# SUPPLEMENTARY MATERIAL: PART I

### MATHEMATICAL DESCRIPTION OF THE NEURAL NETWORK

### **1** Basal structure of the network

The neural network model consists of two chains of N unisensory neurons (Fig. 1, upper panel). Each neuron codes for a particular spatial position in its modality. Moreover, each chain is topologically organized, i.e., proximal neurons code for proximal positions. In the following, we will denote with a first subscript the particular area (auditory or visual) and with a second subscript, after a comma, the neuron position within the area.

Each neuron receives three different kinds of inputs: a sensory input from the environment (say u), a lateral input from neurons of the same modality (say l) and a cross-modal input from neurons of the other modality (say c). The global input (equal to the sum of the previous three contributions) is then passed through a sigmoidal relationship,  $\varphi(\cdot)$ , which accounts for the presence of a lower threshold and upper saturation in neuron activity, and a first-order low-pass filter with time constant  $\tau$ , which accounts for the neuron integrative capacity.

Hence, for the generic k-th neuron in the modality S (S = A or V for the auditory and visual modalities, respectively) we can write

$$\tau \frac{dy_{S,k}}{dt} = -y_{S,k} + \phi \left( u_{S,k} + l_{S,k} + c_{S,k} \right)$$
(S1)

Where  $y_{S,k}$  represents the neuron output, and the sigmoidal relationship is described by the following equation

$$\phi(x) = \frac{1}{1 + \exp(-s(x - x_0))}$$
(S2)

*s* and  $x_0$  are parameters, which set the slope and the position of the sigmoidal relationship. According to Eq. (2), the neuron output activity is normalized between 0 and 1 (zero means a silent neuron, one a maximally activated neuron).

It is worth noting that, for the sake of simplicity, we used the same parameters ( $\tau$ , *s* and *x*<sub> $\theta$ </sub>) for all neurons independently of their modality. This choice was adopted to minimize the number of model assumptions.

The expression for the sensory input is computed as the scalar product of the sensory representation of the stimulus (i.e., the vector  $I_s = [i_{s,1} i_{s,2} \cdots i_{s,k} \cdots i_{s,N}]^T$ ) and the neuron receptive field (i.e., the

vector 
$$R_{s,k} = [r_{S,k1} r_{S,k2} \cdots r_{S,kj} \cdots r_{S,kN}]^T$$
):

$$u_{S,k} = \left\langle R_{s,k}, I_s \right\rangle = \sum_{j=1}^{N} r_{S,kj} i_{S,j}$$
(S3)

We assumed that the neuron receptive field,  $R_{s,k}$ , has initially a large extension, described with a Gaussian function, and then progressively shrinks during training, to fit the width of the external input (see section "Training the network").

The lateral input is computed as follows

$$l_{S,k} = \sum_{j=1}^{N} \lambda_{kj} y_{S,j}$$
(S4)

where  $\lambda_{kj}$  represents a lateral intra-area synapse connecting the presynaptic neuron *j* to the post synaptic neuron *k* in the same area. Here we used the classical Mexican-hat arrangement: a neuron is excited by proximal neurons in the same area, and inhibited by more distal ones

$$\lambda_{kj} = \lambda_{ex} \exp\left(-\frac{d(\theta_j, \theta_k)^2}{2\sigma_{ex}^2}\right) - \lambda_{in} \exp\left(-\frac{d(\theta_j, \theta_k)^2}{2\sigma_{in}^2}\right)$$
(S5)

where  $\lambda_{ex}, \lambda_{in}, \sigma_{ex}, \sigma_{in}$  are parameters which set the strength and width of the excitatory and inhibitory portions of the Mexican hat. In particular, we have  $\lambda_{ex} > \lambda_{in}$  and  $\sigma_{ex} < \sigma_{in}$ . Moreover,  $d(\theta_j, \theta_k)$ represents the distance between neurons' preferred positions, i.e.

$$d(\theta_j, \theta_k) = \left|\theta_j - \theta_k\right| \tag{S6}$$

It is worth noting that we used the same expression of lateral synapses (Eq. S5) in both the auditory and visual areas, to limit the number of model assumptions.

