

Epreuve orale Sélection Internationale ENS 2016, Sciences Cognitives

Choisissez l'un des deux articles proposés pour un rapide commentaire. Naturellement, vu le temps limité pour la préparation, nous ne nous attendons pas à un commentaire détaillé, mais plutôt à une réflexion personnelle prenant comme point de départ l'article. En particulier, vous pourrez vous appuyer sur les questions suivantes :

- 1) Quelle est la tâche ?
- 2) Quel est le résultat principal ?
- 3) Quelle grande fonction est questionnée par cette étude ?
- 4) Comment les résultats enrichissent-ils de façon surprenante les théories actuelles sur cette grande fonction ?

LETTERS

A new perceptual illusion reveals mechanisms of sensory decoding

Mehrdad Jazayeri¹ & J. Anthony Movshon¹

Perceptual illusions are usually thought to arise from the way sensory signals are encoded by the brain, and indeed are often used to infer the mechanisms of sensory encoding¹. But perceptual illusions might also result from the way the brain decodes sensory information², reflecting the strategies that optimize performance in particular tasks. In a fine discrimination task, the most accurate information comes from neurons tuned away from the discrimination boundary^{3,4}, and observers seem to use signals from these 'displaced' neurons to optimize their performance^{5,6,7}. We wondered whether using signals from these neurons might also bias perception. In a fine direction discrimination task using moving random-dot stimuli, we found that observers' perception of the direction of motion is indeed biased away from the boundary. This misperception can be accurately described by a decoding model that preferentially weights signals from neurons whose responses best discriminate those directions. In a coarse discrimination task, to which a different decoding rule applies⁴, the same stimulus is not misperceived, suggesting that the illusion is a direct consequence of the decoding strategy that observers use to make fine perceptual judgments. The subjective experience of motion is therefore not mediated directly by the responses of sensory neurons, but is only developed after the responses of these neurons are decoded.

Subjects viewed a field of moving dots within a circular aperture around a fixation point for 1 s and reported whether the direction of motion was clockwise (CW) or counter-clockwise (CCW) of a decision boundary indicated by a bar outside the edge of the dot-field (Fig. 1a). On each trial, the boundary took a random position around the dot field, and a percentage of dots (3%, 6% or 12%) moved coherently in a randomly chosen direction within 22 degrees

of the boundary; the other dots moved randomly. After each trial, subjects pressed one of two keys to indicate their choice (CW or CCW). For 70% of trials, they were given feedback. On the remaining

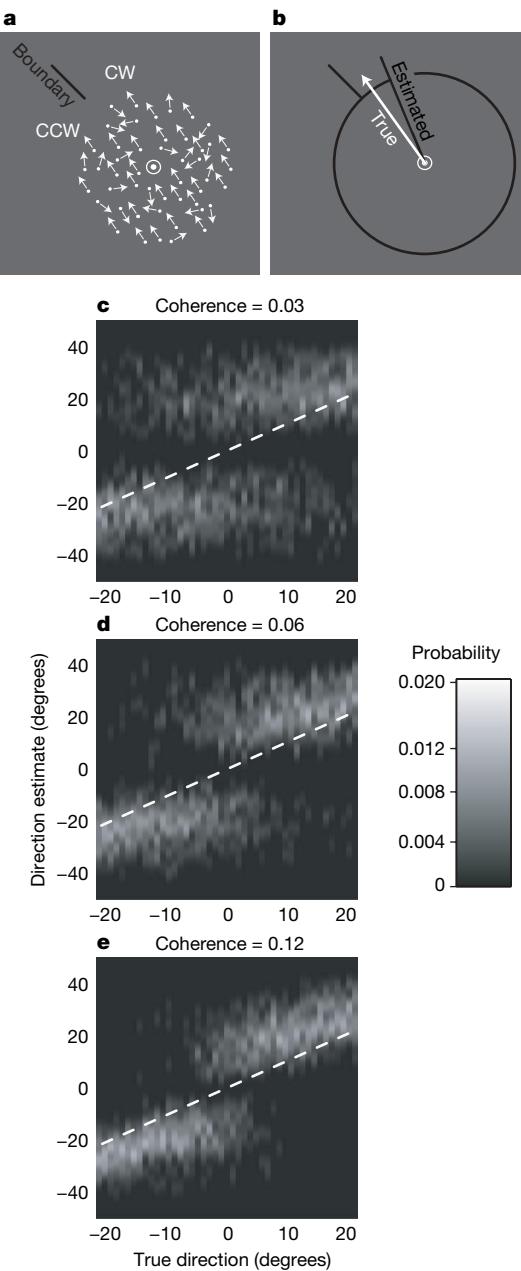


Figure 1 | The combined discrimination-estimation experiment. **a**, The discrimination phase. Subjects viewed a field of moving random dots and indicated whether its direction was clockwise (CW) or counter-clockwise (CCW) with respect to an indicated discrimination boundary that varied randomly from trial to trial. **b**, The estimation phase. On an unpredictable 30% of trials, after discriminating the direction of motion, subjects reported their estimate of the direction of motion by extending a dark line from the centre of the display with the computer mouse. **c–e**, Image maps representing the distribution of estimation responses for one subject at the three coherence levels. Each column of each plot represents the distribution of estimates for a particular true direction of motion, using a nonlinear lightness scale for probability (right). The observed values have been smoothed parallel to the ordinate with a gaussian (s.d. = 2 degrees) for clarity. The black dashed line is the locus of veridical estimates. Responses in the top-left and bottom-right quadrants of each map correspond to error trials, whereas those in the top-right and bottom-left quadrants were for correct trials; the discrimination and estimation responses were concordant throughout. Judgement errors (top-left and bottom-right quadrants) decrease with increasing coherence and as direction becomes more different from the discrimination boundary.

¹Center for Neural Science, New York University, 4 Washington Place, New York, New York 10003, USA.

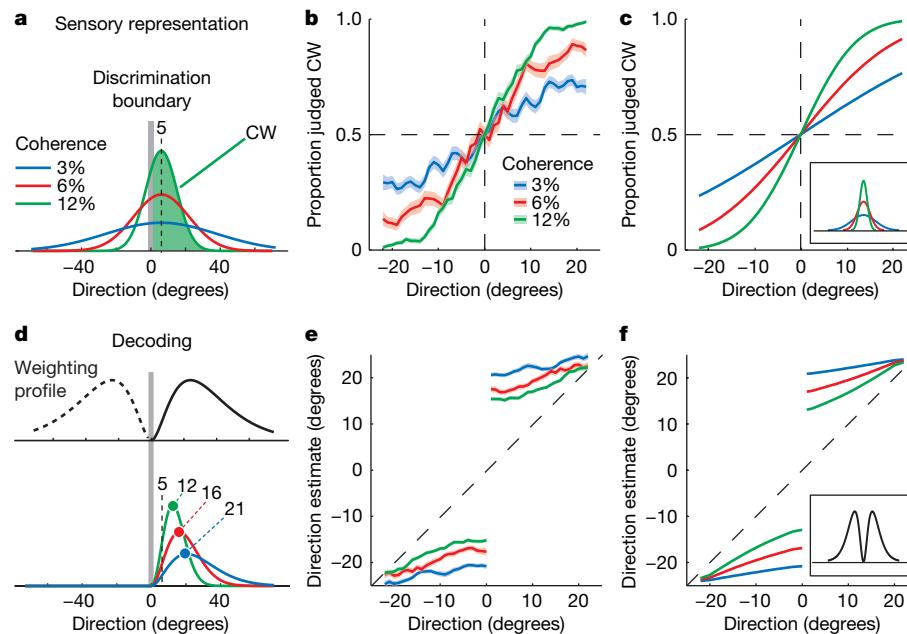


Figure 2 | Discrimination and estimation responses. **a**, The sensory representations evoked by a dot field moving 5° CW vary from one trial to the next. The plot is a cartoon of the distribution of these noise-perturbed sensory representations for different coherences. The distribution becomes more variable with weaker signals. The proportion of CW judgements is the area under the CW part of this distribution, shown for 12% coherence by the shaded green area. **b**, Proportion of CW judgements (thick lines) and their standard errors (shading) as a function of direction of motion for all coherence values for one subject. The CW and CCW portions of the data from Fig. 1b have been pooled and smoothed with a 3 degree boxcar filter. **c**, Fits to the discrimination data in **b**, using the model drawn in **a**; the inset shows the inferred sensory representations. **d**, The decoding model. The

