. (2011). Understanding Voice Perception: Understanding voice perception. *British Journal of Psychology*, 102(4), 711–725. <https://doi.org/10.1111/j.2044-8295.2011.02041.x>
- Belin, P., Bodin, C., & Aglieri, V. (2018). A “voice patch” system in the primate brain for processing vocal information? *Hearing Research*, 366, 65–74. <https://doi.org/10.1016/j.heares.2018.04.010>

- 848 Belin, P., Fecteau, S., & Bédard, C. (2004). Thinking the voice: Neural correlates of  
849 voice perception. *Trends in Cognitive Sciences*, 8(3), 129–135.  
850 <https://doi.org/10.1016/j.tics.2004.01.008>
- 851 Belin, P., & Zatorre, R. J. (2003). Adaptation to speaker's voice in right anterior temporal  
852 lobe: *NeuroReport*, 14(16), 2105–2109. [https://doi.org/10.1097/00001756-](https://doi.org/10.1097/00001756-200311140-00019)  
853 [200311140-00019](https://doi.org/10.1097/00001756-200311140-00019)
- 854 Belin, P., Zatorre, R. J., Lafaille, P., Ahad, P., & Pike, B. (2000). Voice-selective areas in  
855 human auditory cortex. *Nature*, 403(6767), 309–312.  
856 <https://doi.org/10.1038/35002078>
- 857 Bhattacharya, G., Alam, J., & Kenny, P. (2017). Deep Speaker Embeddings for Short-  
858 Duration Speaker Verification. *Interspeech 2017*, 1517–1521.  
859 <https://doi.org/10.21437/Interspeech.2017-1575>
- 860 Blank, H., Wieland, N., & von Kriegstein, K. (2014). Person recognition and the brain:  
861 Merging evidence from patients and healthy individuals. *Neuroscience &*  
862 *Biobehavioral Reviews*, 47, 717–734.  
863 <https://doi.org/10.1016/j.neubiorev.2014.10.022>
- 864 Bodin, C., Trapeau, R., Nazarian, B., Sein, J., Degiovanni, X., Baurberg, J., Rapha, E.,  
865 Renaud, L., Giordano, B. L., & Belin, P. (2021). Functionally homologous  
866 representation of vocalizations in the auditory cortex of humans and macaques.  
867 *Current Biology*, S0960982221011477. <https://doi.org/10.1016/j.cub.2021.08.043>
- 868 Bouthillier, X., Delaunay, P., Bronzi, M., Trofimov, A., Nichyporuk, B., Szeto, J., Sepah,  
869 N., Raff, E., Madan, K., Voleti, V., Kahou, S. E., Michalski, V., Serdyuk, D., Arbel,  
870 T., Pal, C., Varoquaux, G., & Vincent, P. (2021). *Accounting for Variance in*  
871 *Machine Learning Benchmarks* (arXiv:2103.03098). arXiv.  
872 <http://arxiv.org/abs/2103.03098>

