

a specific network in which the parietal and perhaps lateral frontal cortices appear to be optimally situated to mediate the integration and attentional selection of motion information across modalities. In audiovisual face perception, crossmodal attention influences crossmodal binding during speech reading, attention and audiovisual integration interact with each other in a sophisticated manner. However, feature-selective attention in audiovisual conditions and the relationship between feature-selective attention and high-level audiovisual semantic integration remain to be explored.

In a single (visual or auditory) modality, feature-selective attention may lead to selective processing of the attended features of an object in the brain¹⁷. Nobre *et al.* demonstrated that ERPs are modulated by feature-selective attention. a4561 0 9 155.9.5(h)3.61al 22 Tm 2(t)-

Figure 1. (A) Four examples of audiovisual stimuli; the red numbers indicate runs with the number task only. (B) Time course of a trial for the runs with the number task, in which the stimuli included randomly presented numbers and videos/audios/movie clips. (C) Time course of a trial for the runs with the gender, emotion, or bi-feature task. For both (B, C), the presentation of a stimulus (video/audio/movie clip) lasted 1,400 ms and was repeated four times during the first eight seconds in a trial. A visual cue (“+”) appeared at the 8th second and persisted for six seconds.

For each of the three runs with the number task, in addition to the corresponding audiovisual, visual-only, or auditory-only facial stimuli from the movie clips, numbers in red appeared sequentially at the center of the screen (see Fig. 1A). The subject’s task was to attend to the numbers instead of the other stimuli (see Table 1). We designed a difficult number task for the subjects in which they were asked to find and count the repeated numbers to ensure that they fully ignored the features of the visual-only, auditory-only, or audiovisual facial stimuli. Therefore, the subjects performed this task with low accuracy, as shown in Fig. S3. At the beginning of each block, there were four seconds before the first trial, and a short instruction in Chinese (see Table 1) was displayed on the screen in the first two seconds (the last two seconds were used to display numbers, as indicated below). At the beginning of each trial, a visual-only, auditory-only or audiovisual facial stimulus was presented to the subject for 1,400 ms, followed by a 600-ms blank period. This two-second cycle with the same stimulus was repeated four times, followed by a six-second blank period. Therefore, one trial lasted 34 seconds. In addition to the above stimuli, eight numbers in red appeared one by one at the center of the screen, each a random integer from 0 to 9. Each number lasted 900 ms and the interval between two subsequent numbers was 350 ms. The first number appeared 2 seconds before the beginning of this trial. The subjects were asked to find and count the repeated numbers. After the stimulation, a

xation cross appeared on the screen. The subjects then responded by pressing the right-hand keys according to the instruction for this block (see Table 1). The fixation cross changed color at the 12th second, indicating that the next trial would begin shortly (see Fig. 1B). In total, a run lasted 350 s.

The procedure for the three runs with the gender/emotion task was similar to that for the runs with the number task, except that no numbers appeared on the screen and the subjects performed a gender/emotion judgment task (See Table 1). Specifically, the subjects were asked to focus their attention on either the gender or the emotion of the presented stimuli (visual-only, auditory-only, or audiovisual facial stimuli; see Fig. 1A without regard to the numbers) and make a corresponding judgment (male vs. female for the gender task or crying vs. laughing for the emotion task) to each stimulus. At the beginning of each block, a short instruction (see Table 1) was displayed for four seconds on the screen. The time course of each trial was similar to that in the runs with number task (see Fig. 1C). In each trial, the subject was asked to judge the gender/emotion category of the stimulus and press the right-hand keys according to the instruction for this block.

For the three runs with the bi-feature task, the subjects were asked to simultaneously attend to both gender and emotion features (see Table 1). The experimental procedure for each run was similar to that for the runs with the gender/emotion task with the following3(lo)-c1i0364638 21e13.1((a)3.3(s).5(w)-2.6(91(o)12(r to)-c1i03Bh)3.5

voxels, time series detrending, and normalization of the time series in each block to zero mean and unit variance. All preprocessing steps were performed using SPM8 custom functions in MATLAB 7.4 (MathWorks, Natick, Massachusetts, USA).