sensory representations from **a** are multiplied by a displaced weighting profile that is optimal for discriminating CW from CCW alternatives. As a result, the peaks of the distributions shift away from the boundary. The plot shows schematically how this model predicts larger shifts for lower coherence values—the peaks for coherences of 3%, 6% and 12% fall at 12, 16 and 21 degrees respectively, even though the peak of the underlying sensory representation (from **a**) remains at 5 degrees. **e**, Subjective estimates as a function of direction of motion for trials on which motion direction was correctly discriminated. The CW and CCW portions of the data (Fig. 1d) have been pooled and smoothed with a 3 degree boxcar filter. **f**, The model fits for the subjective estimates after estimating and applying the single weighting profile that best matches the data.

30%, feedback was withheld and subjects estimated the direction of motion they had seen by aligning a bar extending from the fixation point to the direction of their estimate (Fig. 1b).

For all subjects, discrimination performance was lawfully related to motion coherence and direction: performance improved for higher coherences and for directions of motion farther away from

the boundary, and there was no systematic bias in the choice behaviour (Fig. 1c–e). However, when subjects were asked to report the direction of motion, their estimates deviated from the direction of motion in the stimulus, and were biased in register with their discrimination choice (Fig. 1c–e). The magnitude of these deviations depended on both the coherence and the direction of motion, being

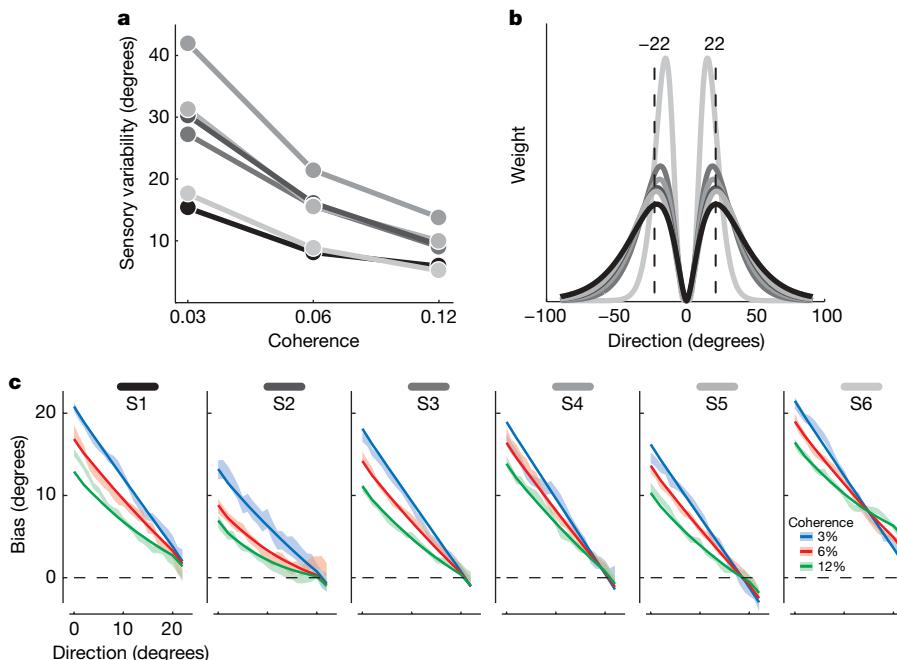


Figure 3 | Summary data for all six subjects. **a**, The variability of the sensory representations as a function of motion coherence for all subjects (computed as the standard deviation of the gaussian fits, for example, Fig. 2c, inset) are shown with different shades of grey (S1 to S6 in **c**). **b**, Recovered weighting functions for all six subjects. The dotted lines delimit the range of directions of motion (-22 to 22 degrees) that were used in the experiment. **c**, The mean bias (the difference between the true and estimated directions) \pm one standard error (shading) and the model fits (thick line) for all subjects and all coherence values.

larger for more uncertain conditions when either coherence was low or the direction was close to the boundary—the conditions in which discrimination performance was worst.

The sensory representation evoked by the random dot stimulus is perturbed by noise^{8,9,10}, and we therefore expect it to be more variable from trial to trial for weak motion signals than for stronger ones (Fig. 2a). How well observers discriminate the alternatives depends on the strength of the motion signal and its direction with respect to the boundary. Assuming that the variability in the sensory representation can be described by a gaussian, we fitted the discrimination performance of each subject (Fig. 2b) with a cumulative gaussian to estimate the spread of the sensory representation for each level of coherence. As expected, the variance of this distribution decreased with increasing coherence (Fig. 2c, inset).

This formulation accounts simply and well for discrimination behaviour (Fig. 2c), but it does not explain why the subjective estimates deviate from the true direction of motion in the stimulus. To understand what causes the perceptual biases, we considered the events that lead to the subjective estimates of the direction of motion. On each trial, before reporting their estimate, observers make a fine perceptual judgement. To do so, they have to transform the sensory responses into a binary decision (CW or CCW). Because subjects did not receive feedback on their subjective estimates, they could only adjust their decoding strategy for the discrimination part of the combined discrimination–estimation task where they did receive feedback. As shown both in theory^{3,4} and experiment^{6,7}, in a fine discrimination task like ours, neurons with direction preferences moderately shifted to the sides of the boundary make the largest contribution, whereas neurons tuned to directions either near or very remote from the boundary are less important. Therefore, to decode the activity of sensory neurons efficiently, the brain must pool their responses with a weighting profile that has maxima moderately shifted to the sides of the boundary⁴ (Fig. 2d, top panel).

If the pattern of direction estimates is explained by such a displaced profile, there should be a weighting function which, when applied to the sensory representation of different stimuli, predicts the corresponding estimates. We computed the product of this weighting profile (Fig. 2d, top panel) with the sensory representation of the stimulus estimated from the discrimination performance (Fig. 2a), and took the peak as the direction estimate (Fig. 2d, bottom panel). For each observer, we fitted the weighting profile that, when combined with that observer's discrimination performance, best predicted the pattern of direction estimates. Remarkably, combining the sensory representation with a single weighting profile (Fig. 2f, inset) accurately captured the observed estimates for all coherence levels and all directions of motion (Fig. 2f). Though observers varied in the accuracy of their sensory representations (Fig. 3a), the inferred weighting functions were similar for all six (Fig. 3b), and the resulting model accurately predicted the estimation bias (the difference between the true and estimated directions) for all six (Fig. 3c).

The misperception of motion can be economically attributed to the decoding strategy that observers adopt to optimize fine perceptual judgements, but other interpretations are possible. For example, the misperception might reflect a change in the sensory representation evoked by the stimulus, and not the way it is decoded. To test this idea, we ran a second experiment that differed from the first only in that the fine discrimination was replaced by coarse discrimination. On every trial, we presented motion in a randomly chosen direction within 22 degrees of a bar presented in the periphery (previously used for discrimination boundary), or within 22 degrees of the direction opposite the bar. Subjects discriminated whether the direction of motion was towards or away from the bar and as before, on a subset of trials reported their estimate of the direction of motion. As shown in theory⁴ and experiment¹¹, the most accurate information now comes from neurons tuned to the two alternatives. Therefore, the bias should, if anything, change from repulsion to attraction. This is exactly the pattern of responses we observed (Fig. 4a–c). The illusion thus

depends entirely on the subject's task—it occurs during fine discrimination but not during coarse discrimination (see Supplementary Information for a more detailed discussion). Changes in the sensory representation therefore cannot explain the effect.