- 873 Capilla, A., Belin, P., & Gross, J. (2013). The Early Spatio-Temporal Correlates and  
874 Task Independence of Cerebral Voice Processing Studied with MEG. *Cerebral*  
875 *Cortex*, 23(6), 1388–1395. <https://doi.org/10.1093/cercor/bhs119>
- 876 Caucheteux, C., Gramfort, A., & King, J.-R. (2022). Deep language algorithms predict  
877 semantic comprehension from brain activity. *Scientific Reports*, 12(1), Article 1.  
878 <https://doi.org/10.1038/s41598-022-20460-9>
- 879 Caucheteux, C., Gramfort, A., & King, J.-R. (2023). Evidence of a predictive coding  
880 hierarchy in the human brain listening to speech. *Nature Human Behaviour*, 1–12.  
881 <https://doi.org/10.1038/s41562-022-01516-2>
- 882 Caucheteux, C., & King, J.-R. (2022). Brains and algorithms partially converge in natural  
883 language processing. *Communications Biology*, 5(1), Article 1.  
884 <https://doi.org/10.1038/s42003-022-03036-1>
- 885 Chang, N., Pyles, J. A., Marcus, A., Gupta, A., Tarr, M. J., & Aminoff, E. M. (2019).  
886 BOLD5000, a public fMRI dataset while viewing 5000 visual images. *Scientific*  
887 *Data*, 6(1), Article 1. <https://doi.org/10.1038/s41597-019-0052-3>
- 888 Charest, I., Pernet, C., Latinus, M., Crabbe, F., & Belin, P. (2013). Cerebral Processing  
889 of Voice Gender Studied Using a Continuous Carryover fMRI Design. *Cerebral*  
890 *Cortex*, 23(4), 958–966. <https://doi.org/10.1093/cercor/bhs090>
- 891 Charest, I., Pernet, C. R., Rousselet, G. A., Quiñones, I., Latinus, M., Fillion-Bilodeau,  
892 S., Chartrand, J.-P., & Belin, P. (2009). Electrophysiological evidence for an early  
893 processing of human voices. *BMC Neuroscience*, 10(1), 127.  
894 <https://doi.org/10.1186/1471-2202-10-127>
- 895 Chen, Z., Qing, J., Xiang, T., Yue, W., & Zhou, J. (2022). *Seeing Beyond the Brain:*  
896 *Conditional Diffusion Model with Sparse Masked Modeling for Vision Decoding.*  
897 <https://doi.org/10.48550/arXiv.2211.06956>

- 898 Cowen, A. S., Chun, M. M., & Kuhl, B. A. (2014). Neural portraits of perception:  
899 Reconstructing face images from evoked brain activity. *NeuroImage*, 94, 12–22.  
900 <https://doi.org/10.1016/j.neuroimage.2014.03.018>
- 901 Défossez, Alexandre, Charlotte Caucheteux, Jérémy Rapin, Ori Kabeli, et Jean-Rémi  
902 King. 2023. « Decoding Speech Perception from Non-Invasive Brain Recordings ». *Nature Machine Intelligence* 5(10):1097-1107. doi: [10.1038/s42256-023-00714-5](https://doi.org/10.1038/s42256-023-00714-5)  
903
- 904 Diedrichsen, J., & Kriegeskorte, N. (2017). Representational models: A common  
905 framework for understanding encoding, pattern-component, and representational-  
906 similarity analysis. *PLOS Computational Biology*, 13(4), e1005508.  
907 <https://doi.org/10.1371/journal.pcbi.1005508>
- 908 Dosenbach, N. U. F., Koller, J. M., Earl, E. A., Miranda-Dominguez, O., Klein, R. L.,  
909 Van, A. N., Snyder, A. Z., Nagel, B. J., Nigg, J. T., Nguyen, A. L., Wesevich, V.,  
910 Greene, D. J., & Fair, D. A. (2017). Real-time motion analytics during brain MRI  
911 improve data quality and reduce costs. *NeuroImage*, 161, 80–93.  
912 <https://doi.org/10.1016/j.neuroimage.2017.08.025>
- 913 Friston, K. J., Glaser, D. E., Henson, R. N. A., Kiebel, S., Phillips, C., & Ashburner, J.  
914 (2002). Classical and Bayesian inference in neuroimaging: Applications.  
915 *NeuroImage*, 16(2), 484–512. <https://doi.org/10.1006/nimg.2002.1091>
- 916 Friston, K. J., Holmes, A. P., Worsley, K. J., Poline, J.-P., Frith, C. D., & Frackowiak, R.  
917 S. J. (1994). Statistical parametric maps in functional imaging: A general linear  
918 approach. *Human Brain Mapping*, 2(4), 189–210.  
919 <https://doi.org/10.1002/hbm.460020402>
- 920 Gaziv, G., Beliy, R., Granot, N., Hoogi, A., Strappini, F., Golan, T., & Irani, M. (2022).  
921 Self-supervised Natural Image Reconstruction and Large-scale Semantic

Classification from Brain Activity. *NeuroImage*, 254, 119121.