Finally, the cross-modal term in Eq. (1) is computed as the convolution of the vector of cross modal synapses and the activity in the other unisensory area, i.e.

$$c_{S,k} = \sum_{j=1}^{N} w_{SQ,kj} y_{Q,j} \quad \text{with} \quad S = A \text{ or } V \quad Q = A \text{ or } V \text{ with } S \neq Q$$
(S7)

where  $w_{SQ,kj}$  represents a cross-modal synapse from the pre-synaptic neuron j in the area Q to the post-synaptic neuron k in the area S. We assumed that the cross-modal synapses are initially ineffective and are progressively reinforced during the training phase.

### 2 Training the network

Starting from the initial basal value of synapses, the network has been trained during a training period in which the sensory input representations (i.e.,  $I_A$  and  $I_V$ ) have been given with a random distribution.

The synapses describing the receptive field,  $r_{S,kj}$ , and those describing the cross-modal link between the two areas,  $w_{SQ,kj}$ , have been trained using a learning rule with a classical Hebbian potentiation factor and a decay term. We can write, in scalar form

$$\Delta r_{S,kj} = \gamma y_{S,k} \left( i_{S,j} - r_{S,kj} \right) \qquad \text{with} \quad S = A, V \tag{S8}$$

$$\Delta w_{SQ,kj} = \gamma \, y_{S,k} \left( y_{Q,j} - w_{SQ,kj} \right) \qquad \text{with} \quad S = A, \ V \quad Q = A, \ V \quad Q \neq S \tag{S9}$$

Eqs. (8) and (9) have been applied, at each step, using the final steady state values of the neuron output (i.e., when transient phenomena are exhausted).

At the beginning of training all cross-modal synapses are assumed equal to zero. Conversely, the receptive-field synapses have a broad spatial extension, and moderate amplitude, identical for the two modalities, i.e.

$$r_{S,kj} = r_0 \exp\left(-\frac{\left(d\left(\theta_j, \theta_k\right)\right)^2}{2\sigma_R^2}\right) \text{ with } S = A, V$$
(S10)

where  $r_0$  sets the initial strength of the receptive field, and  $\sigma_R$  establishes its initial spatial extension (we assume  $\sigma_R > \sigma_A$  and  $\sigma_R > \sigma_V$  i.e. a wide initial receptive fields). Of course, Eq. (10) holds only at the first step of training.

#### **3** Probability distribution and spatial accuracy of the inputs

According to the previous section, we assumed that the sensory inputs are composed of a deterministic term, which represents the spatial distribution of the input, centered on the stimulus spatial position, and a Gaussian white noise term (zero mean value and assigned standard deviation). Hence

$$i_{S,k}(\theta) = \frac{i_{S,Strength}}{\sqrt{2\pi\sigma_S^2}} \exp\left(-\frac{\left(d(\theta_s, \theta_k)\right)^2}{2\sigma_s^2}\right) + n_{S,k}(\theta) \quad S = A \text{ or } V$$
(S11)

where  $\theta_s$  represents the spatial position of the stimulus,  $i_{s,Strength}$  is the stimulus strength in the absence of noise, and  $\sigma_s$  is the standard deviation of the spatial representation. According to physiology, we assumed that the visual inputs are spatially more accurate than the auditory ones, hence we set  $\sigma_v < \sigma_A$ . Conversely, we assumed that the standard deviation of noise (say  $v_s$ ) is a given fraction of the input strength, to set the signal to noise ratio (see table 1 in the text).

In order to simulate the presence of better acuity at the center, and reduced acuity at the periphery, we assumed that the SDs of the visual and auditory inputs increase with the eccentricity of the stimulus.