One other possibility is that observers did not 'truly' misperceive the motion, but when uncertain about its direction, adopted a biased response strategy to ensure that they would not disagree with their immediately preceding discrimination choice. Simple models of response bias that only depend on the preceding choices can easily be discarded because they cannot account for the systematic relationship between the subjective estimates and the strength and direction of motion (Figs 1c–e, 2e). As detailed in the Supplementary Discussion, biased response strategies that are rich enough to account for our data have to incorporate computations effectively the same as those that our decoding model employs, applying a weighting profile to the sensory representation. Should we then view these computations in the framework of sensory decoding or complex response bias? In our decoding model, all computations serve a well-grounded function that can be inferred from theoretical and experimental observations of fine discrimination^{3,4,6,7}. This model also has the

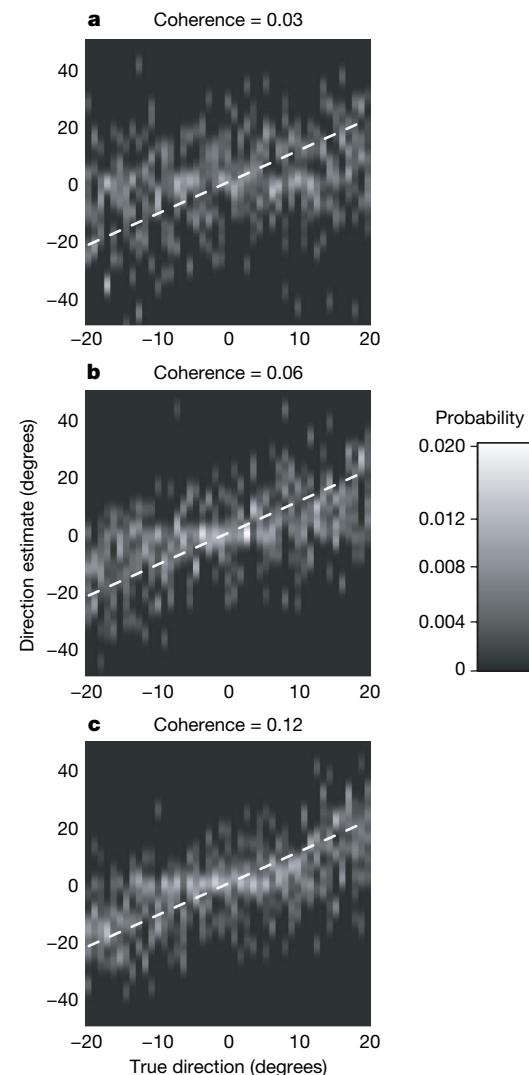


Figure 4 | Subjective estimates in a coarse direction discrimination task. **a–c**, The estimates are represented as in Fig. 1c–e. In the discrimination phase, the subject's task was to indicate whether the direction of motion in the random-dot was towards or away from a peripheral visual cue (the same as the discrimination boundary in Fig. 1b). The distributions for correct discriminations towards and away from the visual cue are pooled (separate plots for the two conditions are shown in Supplementary Fig. 3).

virtue of simplicity, because it accounts for observers' subjective reports using only the machinery that accounts for their objective discrimination performance. Response bias models, on the other hand, postulate two unrelated mechanisms, one to account for discrimination and the other for perceptual reports.

Bias arising from a decoding strategy may also explain some other perceptual distortions, such as those related to repulsion away from the cardinal axes¹² or from other discrimination boundaries¹³. We believe that this 'reference repulsion' phenomenon arises when subjects implicitly discriminate stimulus features against available internal or external references, such as a cardinal direction or the boundary marker in our experiment. This causes their perception of those features to be shifted away from the reference by the mechanism we have described. In other words, these incidences of misperception reveal the optimality of the system—not in perceiving, but in decoding sensory signals to make fine perceptual judgments.

Because the misperception does not seem to reflect the sensory responses to the direction of motion, the subjective experience of motion must be mediated by the machinery that decodes the responses of motion-sensitive neurons. We have argued elsewhere that areas downstream of sensory representations recode sensory responses into sensory likelihoods⁴, and the discrimination model used here is derived directly from that representation. Our results therefore suggest that the subjective experience of sensory events arises from the representation of sensory likelihoods, and not directly from the responses of sensory neuron populations.

METHODS

Eight subjects aged 19 to 35 participated in this study after giving informed consent. All had normal or corrected-to-normal vision, and all except one were naïve to the purpose of the experiment. Subjects viewed all stimuli binocularly from a distance of 71 cm on an Eizo T960 monitor driven by a Macintosh G5 computer at a refresh rate of 120 Hz in a dark, quiet room.

In the main experiment, in which six of the subjects participated, each trial began with the presentation of a fixation point along with a dark bar in the periphery representing the discrimination boundary for the subsequent motion discrimination (Fig. 1a). After 0.5 s, the motion stimulus was presented for 1 s. Subjects were asked to keep fixation during the presentation of the motion stimulus. After the motion stimulus was extinguished, subjects pressed one of two keys to report whether the direction of motion was CW or CCW with respect to the boundary and received distinct auditory feedback for correct and incorrect judgements. On approximately 30% of trials chosen at random, feedback was withheld and a circular ring was presented as a cue for the subject to report the direction of motion in the stimulus (Fig. 1b). The subject reported the estimate by using a mouse to extend a dark bar from the fixation point in the direction of their estimate and terminated the trial by pressing a key. Subjects were asked to estimate accurately but did not receive feedback. The discrimination boundary and the fixation point persisted throughout the trial. Trials were separated with a 1.5 s inter-trial interval during which the screen was blank.

For the main experiment, subjects had ample time to practice and master the task contingencies. Data for the main experiment were collected only after the discrimination thresholds stabilized (changed less than 10% across consecutive sessions). After this period, subjects completed roughly 8,000 trials in 10–12 sessions, each lasting approximately 45 min.

The remaining two subjects participated in the control coarse discrimination experiment. This experiment differed from the main experiment in that motion was within 22 degrees either towards or opposite to the peripherally presented bar (that is, the discrimination boundary in the main experiment), and during the discrimination stage, subjects had to report whether motion was towards or away from the bar. On approximately 30% to 50% of the trials chosen at random, the feedback was withheld and subjects were asked to report the perceived direction of motion (same procedure as in the main experiment).

All stimuli were presented on a dark grey background of 11 cd m^{-2} . The fixation point was a central circular white point subtending 0.5 degrees with a luminance of 77 cd m^{-2} . A gap of 1 degree between the fixation point and the motion stimulus helped subjects maintain fixation. The discrimination boundary was a black bar 0.5 degrees by 0.15 degrees, 3.5 degrees from fixation. The motion stimulus was a field of dots (each 0.12 degrees in diameter with a

luminance of 77 cd m^{-2}) contained within a 5-degree circular aperture centred on the fixation point (Fig. 1). On successive video frames, some dots moved coherently in a designated direction at a speed of 4 deg s^{-1} , and the others were replotted at random locations within the aperture. On each trial, the percentage of coherently moving dots (coherence) was randomly chosen to be 3, 6 or 12%, and their direction was randomly set to a direction within 22 degrees of the discrimination boundary; in the second experiment, half the trials presented motion within 22 degrees of a direction 180 degrees away from the boundary. The dots had an average density of 40 dots $\text{deg}^{-2} \text{ s}^{-1}$. The presentation of a black circular ring with a radius of 3.3 degrees around the fixation cued the subjects to report their estimate, which they did by moving the mouse to extend and align a black bar of width 0.15 deg to the direction of their estimate.