<https://doi.org/10.1016/j.neuroimage.2022.119121>

Giordano, B. L., Esposito, M., Valente, G., & Formisano, E. (2023). Intermediate acoustic-to-semantic representations link behavioral and neural responses to natural sounds. *Nature Neuroscience*, 1–9. <https://doi.org/10.1038/s41593-023-01285-9>

Glover, G. H. (1999). Deconvolution of Impulse Response in Event-Related BOLD fMRI1. *NeuroImage*, 9(4), 416–429. <https://doi.org/10.1006/nimg.1998.0419>

Glover, G. H., Li, T. Q., & Ress, D. (2000). Image-based method for retrospective correction of physiological motion effects in fMRI: RETROICOR. *Magnetic Resonance in Medicine*, 44(1), 162–167. [https://doi.org/10.1002/1522-2594\(200007\)44:1<162::aid-mrm23>3.0.co;2-e](https://doi.org/10.1002/1522-2594(200007)44:1<162::aid-mrm23>3.0.co;2-e)

Griffin, D. & Jae Lim. (1983). Signal estimation from modified short-time Fourier transform. *ICASSP '83. IEEE International Conference on Acoustics, Speech, and Signal Processing*, 8, 804–807. <https://doi.org/10.1109/ICASSP.1983.1172092>

Güçlü, U., Thielen, J., Hanke, M., Gerven, M. A. J. van, & Gerven, M. A. J. van. (2016). Brains on Beats. In *Proceedings of the International Conference on Neural Information Processing Systems*, 2101–2109. <https://doi.org/null>

Gutierrez, M. (2021). Algorithmic Gender Bias and Audiovisual Data: A Research Agenda. *International Journal of Communication*, 15, 439–461.

Han, K., Wen, H., Shi, J., Lu, K.-H., Zhang, Y., Fu, D., & Liu, Z. (2019). Variational autoencoder: An unsupervised model for encoding and decoding fMRI activity in visual cortex. *NeuroImage*, 198, 125–136. <https://doi.org/10.1016/j.neuroimage.2019.05.039>

946 Harris, C. R., Millman, K. J., van der Walt, S. J., Gommers, R., Virtanen, P.,  
947 Cournapeau, D., Wieser, E., Taylor, J., Berg, S., Smith, N. J., Kern, R., Picus, M.,  
948 Hoyer, S., van Kerkwijk, M. H., Brett, M., Haldane, A., del Río, J. F., Wiebe, M.,  
949 Peterson, P., ... Oliphant, T. E. (2020). Array programming with NumPy. *Nature*,  
950 585(7825), 357–362. <https://doi.org/10.1038/s41586-020-2649-2>

951 Higgins, I., Chang, L., Langston, V., Hassabis, D., Summerfield, C., Tsao, D., &  
952 Botvinick, M. (2021). Unsupervised deep learning identifies semantic  
953 disentanglement in single inferotemporal face patch neurons. *Nature*  
954 *Communications*, 12(1), 6456. <https://doi.org/10.1038/s41467-021-26751-5>

955 Horikawa, T., & Kamitani, Y. (2017). Generic decoding of seen and imagined objects  
956 using hierarchical visual features. *Nature Communications*, 8(1), 15037.  
957 <https://doi.org/10.1038/ncomms15037>

958 Howard, A. G., Zhu, M., Chen, B., Kalenichenko, D., Wang, W., Weyand, T., Andreetto,  
959 M., & Adam, H. (2017). MobileNets: Efficient Convolutional Neural Networks for  
960 Mobile Vision Applications. *ArXiv:1704.04861 [Cs]*. <http://arxiv.org/abs/1704.04861>

961 Kasper, L., Bollmann, S., Diaconescu, A. O., Hutton, C., Heinzle, J., Iglesias, S.,  
962 Hauser, T. U., Sebold, M., Manjaly, Z.-M., Pruessmann, K. P., & Stephan, K. E.  
963 (2017). The PhysIO Toolbox for Modeling Physiological Noise in fMRI Data. *Journal*  
964 *of Neuroscience Methods*, 276, 56–72.  
965 <https://doi.org/10.1016/j.jneumeth.2016.10.019>