Univariate GLM analysis. This experiment included four experimental tasks (number, gender, emotion, and bi-feature). For each experimental task, three runs corresponding to the visual-only, the auditory-only, and the audiovisual stimulus conditions were performed. To confirm that audiovisual sensory integration occurred for each experimental task and determine the heteromodal areas associated with audiovisual integration, we performed voxel-wise group analysis of the fMRI data based on a mixed-effect two-level GLM in SPM8. In particular, using the data from the three number runs, we performed GLM analysis to explore the audiovisual integration at the sensory level when the subjects fully ignored the visual-only, auditory-only, or audiovisual facial stimuli while only attending to the numbers. The GLM analysis included the following data processing. The fMRI data for each subject were subjected to a first-level GLM, and the estimated beta coefficients across all subjects were then combined and analyzed using a second-level GLM. The following statistical criterion was used to determine brain areas for audiovisual sensory integration: $\max(A, V)$ ($p < 0.05$, FWE-corrected) $\cap [V > 0$ or $A > 0$ ($p < 0.05$, uncorrected)]^{1,6,24-27}, where \cap denotes the intersection of two sets. For each subject, each task, and each stimulus condition, we also computed the percent signal changes of the pSTS/MTG clusters via region-of-interest (ROI)-based analysis (implemented by the MATLAB toolbox MarsBaR). Specifically, we identified the clusters consisting of significantly activated voxels in the bilateral pSTS/MTG via group GLM

where θ_{ij} is the angle between two pattern vectors

differentiated for different experimental tasks or different semantic features. Thus, audiovisual sensory integration rather than audiovisual semantic integration occurred in the identified heteromodal areas of the pSTS/MTG, consistent with previous results

3. $\frac{R_{AV} - \max(R_{VO}, R_{AO})}{R_{AV}}$ (Reproducibility ratio in the audiovisual condition minus the maximum of the reproducibility ratios in the visual-only and auditory-only conditions). Left/Right: gender/emotion categories; the first 3 rows: audiovisual, visual-only, and auditory-only stimulus conditions, respectively; the 4th row: the reproducibility ratio in the audiovisual condition minus the maximum of the reproducibility ratios in the visual-only and auditory-only conditions.

$p < 10^{-17}$, $F(3, 8) = 68.26$) (Fig. 3A–C,E–G). There was also a significant interaction effect between the two factors of stimulus condition and experimental task (gender categories: $p < 10^{-17}$, $F(6, 8) = 30.07$; emotion categories: $p < 10^{-8}$, $F(6, 8) = 10.05$). Post hoc Bonferroni-corrected paired t-tests on the stimulus conditions revealed the following: (i) for each task-relevant feature (gender categories with the gender or the bi-feature task, left panel of Fig. 3; emotion categories with the emotion or the bi-feature task, right panel of Fig. 3), the reproducibility ratios were significantly higher for the audiovisual stimulus condition than for the visual- or auditory-only stimulus condition (all $p < 0.001$ corrected); and (ii) for each task-irrelevant feature (gender categories with the number or the emotion task, left panel of Fig. 3; emotion categories with the number or the gender task, right panel of Fig. 3), there were no significant differences between the audiovisual and the visual-only or auditory-only stimulus condition (all $p > 0.05$). Furthermore, post hoc Bonferroni-corrected paired t-tests on the experimental tasks revealed that (i) in each of the audiovisual, visual-only and auditory-only stimulus conditions, the reproducibility ratios for gender/emotion categories were significantly higher for each relevant task (gender categories: the gender or the bi-feature task, left panel of Fig. 3; emotion categories: the emotion or the bi-feature task, right panel of Fig. 3) than for each irrelevant task (gender categories: the number or the emotion task, left panel of Fig. 3; emotion categories: the number or the gender task, right panel of Fig. 3) (all $p < 0.05$, corrected) and that (ii) in each of the audiovisual, visual-only and auditory-only stimulus conditions, there were no significant differences in the reproducibility ratios.