We modelled the sensory representation with a gaussian probability density function centred at the true direction of motion. The variance of this distribution for each subject was estimated by fitting a cumulative gaussian to his/her discrimination performance to maximize the likelihood of observing the subjects' choices for each level of motion coherence. To predict the direction estimates, we multiplied this sensory representation by a weighting function, and took the peak of the result. An additional additive constant accounted for any motor bias independent of sensory evidence. We chose a gamma probability density function as a convenient parametric form for the weighting profile. For each subject, we obtained parametric fits by minimizing the squared error of the model's prediction for the observed mean direction estimates for that subject. The gamma distribution provided a good fit for our data, but our conclusions do not depend on the exact form of the weighting profile. This simple procedure of first finding the sensory likelihoods from discrimination data, and then computing the weighting profile from the estimation data, crystallizes the contrast between encoding and decoding in our model. In detail, however, it neglects the subtle effect of the weighting profile on discrimination behaviour. In Supplementary Methods, we detail a complete model that simultaneously accounts for both the discrimination choices and the direction estimates in a single step (Supplementary Fig. 1), and present the fits for the correct as well as error trials (Supplementary Fig. 2).

Received 7 November 2006; accepted 9 March 2007.

Published online 4 April 2007.

- Eagleman, D. M. Visual illusions and neurobiology. *Nature Rev. Neurosci.* **2**, 920–926 (2001).
- Gregory, R. L. *Eye and Brain: The Psychology of Seeing* 5th edn (Oxford Univ. Press, 1997).
- Seung, H. S. & Sompolinsky, H. Simple models for reading neuronal population codes. *Proc. Natl. Acad. Sci. USA* **90**, 10749–10753 (1993).
- Jazayeri, M. & Movshon, J. A. Optimal representation of sensory information by neural populations. *Nature Neurosci.* **9**, 690–696 (2006).
- Patterson, R. D. Auditory filter shapes derived with noise stimuli. *J. Acoust. Soc. Am.* **59**, 640–654 (1976).
- Regan, D. & Beverley, K. I. Postadaptation orientation discrimination. *J. Opt. Soc. Am. A* **2**, 147–155 (1985).
- Hol, K. & Treue, S. Different populations of neurons contribute to the detection and discrimination of visual motion. *Vision Res.* **41**, 685–689 (2001).
- Parker, A. J. & Newsome, W. T. Sense and the single neuron: probing the physiology of perception. *Annu. Rev. Neurosci.* **21**, 227–277 (1998).
- Dean, A. F. The variability of discharge of simple cells in the cat striate cortex. *Exp. Brain Res.* **44**, 437–440 (1981).
- Tolhurst, D. J., Movshon, J. A. & Dean, A. F. The statistical reliability of signals in single neurons in cat and monkey visual cortex. *Vision Res.* **23**, 775–785 (1983).
- Britten, K. H., Shadlen, M. N., Newsome, W. T. & Movshon, J. A. The analysis of visual motion: a comparison of neuronal and psychophysical performance. *J. Neurosci.* **12**, 4745–4765 (1992).
- Huttenlocher, J., Hedges, L. V. & Duncan, S. Categories and particulars: Prototype effects in estimating spatial location. *Psychol. Rev.* **98**, 352–376 (1991).
- Rauber, H. & Treue, S. Reference repulsion when judging the direction of visual motion. *Perception* **27**, 393–402 (1998).

Supplementary Information is linked to the online version of the paper at www.nature.com/nature.

Acknowledgements This work was supported by a research grant from the NIH. We are grateful to B. Lau, E. Simoncelli, D. Heeger, M. Landy and N. Graham for advice and discussion.

Author Information Reprints and permissions information is available at www.nature.com/reprints. The authors declare no competing financial interests. Correspondence and requests for materials should be addressed to M.J. (mjaz@cnr.nyu.edu).

Selective cortical representation of attended speaker in multi-talker speech perception

Nima Mesgarani¹ & Edward F. Chang¹

Humans possess a remarkable ability to attend to a single speaker's voice in a multi-talker background^{1–3}. How the auditory system manages to extract intelligible speech under such acoustically complex and adverse listening conditions is not known, and, indeed, it is not clear how attended speech is internally represented^{4,5}. Here, using multi-electrode surface recordings from the cortex of subjects engaged in a listening task with two simultaneous speakers, we demonstrate that population responses in non-primary human auditory cortex encode critical features of attended speech: speech spectrograms reconstructed based on cortical responses to the mixture of speakers reveal the salient spectral and temporal features of the attended speaker, as if subjects were listening to that speaker alone. A simple classifier trained solely on examples of single speakers can decode both attended words and speaker identity. We find that task performance is well predicted by a rapid increase in attention-modulated neural selectivity across both single-electrode and population-level cortical responses. These findings demonstrate that the cortical representation of speech does not merely reflect the external acoustic environment, but instead gives rise to the perceptual aspects relevant for the listener's intended goal.

Separating out a speaker of interest from other speakers in a noisy, crowded environment is a perceptual feat that we perform routinely. The ease with which we hear under these conditions belies the intrinsic complexity of this process, known as the cocktail party problem^{1–3,6}: concurrent complex sounds, which are completely mixed upon entering the ear, are re-segregated and selected from within the auditory system. The resulting percept is that we selectively attend to the desired speaker while tuning out the others.

Although previous studies have described neural correlates of masking and selective attention to speech^{4,5,7–9}, fundamental questions remain unanswered regarding the precise nature of speech representation at the juncture where competing signals are resolved. In particular, when attending to a speaker within a mixture, it is unclear what key aspects (for example, spectrotemporal profile, spoken words and speaker identity) are represented in the auditory system and how they compare to representations of that speaker alone; how rapidly a selective neural representation builds up when one attends to a specific speaker; and whether breakdowns in these processes can explain distinct perceptual failures, such as the inability to hear the correct words, or follow the intended speaker.

To answer these questions, we recorded cortical activity from human subjects implanted with customized high-density multi-electrode arrays as part of their clinical work-up for epilepsy surgery¹⁰. Although limited to this clinical setting, these recordings provide simultaneous high spatial and temporal resolution while sampling the population neural activity from the non-primary auditory speech cortex in the posterior superior temporal lobe. We focused our analysis on high gamma (75–150 Hz) local field potentials¹¹, which have been found to correlate well with the tuning of multi-unit spike recordings¹². In humans, the posterior superior temporal gyrus has been heavily implicated in speech perception¹³, and is anatomically defined as the

lateral parabelt auditory cortex (including Brodmann areas 41, 42 and 22)¹⁴.

Subjects listened to speech samples from a corpus commonly used in multi-talker communication research^{15,16}. A typical sentence was “ready tiger go to red two now” where “tiger” is the call sign, and “red two” is the colour-number combination. One male and one female speaker were selected, each speaking the same 12 unique combinations of two call signs (ringo or tiger), three colours (red, blue or green) and three numbers (two, five or seven). Example acoustic spectrograms from two individual speakers are shown in Fig. 1a, b. The two voices differ along several dimensions including pitch (male versus female), spectral profile (different vocal tract shapes) and temporal characteristics (speaking rate). Subjects first listened to each of the speakers alone and were able to report the colour and number with 100% accuracy. Subjects then listened to a monaural, simultaneous mixture of the two speakers' phrases with different call signs, colours and numbers. The subjects were instructed to respond by indicating the colour and number spoken by the talker who uttered the target call sign. The target call sign (ringo or tiger) was fixed and shown visually on a monitor during each trial block, which contained 28 different mixture sounds. As the target speaker was changed randomly from trial to trial, the subjects were required to monitor both voices initially (divided attention) to identify the target speaker. The target call sign was switched after each block, turning the previous target speaker in each mixture into a masker. This resulted in two sets of behavioural and neural responses for each identical mixture sound, which differed only in the focus of attention. Subjects reported correct responses in 74.8% of trials.

Figure 1c illustrates the mixture spectrogram and how difficult it is to tell which sound parts belong to one speaker versus the other. The energy for both speakers is distributed broadly across the spectral and temporal domains, with overlap in some areas and isolated sound parts in others, as shown in their difference spectrogram (Fig. 1d; average spectrograms in Supplementary Fig. 1a).