966 Kell, A. J. E., Yamins, D. L. K., Shook, E. N., Norman-Haignere, S. V., & McDermott, J.  
967 H. (2018). A Task-Optimized Neural Network Replicates Human Auditory Behavior,  
968 Predicts Brain Responses, and Reveals a Cortical Processing Hierarchy. *Neuron*,  
969 98(3), 630-644.e16. <https://doi.org/10.1016/j.neuron.2018.03.044>

970 Khaligh-Razavi, S.-M., & Kriegeskorte, N. (2014). Deep Supervised, but Not  
 971 Unsupervised, Models May Explain IT Cortical Representation. *PLoS*  
 972 *Computational Biology*, 10(11), e1003915.  
 973 <https://doi.org/10.1371/journal.pcbi.1003915>

974 Kingma, D. P., & Welling, M. (2014). Auto-Encoding Variational Bayes. *ArXiv:1312.6114*  
 975 *[Cs, Stat]*. <http://arxiv.org/abs/1312.6114>

976 Kingma, D. P., & Welling, M. (2019). An Introduction to Variational Autoencoders.  
 977 *Foundations and Trends® in Machine Learning*, 12(4), 307–392.  
 978 <https://doi.org/10.1561/22000000056>

979 Kriegeskorte, N. (2008). Representational similarity analysis – connecting the branches  
 980 of systems neuroscience. *Frontiers in Systems Neuroscience*.  
 981 <https://doi.org/10.3389/neuro.06.004.2008>

982 Kriegstein, K. V., & Giraud, A.-L. (2004). Distinct functional substrates along the right  
 983 superior temporal sulcus for the processing of voices. *NeuroImage*, 22(2), 948–955.  
 984 <https://doi.org/10.1016/j.neuroimage.2004.02.020>

985 Lavan, Nadine. 2023. « The Time Course of Person Perception From Voices: A  
 986 Behavioral Study ». *Psychological Science* 34(7):771-83. doi:  
 987 [10.1177/09567976231161565](https://doi.org/10.1177/09567976231161565).

988 Le, L., Ambrogioni, L., Seeliger, K., Güçlütürk, Y., van Gerven, M., & Güçlü, U. (2022).  
 989 Brain2Pix: Fully convolutional naturalistic video frame reconstruction from brain  
 990 activity. *Frontiers in Neuroscience*, 16.  
 991 <https://www.frontiersin.org/articles/10.3389/fnins.2022.940972>

992 LeCun, Y., Bengio, Y., & Hinton, G. (2015). Deep learning. *Nature*, 521(7553), Article  
 993 7553. <https://doi.org/10.1038/nature14539>



994       Liu, D., Li, Y., Lin, J., Li, H., & Wu, F. (2020). Deep Learning-Based Video Coding: A  
995       Review and a Case Study. *ACM Computing Surveys*, 53(1), 11:1-11:35.  
996       <https://doi.org/10.1145/3368405>

997       Maris, E., & Oostenveld, R. (2007). Nonparametric statistical testing of EEG- and MEG-  
998       data. *Journal of Neuroscience Methods*, 164(1), 177–190.  
999       <https://doi.org/10.1016/j.jneumeth.2007.03.024>

1000       Martin, F. N., & Champlin, C. A. (2000). Reconsidering the limits of normal hearing.  
1001       *Journal of the American Academy of Audiology*, 11(2), 64–66.

1002       Millet, J., Caucheteux, C., Orhan, P., Boubenec, Y., Gramfort, A., Dunbar, E., Pallier, C.,  
1003       & King, J.-R. (2022). *Toward a realistic model of speech processing in the brain*  
1004       *with self-supervised learning* (arXiv:2206.01685). arXiv.  
1005       <https://doi.org/10.48550/arXiv.2206.01685>

1006       Morillon, B., Arnal, L. H., & Belin, P. (2022). The path of voices in our brain. *PLOS*  
1007       *Biology*, 20(7), e3001742. <https://doi.org/10.1371/journal.pbio.3001742>

1008       Mozafari, M., Reddy, L., & VanRullen, R. (2020). Reconstructing Natural Scenes from  
1009       fMRI Patterns using BigBiGAN. *2020 International Joint Conference on Neural*  
1010       *Networks (IJCNN)*, 1–8. <https://doi.org/10.1109/IJCNN48605.2020.9206960>

1011       Nagrani, Arsha, Joon Son Chung, et Andrew Zisserman. 2017. « VoxCeleb: a large-  
1012       scale speaker identification dataset ». P. 2616-20 in *Interspeech 2017*.