($p < 10^{-9}$, $F(2, 8) = 36.97$ for gender categories; $p < 10^{-11}$, $F(2, 8) = 46.13$ for emotion categories). Furthermore, post hoc Bonferroni-corrected paired t-tests demonstrated that the cross-reproducibility ratios were significantly higher for the relevant task than for the irrelevant tasks (gender categories: $p < 0.001$ corrected, $t(8) = 16.23$ for gender task vs. number task; $p < 0.001$ corrected, $t(8) = 15.49$ for gender task vs. emotion task; emotion categories: $p < 0.001$ corrected, $t(8) = 16.05$ for emotion task vs. number task; $p < 0.001$ corrected, $t(8) = 14.36$ for emotion

the group level (see Materials and Methods). As shown in Fig. 5, there were more functional connections from the heteromodal areas to the brain areas encoding the gender/emotion feature (Table 2/Table 3) for the relevant task (gender/emotion task) than for the irrelevant tasks (number and emotion/gender tasks). We thus observed that in the audiovisual condition, feature-selective attention enhanced the functional connectivity and thus regulated the information flows from the heteromodal areas to the brain areas encoding the attended feature. Furthermore,

results were still obtained. Second, only visual-only, auditory-only and audiovisual facial stimuli were considered in this study. In the future, we must simplify our experimental design, increase the number of subjects, and further consider non-facial stimuli to extend our conclusions.