To determine the spectrotemporal encoding of the attended speaker, the method of stimulus reconstruction was used^{17–19} to estimate the speech spectrogram represented by the population neural responses. Reconstructed spectrograms provide an intuitive way to examine how the population neural responses encode the spectrotemporal features of speech, and more importantly, can be compared with the original acoustic spectrograms as well as across attentional conditions. We first calculated the reconstruction filters from a passive listening task using a separate continuous speech corpus (TIMIT²⁰) that consisted of 499 unique short sentences spoken by 402 different speakers. The filters were then fixed and applied to a novel set of population neural responses to the single and attended mixture speech for spectrogram reconstruction.

When listening to a single speaker alone, the reconstructed spectrograms from population neural activity corresponded well to the spectrotemporal features of the original acoustic spectrograms (Fig. 1e, f compared to Fig. 1a, b, respectively), exhibiting fairly precise temporal features and spectral selectivity (for example, correspondence between the high frequency bursts of energy in “tiger” and “two”, in Fig. 1a, b, e, f).

¹Departments of Neurological Surgery and Physiology, UCSF Center for Integrative Neuroscience, University of California, San Francisco, California 94143, USA.

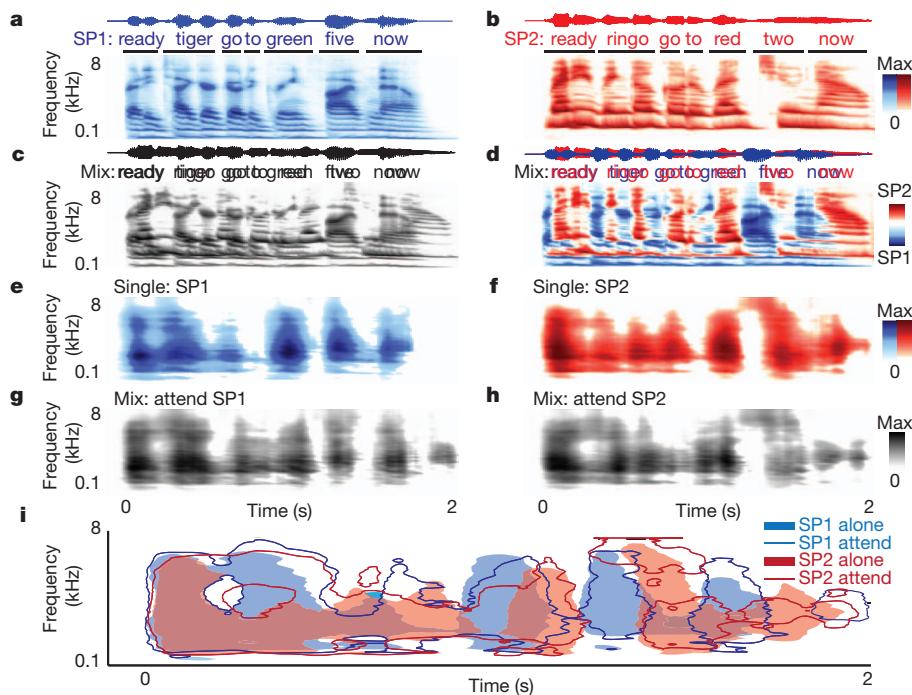


Figure 1 | Acoustic and neural reconstructed spectrograms for speech from a single speaker or a mixture of speakers. **a, b,** Example acoustic waveform and auditory spectrograms of speaker one (male; **a**) and speaker two (female; **b**). **c, d,** Waveform and spectrogram of the mixture of the two shows highly overlapping energy distributions. **d,** Difference spectrogram highlights the mixture regions where speaker one (blue) or two (red) has more acoustic energy. **e, f,** Neural-population-based stimulus reconstruction of speaker one (**e**) and speaker two (**f**) alone shows similar spectrotemporal features as the original spectrograms in **a** and **b**. **g, h,** The reconstructed spectrograms from the same mixture sound when attending to either speaker one (**g**) or two (**h**) highly resemble the single speaker reconstructions, shown in **e** and **f**, respectively. **i,** Overlay of the spectrogram contours at 50% of maximum energy from the reconstructed spectrograms in **e, f** and **h**.

The average and standard deviation of the correlation between reconstructed and original spectrograms over 24 sentences were 0.60 ± 0.034 (0.60 and 0.62 for the examples in Fig. 1e, f). When attending to each of the two speakers, the reconstructed spectrograms from the same speech mixture showed a marked difference depending upon which speaker was attended (Fig. 1g, h). For each pair, the key temporal and spectral features of the target speaker are enhanced relative to the masker speaker (Fig. 1g, h compared to Fig. 1e, f, respectively). To compare directly, the energy contours from these reconstructed spectrograms are overlaid in Fig. 1i. Important spectrotemporal details of the attended speaker were extracted, while the masker speech was effectively suppressed.

Attentional modulation of the neural representation was quantified, separately for correct and error trials, by measuring the correlation of the reconstructed spectrograms from the mixture in two attended conditions with original acoustic spectrograms of the speakers alone (Fig. 2a–d). During correct trials (Fig. 2a, c), we observed a significant shift of average correlation values towards the target speaker representation. During error trials, in contrast, no significant shift was

observed (Fig. 2b, d). Furthermore, the correlations between the reconstructed mixture and the masker speaker were higher than the average intrinsic correlation between randomly chosen original acoustic speech phrases (Fig. 2c, d, dashed lines), revealing a weak presence of the masker speaker in mixture reconstructions, even in correct trials.

The difference in speaking rate of the two speakers, coupled with the stereotyped structure of the carrier phrases, results in specific average temporal modulation profiles for each speaker (average spectrogram for each speaker is shown in Supplementary Fig. 1a, b). To investigate encoding of the distinct spectral profile and characteristic temporal rhythm of the target compared to the masker speaker, we estimated the average difference between reconstructed spectrograms of the two speakers, when presented alone and in the attended mixture (Fig. 2e, f). The comparison between the two average difference reconstructed spectrograms reveals enhanced encoding of both temporal and spectral aspects of the attended speaker (Supplementary Fig. 1c, d). To study the time course of attention-induced modulation of reconstructed mixture spectrograms towards the attended speaker, we

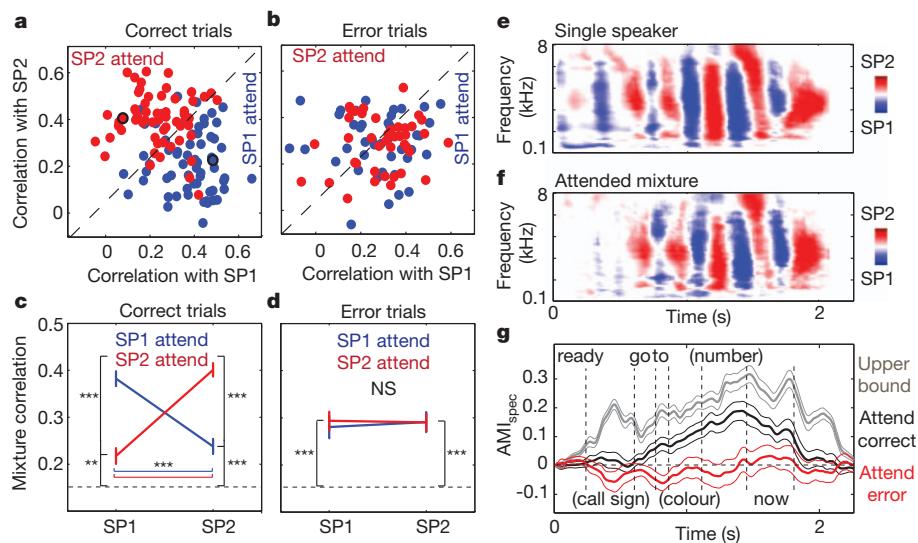


Figure 2 | Quantifying the attentional modulation of neural responses. **a, b,** Correlation coefficients of reconstructed mixture spectrograms under attentional control and the corresponding single speaker original spectrograms in correct and error trials (examples in Fig. 1g, h shown with black outline). **c, d,** Mean and standard error of correlation values for correct and error trials (28 mixtures). The dashed line corresponds to the average intrinsic correlation between randomly chosen original speech phrases. Brackets indicate pairwise statistical comparisons. NS, not significant. **e, f,** Average difference reconstructed spectrograms of speakers one and two from responses to single speaker (**e**) and attended mixture (**f**). **g,** Time course of average and standard error of AMI_{spec} of 28 mixtures for correct (black) and error (red) trials. Grey curve shows the upper bound of AMI_{spec} .