1013       Naselaris, T., Kay, K. N., Nishimoto, S., & Gallant, J. L. (2011). Encoding and decoding  
1014       in fMRI. *NeuroImage*, 56(2), 400–410.  
1015       <https://doi.org/10.1016/j.neuroimage.2010.07.073>

1016       Pasley, B. N., David, S. V., Mesgarani, N., Flinker, A., Shamma, S. A., Crone, N. E.,  
1017       Knight, R. T., & Chang, E. F. (2012). Reconstructing Speech from Human Auditory

1018 Cortex. *PLOS Biology*, 10(1), e1001251.  
 1019 <https://doi.org/10.1371/journal.pbio.1001251>

1020 Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin,  
 1021 Z., Gimelshein, N., Antiga, L., Desmaison, A., Köpf, A., Yang, E., DeVito, Z.,  
 1022 Raison, M., Tejani, A., Chilamkurthy, S., Steiner, B., Fang, L., ... Chintala, S.  
 1023 (2019). *PyTorch: An Imperative Style, High-Performance Deep Learning Library*  
 1024 (arXiv:1912.01703). arXiv. <https://doi.org/10.48550/arXiv.1912.01703>

1025 Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel,  
 1026 M., Müller, A., Nothman, J., Louppe, G., Prettenhofer, P., Weiss, R., Dubourg, V.,  
 1027 Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., & Duchesnay,  
 1028 É. (2018). *Scikit-learn: Machine Learning in Python* (arXiv:1201.0490). arXiv.  
 1029 <https://doi.org/10.48550/arXiv.1201.0490>

1030 Penhune, V. B., Zatorre, R. J., MacDonald, J. D., & Evans, A. C. (1996).  
 1031 Interhemispheric anatomical differences in human primary auditory cortex:  
 1032 Probabilistic mapping and volume measurement from magnetic resonance scans.  
 1033 *Cerebral Cortex (New York, N.Y.: 1991)*, 6(5), 661–672.  
 1034 <https://doi.org/10.1093/cercor/6.5.661>

1035 Pernet, C. R., McAleer, P., Latinus, M., Gorgolewski, K. J., Charest, I., Bestelmeyer, P.  
 1036 E. G., Watson, R. H., Fleming, D., Crabbe, F., Valdes-Sosa, M., & Belin, P. (2015).  
 1037 The human voice areas: Spatial organization and inter-individual variability in  
 1038 temporal and extra-temporal cortices. *NeuroImage*, 119, 164–174.  
 1039 <https://doi.org/10.1016/j.neuroimage.2015.06.050>

1040 Petkov, C. I., Kayser, C., Steudel, T., Whittingstall, K., Augath, M., & Logothetis, N. K.  
 1041 (2008). A voice region in the monkey brain. *Nature Neuroscience*, 11(3), Article 3.  
 1042 <https://doi.org/10.1038/nn2043>

1043 Rupp, K., Hect, J. L., Remick, M., Ghuman, A., Chandrasekaran, B., Holt, L. L., & Abel,  
1044 T. J. (2022). Neural responses in human superior temporal cortex support coding of  
1045 voice representations. *PLOS Biology*, 20(7), e3001675.  
1046 <https://doi.org/10.1371/journal.pbio.3001675>

1047 Santoro, R., Moerel, M., De Martino, F., Valente, G., Ugurbil, K., Yacoub, E., &  
1048 Formisano, E. (2017). Reconstructing the spectrotemporal modulations of real-life  
1049 sounds from fMRI response patterns. *Proceedings of the National Academy of*  
1050 *Sciences*, 114(18), 4799–4804. <https://doi.org/10.1073/pnas.1617622114>