References

- Calvert, G. A. & Pesenti, T. Multisensory integration: methodological approaches and emerging principles in the human brain. *J. Physiol. Paris* , 191–205 (2004).
- Campanella, S. & Belin, P. Integrating face and voice in person perception. *Trends Cogn. Sci.* 535–543 (2007).
- Schweinberger, S., Ebner, F. F., Obertson, D. & Kaufmann, J. M. Hearing facial identities. *Q. J. Exp. Psych.* 1446–1456 (2007).
- Bushara, C. O. et al. Neural correlates of cross-modal binding. *Nat. Neurosci.* 190–195 (2003).
- Macaluso, E., Frith, C. D. & Driver, J. Multisensory stimulation with or without saccades: fMRI evidence for crossmodal effects on sensory-specific cortices that reflect multisensory location-congruence rather than task-relevance. *NeuroImage* 144–145 (2005).
- Macaluso, E., George, N., Dolan, R. J., Spence, C. & Driver, J. Spatial and temporal factors during processing of audiovisual speech: PET study. *NeuroImage* 21, 725–732 (2004).
- McClurkin, J. W. & Optican, L. M. Primate striate and prestriate cortical neurons during discrimination. I. Simultaneous temporal encoding of information about color and pattern. *J. Neurophysiol.* 481–495 (1996).
- Nobre, A. C., Pessoa, A. & Chelazzi, L. Selective attention to specific features within objects: Behavioral and electrophysiological evidence. *J. Cognitive Neurosci.* 18, 539–561 (2006).
- Woodman, G. F. & Vogel, E. K. Selective storage and maintenance of an object's features in visual working memory. *Psychon. B. Rev.* 15, 223–229 (2008).
- Taylor, R. I., Moss, H. E., Stamataki, E. A. & Tyler, L. W. Binding crossmodal object features in perirhinal cortex. *Proc. Natl. Acad. Sci. U.S.A.* 103, 8239–8244 (2006).
- Talsma, D., Senkowski, D., Soto-Faraco, S. & Woldorff, M. G. The multifaceted interplay between attention and multisensory integration. *Trends Cogn. Sci.* 14, 400–410 (2010).
- Lewis, J. W., Beauchamp, M. S. & DeYoe, E. A. A comparison of visual and auditory motion processing in human cerebral cortex. *Cereb. Cortex* 10, 873–888 (2000).
- Joassin, F. et al. Cross-modal interactions between human faces and voices involved in person recognition. *Cognition* , 367–376 (2011).
- Saito, D. et al. Cross-modal binding and activated attentional networks during audio-visual speech integration: a functional MRI study. *Cereb. Cortex* 15, 1750–1760 (2005).
- Ahveninen, J. et al. Task-modulated “what” and “where” pathways in human auditory cortex. *Proc. Natl. Acad. Sci. U.S.A.* 103, 14608–14613 (2006).
- Maunsell, J. H. R. & Hochstein, S. Effects of behavioral state on the stimulus selectivity of neurons in area V4 of the macaque monkey. In: *Channels in the visual nervous system: neurophysiology, psychophysics and models*, (ed. Blum B), 447–470. London: Freeman (1991).
- Mirabella, G. et al. Neurons in area V4 of the macaque translate attended visual features into behaviorally relevant categories. *Neuron* 52, 303–318 (2007).
- Jeong, J. W. et al. Congruence of happy and sad emotion in music and faces modulates cortical audiovisual activation. *NeuroImage* 54, 2973–2982 (2011).
- Reifels, B., Ethofer, T., Grodd, W., Erb, M. & Wildgruber, D. Audiovisual integration of emotional signals in voice and face: an event-related fMRI study. *NeuroImage* 36, 1445–1456 (2007).
- Müller, V. I., Ciesli, E. C., Turetsky, B. I. & Eichler, S. B. Crossmodal interactions in audiovisual emotion processing. *NeuroImage* 54, 553–561 (2011).
- Müller, V. I. et al. Incongruence effects in crossmodal emotional integration. *NeuroImage* 54, 2257–2266 (2011).
- Li, Y. et al. Crossmodal Integration Enhances Neural Representation of Task-Relevant Features in Audiovisual Face Perception. *Cereb. Cortex* 25, 384–395 (2015).
- Friston, K. J. et al. Statistical parametric maps in functional imaging: a general linear approach. *Hum. Brain Mapp.* 2, 189–210 (1994).
- Calvert, G. A., Campbell, K. L. & Brammer, M. J. Evidence from functional magnetic resonance imaging of crossmodal binding in the human heteromodal cortex. *Curr. Biol.* 10, 649–657 (2000).
- Frassinetti, F., Bolognini, N. & Ladavas, E. Enhancement of visual perception by crossmodal visuo-auditory interaction. *Exp. Brain Res.* 171, 332–343 (2002).
- Macaluso, E. & Driver, J. Multisensory spatial interactions: a window onto functional integration in the human brain. *TRENDS Neurosci.* 28, 264–271 (2005).
- Beauchamp, M. S. Statistical criteria in fMRI studies of multisensory integration. *Neuroinformatics* 3, 105–113 (2005).
- Brett, M., Anton, J.-L., Valabregue, K. & Poline, J.-B. Region of interest analysis using the MarsBar toolbox for SPM 99. *NeuroImage* 18, 1140–1141 (2002).
- DiGirolamo, N., Goebel, R. & Bandettini, P. Information-based functional brain mapping. *Proc. Natl. Acad. Sci. U.S.A.* 103, 3863–3868 (2006).
- Nichols, T. & Hayasaka, S. Controlling the familywise error rate in functional neuroimaging: a comparative review. *Stat. Methods Med. Res.* 20, 125–142 (2006).

Acknowledgements

This work was supported by the National Key Basic Research Program of China (973 Program) under Grant 2015CB351703, the National High-tech R&D Program of China (863 Program) under Grant 2012AA011601, the National Natural Science Foundation of China under Grants 91420302, 81471654 and 61403147, and Guangdong Natural Science Foundation under Grant 2014A030312005.

Author Contributions

Y.L. designed research and wrote the paper; J.L. and W.W. analyzed the data; B.H., T.Y. and P.L. performed the research; F.F. and P.S. revised the paper; all authors reviewed the manuscript.

Additional Information

Supplementary Information accompanies this paper at <http://www.nature.com/srep>

The authors declare no competing financial interests.

Correspondence and requests for materials should be addressed to Li, Y. *et al.* Selective Audiovisual Semantic Integration Enabled by Feature-Selective Attention. *Sci. Rep.* 6, 18914; doi: 10.1038/srep18914 (2016).