calculated an attentional modulation index (AMI_{spec}), using a sliding window of 250 ms throughout the trial duration:

$$AMI_{spec} = \text{Corr}(SP1_{spec}, SP1_{attend}) - \text{Corr}(SP1_{spec}, SP2_{attend}) + \text{Corr}(SP2_{spec}, SP2_{attend}) - \text{Corr}(SP2_{spec}, SP1_{attend}) \quad (1)$$

where $SP1_{spec}$ and $SP2_{spec}$ are the original acoustic spectrograms of speakers one and two, respectively, and $SP1_{attend}$ and $SP2_{attend}$ are the spectrograms reconstructed from neural responses to the mixture with attended targets, speaker one and two, respectively. Positive values of this index reflect shifts towards the target, negative values reflect shifts to the masker representation, and values around zero reflect no shift ($AMI_{spec} = 0.58$ for the example in Fig. 1). An upper bound for the AMI_{spec} was calculated by assuming that attention, at best, restores the single speaker reconstructions of the target speaker (replacing $SP1_{attend}$ and $SP2_{attend}$ in equation (1) with $SP1_{alone}$ and $SP2_{alone}$; Fig. 2g, grey line). The AMI_{spec} from the mixture was first estimated from correct trials (Fig. 2g, black line), and could resolve the time point at which the reconstructed spectrograms were modulated by attention. After the end of the call sign, which cues the speaker that should be attended, a rapid positive shift in the AMI_{spec} was observed, implying the enhanced representation of the target speaker. In error trials, this effect shows a bias towards the masker speaker, which, in contrast, occurred far earlier in the time course. The neural response shift towards the masker, which occurs as early as the call sign, suggests that listeners had prematurely attended to the wrong speaker during those error trials.

Although the reconstruction analyses showed clear attention-based spectrotemporal modulation, we wanted to determine explicitly whether the attended speech in a mixture could be decoded from a model of a single speaker. A regularized linear classifier²¹ was trained on neural responses to the single speakers and then used to decode both the spoken words and speaker identity of the attended speech mixture. To keep the chance performance at 50% across all comparisons, classification results were limited only to the choices that were present in each mixture. For correct trials, the colour and number of the attended speech were decoded with high accuracy (77.2% and 80.2%, $P < 10 \times 10^{-4}$, t -test; Fig. 3a). However, the decoding performance during error trials was significantly below chance (30.0%, 30.1%, $P < 10 \times 10^{-4}$, t -test; Fig. 3b), indicating a systematic bias towards decoding the words of the masker speaker. In addition, for correct trials, the call sign was classified at chance performance (Fig. 3a). However, for incorrect trials the classifier detected the masker call sign significantly more often than the target call sign (34.1%, $P < 10 \times 10^{-4}$, t -test; Fig. 3b), which again shows errors due to an early selection of the masker (incorrect) speaker.

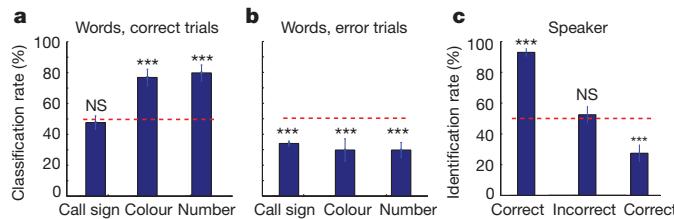


Figure 3 | Decoding spoken words and the identity of the attended speaker. **a**, Classification rate and standard deviation for spoken words (call sign, colour and number) of the attended speaker from the neural responses to the 28 mixtures. Classifiers were trained on single speaker examples only. Colour and number of the attended speech are decoded with high accuracy (77.2% and 80.2%, $P < 10 \times 10^{-4}$, t -test) in correct trials, but not the call sign (48.0%, not significant (NS), t -test). **b**, In error trials, the classifier showed a systematic bias towards the words of the masker speaker (34.1%, 30.0%, 30.1%, $P < 10 \times 10^{-4}$, t -test). **c**, Attended speaker identification rate and standard deviation in correct for target, incorrect (for both target and masker), and correct for masker trials.

For the speaker identification analyses, we divided the behavioural error types into two subsets. The first type occurred when the reported colour-number combination was incorrect for either speaker ('incorrect'; 16.5% of trials). The second type occurred when subjects reported the correct colour-number for the masker instead of the target speaker ('correct for masker'; 8.6% of trials).

In correct trials, the classifier identified the target speaker 93.0% of the time ($P < 10 \times 10^{-4}$, t -test; Fig. 3c). During incorrect trials, the classifier performance was at chance. However, during correct for masker trials, the classifier identified the masker rather than the target speaker (27.3%; $P < 10 \times 10^{-4}$, t -test; Fig. 3c). These classification results confirm the observed restoration seen in spectrotemporal reconstruction, without necessarily assuming a linear relationship between the neural responses and the stimulus. Furthermore, they extend recent findings using similar methods to decode speech sounds presented in isolation²² to full words and sentences under complex listening conditions.

We next asked whether the observed robust encoding of attended speech results as an emergent property of the distributed population activity or is driven by a few spatially discrete sites. The cortical regions with reliable evoked responses to speech stimuli were found using a t -test between neural responses during speech and silence ($P < 0.01$), and were confined to the posterior superior and middle temporal gyri (Fig. 4a). An example of the attentional response modulation at a single electrode is shown in Fig. 4b–d. The spectrotemporal receptive field (STRF, estimated using the <http://www.strflab.berkeley.edu> package) of this electrode in passive listening to speech (TIMIT²⁰) showed a strong preference for high frequency sounds (Fig. 4b) (STRFs for all electrodes of one subject are provided in Supplementary Fig. 2b). This tuning was also evident in the increased neural response at this electrode (Fig. 4d, dashed lines) to each of the single speakers' high frequency sound components (circled in Fig. 4c, responses are delayed about 120 ms from the stimulus). However, the responses to the same speech mixture sound (Fig. 4d, solid lines) were significantly modulated by attention. The responses to high frequency components were enhanced for the attended speaker, but suppressed for similar sounds in the masker speaker (Fig. 4d, solid lines compared to dashed lines). This highly modulated yet fixed feature selectivity probably contributes to the constancy of the single speaker representation observed in our previous analyses. To quantify this effect for each individual electrode, we measured the correlation between the neural responses to the attended mixture and to those of the speakers in isolation (AMI_{elec} , equation (2) in Methods). We found a varying degree of bias towards the attended speaker distributed across the population (Supplementary Fig. 3d; $AMI_{elec} = 0.28$ for the example

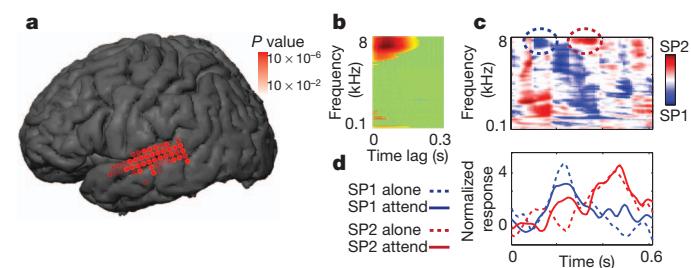


Figure 4 | Attentional modulation of individual electrode sites. **a**, Electrodes picking up a significant difference between responses to silence and speech sounds ($P < 0.01$, t -test). **b**, STRF of this representative electrode site shows a preference for high frequency sounds. **c**, Mixture difference spectrogram for a selected duration containing a high frequency component for each speaker (circled). **d**, The electrode shows an increased response to high frequency sounds of single speakers (dashed lines, peak neural response is delayed by about 120 ms). However, the neural response to the same mixture sound in two attention conditions (solid lines) showed an enhanced response to high frequency sounds only for the target, but with responses for similar sounds in the masker speaker suppressed.

in Fig. 4), which gradually builds up after the end of the call sign (Supplementary Fig. 3e). We did not observe any particular anatomical pattern for the attentional modulation across sites (Supplementary Fig. 3f). Rather, it appeared to be distributed over responsive sites, consistent with previous findings of higher-order sound processing²³.