1051 Schrimpf, M., Blank, I. A., Tuckute, G., Kauf, C., Hosseini, E. A., Kanwisher, N.,  
1052 Tenenbaum, J. B., & Fedorenko, E. (2021). The neural architecture of language:  
1053 Integrative modeling converges on predictive processing. *Proceedings of the*  
1054 *National Academy of Sciences*, 118(45), e2105646118.  
1055 <https://doi.org/10.1073/pnas.2105646118>

1056 Schrimpf, M., Kubilius, J., Hong, H., Majaj, N. J., Rajalingham, R., Issa, E. B., Kar, K.,  
1057 Bashivan, P., Prescott-Roy, J., Geiger, F., Schmidt, K., Yamins, D. L. K., & DiCarlo,  
1058 J. J. (2018). *Brain-Score: Which Artificial Neural Network for Object Recognition is*  
1059 *most Brain-Like?* [Preprint]. Neuroscience. <https://doi.org/10.1101/407007>

1060 Schütt, H. H., Kipnis, A. D., Diedrichsen, J., & Kriegeskorte, N. (2021). *Statistical*  
1061 *inference on representational geometries* (arXiv:2112.09200). arXiv.  
1062 <http://arxiv.org/abs/2112.09200>

1063 Schweinberger, S. R., A. Herholz, et W. Sommer. 1997. « Recognizing Famous Voices:  
1064 Influence of Stimulus Duration and Different Types of Retrieval Cues ». *Journal of*  
1065 *Speech, Language, and Hearing Research: JSLHR* 40(2):453-63. doi:  
1066 10.1044/jslhr.4002.453.

1067 Smith, S. M., Jenkinson, M., Woolrich, M. W., Beckmann, C. F., Behrens, T. E. J.,  
 1068 Johansen-Berg, H., Bannister, P. R., De Luca, M., Drobnjak, I., Flitney, D. E.,  
 1069 Niazy, R. K., Saunders, J., Vickers, J., Zhang, Y., De Stefano, N., Brady, J. M., &  
 1070 Matthews, P. M. (2004). Advances in functional and structural MR image analysis  
 1071 and implementation as FSL. *NeuroImage*, 23, S208–S219.  
 1072 <https://doi.org/10.1016/j.neuroimage.2004.07.051>

1073 Snoek, L., van der Miesen, M. M., Beemsterboer, T., van der Leij, A., Eigenhuis, A., &  
 1074 Steven Scholte, H. (2021). The Amsterdam Open MRI Collection, a set of  
 1075 multimodal MRI datasets for individual difference analyses. *Scientific Data*, 8(1),  
 1076 Article 1. <https://doi.org/10.1038/s41597-021-00870-6>

1077 Trapeau, R., Thoret, E., & Belin, P. (2022). The Temporal Voice Areas are not “just”  
 1078 Speech Areas. *Frontiers in Neuroscience*, 16, 1075288.  
 1079 <https://doi.org/10.3389/fnins.2022.1075288>

1080 Vallat, R. (2018). Pingouin: Statistics in Python. *Journal of Open Source Software*,  
 1081 3(31), 1026. <https://doi.org/10.21105/joss.01026>

1082 VanRullen, R., & Reddy, L. (2019). Reconstructing faces from fMRI patterns using deep  
 1083 generative neural networks. *Communications Biology*, 2(1), 193.  
 1084 <https://doi.org/10.1038/s42003-019-0438-y>

1085 von Kriegstein, K., Eger, E., Kleinschmidt, A., & Giraud, A. L. (2003). Modulation of  
 1086 neural responses to speech by directing attention to voices or verbal content.  
 1087 *Cognitive Brain Research*, 17(1), 48–55. [https://doi.org/10.1016/S0926-](https://doi.org/10.1016/S0926-6410(03)00079-X)  
 1088 [6410\(03\)00079-X](https://doi.org/10.1016/S0926-6410(03)00079-X)