In summary, we demonstrate that the human auditory system restores the representation of the attended speaker while suppressing irrelevant competing speech. Speech restoration occurs at a level where neural responses still show precise phase-locking to spectrotemporal features of speech. Population responses revealed the emergent representation of speech extracted from a mixture, including the moment-by-moment allocation of attentional focus.

These results have implications for models of auditory scene analysis. In agreement with recent studies, the cortical representation of speech in the posterior temporal lobe does not merely reflect the acoustical properties of the stimulus, but instead relates strongly to the perceived aspects of speech¹⁰. Although the exact mechanisms are not fully known, multiple processes in addition to attention are likely to enable this high-order auditory processing, including grouping of predictable regularities in speech acoustics²⁴, feature binding^{3,25} and phonemic restoration²⁶. Conversely, behavioural errors seem to result from degradation of the neural representation, a direct result of inherent sensory interference such as energetic masking¹⁶ (Supplementary Fig. 3g, h) and/or the allocation of attention²⁷.

In speech, the end result represented in the posterior temporal lobe appears to be unaffected by perceptually irrelevant sounds, which is ideal for subsequent linguistic and cognitive processing. Following one speaker in the presence of another can be trivial for a normal human listener, but remains a major challenge for state-of-the-art automatic speech recognition algorithms²⁸. Understanding how the brain solves this problem may inspire more efficient and generalizable solutions than current engineering approaches²⁹. It will also shed light on how these processes become impaired during ageing and in disorders of speech perception in real-world hearing conditions⁷.

METHODS SUMMARY

Three human subjects with normal hearing underwent the placement of a subdural electrode array as part of their clinical treatment for epilepsy. We used speech samples from a publicly available database called Coordinate Response Measure (CRM¹⁵). One male and one female speaker were selected with two call signs (ringo and tiger), three colours (red, blue or green) and three numbers (two, five or seven). We generated 12 unique combinations of call sign, colour and number per speaker (total of 24 single speaker phrases) and 28 mixture speech samples by selecting from combinations of the 24 single speaker sentences (0 dB target-to-masker ratio). Speech sounds were presented monaurally from a loud speaker. We used stimulus reconstruction^{17–19} to map the population electrocorticographic response to the spectrogram of the speech stimulus. Reconstruction filters were estimated from neural responses to a separate speech corpus (TIMIT²⁰). Test speakers were not used in the estimation of filters. For word and speaker decoding analysis, a regularized linear classifier²¹ was trained on neural responses of the single speakers and then used to decode the spoken words and speaker identity of the attended speech mixture.

Full Methods and any associated references are available in the online version of the paper at www.nature.com/nature.

Received 30 August 2011; accepted 5 March 2012.

Published online 18 April 2012.

- Cherry, E. C. Some experiments on the recognition of speech, with one and with two ears. *J. Acoust. Soc. Am.* **25**, 975–979 (1953).
- Shinn-Cunningham, B. G. Object-based auditory and visual attention. *Trends Cogn. Sci.* **12**, 182–186 (2008).

- Bregman, A. S. *Auditory Scene Analysis: The Perceptual Organization of Sound* (MIT Press, 1994).
- Kerlin, J., Shahin, A. & Miller, L. Attentional gain control of ongoing cortical speech representations in a “cocktail party”. *J. Neurosci.* **30**, 620–628 (2010).
- Besle, J. et al. Tuning of the human neocortex to the temporal dynamics of attended events. *J. Neurosci.* **31**, 3176–3185 (2011).
- Bee, M. & Micheyl, C. The cocktail party problem: what is it? How can it be solved? And why should animal behaviorists study it? *J. Comparative Psychol.* **122**, 235–252 (2008).
- Shinn-Cunningham, B. G. & Best, V. Selective attention in normal and impaired hearing. *Trends Amplif.* **12**, 283–299 (2008).
- Scott, S. K., Rosen, S., Beaman, C. P., Davis, J. P. & Wise, R. J. S. The neural processing of masked speech: evidence for different mechanisms in the left and right temporal lobes. *J. Acoust. Soc. Am.* **125**, 1737–1743 (2009).
- Elhilali, M., Xiang, J., Shamma, S. A. & Simon, J. Z. Interaction between attention and bottom-up saliency mediates the representation of foreground and background in an auditory scene. *PLoS Biol.* **7**, e1000129 (2009).
- Chang, E. F. et al. Categorical speech representation in human superior temporal gyrus. *Nature Neurosci.* **13**, 1428–1432 (2010).
- Crone, N. E., Boatman, D., Gordon, B. & Hao, L. Induced electrocorticographic gamma activity during auditory perception. *Clin. Neurophysiol.* **112**, 565–582 (2001).
- Steinschneider, M., Fishman, Y. I. & Arezzo, J. C. Spectrotemporal analysis of evoked and induced electroencephalographic responses in primary auditory cortex (A1) of the awake monkey. *Cereb. Cortex* **18**, 610–625 (2008).
- Scott, S. K. & Johnsrude, I. S. The neuroanatomical and functional organization of speech perception. *Trends Neurosci.* **26**, 100–107 (2003).
- Hackett, T. A. Information flow in the auditory cortical network. *Hear. Res.* **271**, 133–146 (2011).
- Bolia, R. S., Nelson, W. T., Ericson, M. A. & Simpson, B. D. A speech corpus for multitalker communications research. *J. Acoust. Soc. Am.* **107**, 1065–1066 (2000).
- Brungart, D. S. Informational and energetic masking effects in the perception of two simultaneous talkers. *J. Acoust. Soc. Am.* **109**, 1101–1109 (2001).
- Mesgarani, N., David, S. V., Fritz, J. B. & Shamma, S. A. Influence of context and behavior on stimulus reconstruction from neural activity in primary auditory cortex. *J. Neurophysiol.* **102**, 3329–3339 (2009).
- Bialek, W., Rieke, F., de Ruyter van Steveninck, R. R. & Warland, D. Reading a neural code. *Science* **252**, 1854–1857 (1991).
- Pasley, B. N. et al. Reconstructing speech from human auditory cortex. *PLoS Biol.* **10**, e1001251 (2012).
- Garofolo, J. S. et al. *TIMIT Acoustic-Phonetic Continuous Speech Corpus* (Linguistic Data Consortium, 1993).
- Rifkin, R., Yeo, G. & Poggio, T. Regularized least-squares classification. *Nato Science Series Sub Series III Computer and Systems Sciences* **190**, 131–154 (2003).
- Formisano, E., De Martino, F., Bonte, M. & Goebel, R. “Who” is saying “what”? Brain-based decoding of human voice and speech. *Science* **322**, 970–973 (2008).
- Staeren, N., Renvall, H., De Martino, F., Goebel, R. & Formisano, E. Sound categories are represented as distributed patterns in the human auditory cortex. *Curr. Biol.* **19**, 498–502 (2009).
- Shamma, S. A., Elhilali, M. & Micheyl, C. Temporal coherence and attention in auditory scene analysis. *Trends Neurosci.* **34**, 114–123 (2010).
- Darwin, C. J. Auditory grouping. *Trends Cogn. Sci.* **1**, 327–333 (1997).
- Warren, R. M. Perceptual restoration of missing speech sounds. *Science* **167**, 392–393 (1970).
- Kidd, G. Jr, Arbogast, T. L., Mason, C. R. & Gallun, F. J. The advantage of knowing where to listen. *J. Acoust. Soc. Am.* **118**, 3804–3815 (2005).
- Shen, W., Olive, J. & Jones, D. Two protocols comparing human and machine phonetic discrimination performance in conversational speech. *INTERSPEECH* 1630–1633 (2008).
- Cooke, M., Hershey, J. R. & Rennie, S. J. Monaural speech separation and recognition challenge. *Comput. Speech Lang.* **24**, 1–15 (2010).

Supplementary Information is linked to the online version of the paper at www.nature.com/nature.

Acknowledgements The authors would like to thank A. Ren for technical help, and C. Micheyl, S. Shamma and C. Schreiner for critical discussion and reading of the manuscript. E.F.C. was funded by National Institutes of Health grants R00-NS065120, DP2-OD00862, R01-DC012379, and the Ester A. and Joseph Klingenstein Foundation.

Author Contributions N.M. and E.F.C. designed the experiment, collected the data, evaluated results and wrote the manuscript.

Author Information Reprints and permissions information is available at www.nature.com/reprints. The authors declare no competing financial interests. Readers are welcome to comment on the online version of this article at www.nature.com/nature. Correspondence and requests for materials should be addressed to E.F.C. (changed@neurosurgeon.ucsf.edu).

METHODS

The experimental protocol was approved by the Committee for Human Research at the University of California, San Francisco.

Subjects. Three human subjects underwent the placement of a high-density subdural electrode array (4 mm pitch) over the language-dominant hemisphere as part of routine clinical treatment for epilepsy. Subjects gave their written informed consent before surgery. All subjects had self-reported normal hearing and underwent neuropsychological language testing (including the Boston naming and verbal fluency tests) and were found to be normal. The intracarotid sodium amobarbital (Wada) test was used for language dominance assessment. The electrodes in the study were located over the posterior dorsolateral temporal lobe. The location and corresponding spectrotemporal receptive fields of all the included electrodes for a subject are shown in Supplementary Fig. 2.

Data acquisition and pre-processing. The electrocorticography signal was recorded with a multichannel amplifier optically connected to a digital signal processor (TuckerDavis Technologies). Each channel time series was visually and quantitatively inspected for artefacts or excessive noise. The data were then segmented with a 100 ms pre-stimulus baseline and a 400 ms post-stimulus interval. The common mode signal was estimated using principal component analysis with channels as repetitions and was removed from each channel time series using vector projection.

Task design and behavioural testing. We used speech samples from a publicly available database called Coordinate Response Measure (CRM¹⁵) containing sentences in the form “ready (call sign) go to (colour) (number) now”. One male and one female speaker (speakers one and five in CRM corpus) were selected with two call signs (ringo and tiger), three colours (blue (B), red (R) or green (G)) and three numbers (two, five or seven). For each of the two call signs, we generated six colour-number combinations (B2, B5, R2, R7, G5, G7), resulting in 12 different phrases. We chose the same phrases for each of the two speakers, resulting in 24 single speaker sentences. We then produced 28 unique mixture speech samples by selecting from combinations of the 24 single speaker sentences at 0 dB target-to-masker ratio. Each mixture sample was chosen such that there was no overlap between call signs, colours or the numbers of the two phrases. In addition, each speaker had the same number of call signs (ringo or tiger) in each trial block. The sounds were presented monaurally from a loudspeaker connected to a laptop, which was also used to collect subjects’ responses through a customized graphical user interface. Each trial block consisted of 28 trials and the target call sign was fixed for each block. The target call sign was displayed visually before and during the trial block. Subjects first listened to each of the speakers alone and were able to report the colour and number with 100% accuracy. Subjects then listened to a monaural, simultaneous mixture of the two speakers’ phrases with different call signs, colours and numbers. The subjects were instructed to respond by indicating the colour and number spoken by the talker who uttered the target call sign. The target speaker changed from trial to trial pseudo-randomly, requiring the subjects to initially monitor both speakers until they detect the target call sign. After each trial block, the target call sign was changed, switching the role of target and masker speakers in each mixture sound.

Electrode selection. The cortical sites on the superior and middle temporal gyri with reliable evoked responses to speech stimuli were selected for all the subsequent analysis. Our inclusion criteria consisted of a *t*-test between responses to randomly selected time frames during passive speech presentation (TIMIT) and in silence ($P < 0.01$, resulting in 83, 92 and 102 electrodes for subjects one to three. One example subject is shown in Supplementary Fig. 2a). Solely for visualization, we also estimated the STRFs of these selected sites from passive

listening to TIMIT using normalized reverse correlation algorithm (STRFLab software package, <http://www.strflab.berkeley.edu>; Supplementary Fig. 2b). Correlation histogram of STRF predictions for all 275 electrode sites is shown in Supplementary Fig. 1c.

Stimulus reconstruction. We used stimulus reconstruction to map the population neural responses to the spectrogram of the speech stimulus^{17–19}. Reconstruction filters were estimated from neural responses to a separate speech corpus (TIMIT²⁰) containing a total of 499 unique short sentences from 402 different speakers. Filters were obtained using normalized reverse correlation to minimize the mean squared error of the reconstructed spectrograms¹⁷ with filter time lags from −420 to 0 ms (causal filters). The filters were then fixed in all subsequent conditions and were applied to the neural responses to CRM samples. Neither of the speakers or phrases in the CRM data set was used in estimation of the filters. The output of the reconstruction algorithm was further processed with a band-pass filter applied to each frequency channel of reconstructed spectrograms to remove the baseline. All the processing steps for stimulus reconstruction were identical in all conditions (single and mixture speakers).

AMI. To quantify the change in similarity between the representation of single and attended speaker in mixture speech, we defined the AMI_{spec} in equation (1). The stereotypical format of the CRM phrases results in an intrinsic correlation between the neural responses to different sentences, particularly at the beginning (“ready”) and middle of the carrier phrase (“go to”), which results in reduced possible AMI_{spec} values for these segments. To estimate an upper bound for unbiased comparison, AMI_{spec} was calculated where the representation of an attended speaker in a mixture is ideally assumed to be identical to the representation of that speaker when presented alone; therefore, replacing $\text{SP}_{\text{attend}}$ in equation (1) with the reconstructed spectrogram of single speaker SP_{alone} . The upper bound peaks at the call sign, colour and number where different phrases are most dissimilar. The overall increase in the upper bound is due to the progressive asynchrony between the two speakers.

The same statistics can be used to estimate the AMI of an individual electrode site by calculating the correlation values between the neural response of that site to attended mixture and single speaker presentations:

$$\text{AMI}_{\text{elec}} = \frac{\text{Corr}(\text{R-SP1}_{\text{alone}}, \text{R-SP1}_{\text{attend}}) - \text{Corr}(\text{R-SP1}_{\text{alone}}, \text{R-SP2}_{\text{attend}})}{\text{Corr}(\text{R-SP2}_{\text{alone}}, \text{R-SP2}_{\text{attend}}) - \text{Corr}(\text{R-SP2}_{\text{alone}}, \text{R-SP1}_{\text{attend}})} \quad (2)$$

where $\text{R-SP1}_{\text{alone}}$ and $\text{R-SP2}_{\text{alone}}$ are the responses of an electrode to speakers one and two alone, respectively, and $\text{R-SP1}_{\text{attend}}$ and $\text{R-SP2}_{\text{attend}}$ are the responses of the same electrode to the mixture of the two when the attended target is speaker one and two, respectively.

Classification of spoken words and speaker identity. A linear-frame-based regularized-least-square classifier²¹ was used to investigate the discriminability of the spoken words and speaker identity from electrocorticographic responses. Two binary classifiers were trained to classify the call sign and speaker identity, and two separate three-way classifiers were used for colour and for number classification. Classifiers were trained only on the neural responses of single speakers (24 sentences) and tested on the mixtures. The classifiers produced a linear weighted sum of the neural responses at each time instance and the classifier that produced the maximum average output over the duration of words was chosen as classification result. The classifier decision was limited to only the colours and numbers that occurred in each mixture, therefore resulting in same 50% chance performance in all cases.