1089 Walther, A., Nili, H., Ejaz, N., Alink, A., Kriegeskorte, N., & Diedrichsen, J. (2016).  
 1090 Reliability of dissimilarity measures for multi-voxel pattern analysis. *NeuroImage*,  
 1091 137, 188–200. <https://doi.org/10.1016/j.neuroimage.2015.12.012>

- 1092 Wang, X., Xu, Y., Wang, Y., Zeng, Y., Zhang, J., Ling, Z., & Bi, Y. (2018).  
1093 Representational similarity analysis reveals task-dependent semantic influence of  
1094 the visual word form area. *Scientific Reports*, 8(1), 3047.  
1095 <https://doi.org/10.1038/s41598-018-21062-0>
- 1096 Wetzel, S. J. (2017). Unsupervised learning of phase transitions: From principal  
1097 component analysis to variational autoencoders. *Physical Review E*, 96(2), 022140.  
1098 <https://doi.org/10.1103/PhysRevE.96.022140>
- 1099 Woods, K. J. P., Siegel, M. H., Traer, J., & McDermott, J. H. (2017). Headphone  
1100 screening to facilitate web-based auditory experiments. *Attention, Perception, &*  
1101 *Psychophysics*, 79(7), 2064–2072. <https://doi.org/10.3758/s13414-017-1361-2>
- 1102 Wu, M. C.-K., David, S. V., & Gallant, J. L. (2006). COMPLETE FUNCTIONAL  
1103 CHARACTERIZATION OF SENSORY NEURONS BY SYSTEM IDENTIFICATION.  
1104 *Annual Review of Neuroscience*, 29(1), 477–505.  
1105 <https://doi.org/10.1146/annurev.neuro.29.051605.113024>
- 1106 Xing, C., Wang, D., Liu, C., & Lin, Y. (2015). Normalized Word Embedding and  
1107 Orthogonal Transform for Bilingual Word Translation. *Proceedings of the 2015*  
1108 *Conference of the North American Chapter of the Association for Computational*  
1109 *Linguistics: Human Language Technologies*, 1006–1011.  
1110 <https://doi.org/10.3115/v1/N15-1104>
- 1111 Yamins, D. L. K., & DiCarlo, J. J. (2016). Using goal-driven deep learning models to  
1112 understand sensory cortex. *Nature Neuroscience*, 19(3), 356–365.  
1113 <https://doi.org/10.1038/nn.4244>
- 1114 Zäske, R., Perlich, M.-C., & Schweinberger, S. R. (2016). To hear or not to hear: Voice  
1115 processing under visual load. *Attention, Perception, & Psychophysics*, 78(5), 1488–  
1116 1495. <https://doi.org/10.3758/s13414-016-1119-2>

1117 Zhang, Y., Ding, Y., Huang, J., Zhou, W., Ling, Z., Hong, B., & Wang, X. (2021).  
 1118 Hierarchical cortical networks of “voice patches” for processing voices in human  
 1119 brain. *Proceedings of the National Academy of Sciences*, 118(52), e2113887118.  
 1120 <https://doi.org/10.1073/pnas.2113887118>

## 1121 **Acknowledgements**

1122 We thank Bruno Nazarian for the design of an MRI-compatible button. We thank Jean-Luc  
 1123 Anton and Kepkee Loh for useful discussions. This work was funded by the European  
 1124 Research Council (ERC) under the European Union’s Horizon 2020 research and innovation  
 1125 program (grant agreement no. 788240). This work was performed in the Center IRM-  
 1126 INT@CERIMED (UMR 7289, AMU-CNRS), platform member of France Life Imaging network  
 1127 (grant ANR-11-INBS-0 0 06). This work, carried out within the Institute of Convergence ILCB  
 1128 (ANR-16-CONV-0002), has benefited from support from the French government (*France*  
 1129 *2030*), managed by the French National Agency for Research (ANR) and the Excellence  
 1130 Initiative of Aix-Marseille University (A\*MIDEX).