The expression of  $\sigma_{\nu}$  has been taken from an empirical curve on visual acuity by Dacey (1993) (see also Ursino et al., 2017 for more details). By denoting with  $e_{\nu} = \theta_{\nu}$  the eccentricity with respect to the fovea, we have

$$\sigma_{\nu}(e_{\nu}) = \sigma_{\nu_0} + \frac{\varepsilon\sqrt{3}}{60} \left( 0.058e_{\nu} + 0.022e_{\nu}^2 - 0.00022e_{\nu}^3 \right)$$
(S12)

 $\sigma_{V0}$  represents the SD of the visual inputs at the fovea (i.e., at zero eccentricity). We used the same value as in the previous paper, i.e.  $\sigma_{V0} = 4$  deg. Finally, we use a parameter,  $\varepsilon$ , to adapt the function so that,  $\sigma_V$  ranges between 4 deg, at 0 eccentricity, to about 12 deg at maximum eccentricity.

The auditory acuity also decreases from the center to the periphery, although it is difficult to quantify this effect being influenced by many factors, such as the stimulus intensity and frequency (Middlebrooks and Green, 1991; Wood and Bizley, 2015). However, this effect is less evident and of smaller entity compared with the visual one (Perrott and Saberi, 1990). Hence, we used a simpler linear relationship, assuming that  $\sigma_{A0}$  linearly increases from about 20 deg at the fovea to 30 deg at the periphery:

$$\sigma_{A}(e_{A}) = \sigma_{A0} + 10|e_{A}|/90 \tag{S13}$$

with  $\sigma_0^A = 20$  and  $e_A = \theta_A$  is the eccentricity of the auditory position with respect to the head center.

The positions of the two stimuli (i.e.,  $\theta_A$  and  $\theta_v$  in Eq. S11) have been randomly generated from the prior probability distribution described below.

We assume that both the visual and auditory input have a greater probability close to the fovea, and smaller probability at the periphery. This corresponds to have a non-uniform prior in visual unisensory conditions. The following probabilities have been used to generate the position of the visual and auditory inputs during training. <u>Visual unisensory prior</u>: the visual position follows a Gaussian distribution, centered at the fovea. Hence

$$p(\theta_{\nu}) = \frac{1}{\sqrt{2\pi s_{\nu}^2}} \exp\left(-\frac{\theta_{\nu}^2}{2 s_{\nu}^2}\right)$$
(S14)

The standard deviation  $s_V$  (which here plays the role of a space constant) has been set at 7 deg; i.e., the visual stimuli becomes very rare at ±20 deg eccentricity.

<u>Auditory unisensory prior</u>: the auditory position follows a Gaussian distribution, centered at the head center.

$$p(\theta_{A}) = \frac{1}{\sqrt{2\pi s_{A}^{2}}} \exp\left(-\frac{\theta_{A}^{2}}{2 s_{A}^{2}}\right)$$
(S15)

The standard deviation is assumed higher than in the visual case: we have  $s_A = 30$  deg assuming that, head movements in auditory unimodal conditions are less efficient than eye movement in visual unimodal conditions to maintain the stimulus close to the center.

<u>Cross modal prior</u>: in the cross modal case during training, we assumed that the visual and auditory inputs originate from independent causes with a given probability (say  $\beta$ ) but are produced by the same cause, hence originate from proximal spatial positions, with the complementary probability (1 –  $\beta$ ). According to the Bayes rule, the joined prior probability can be computed from knowledge of the individual probability of one stimulus, and the conditional probability of the other. A problem is whether, in cross modal conditions, the distribution is dominated by the visual prior (more sharply close to the center) or by the auditory one (less sharply close to the center). We assumed that, in 50% of cases, the cross-modal stimuli follow the visual distribution and in the other 50% of cases follow the auditory one. Hence

$$p(\theta_{\nu},\theta_{\lambda}) = 0.5p(\theta_{\nu})p(\theta_{\lambda}|\theta_{\nu}) + 0.5p(\theta_{\lambda})p(\theta_{\nu}|\theta_{\lambda})$$
(S16)

where we used equations (S14) and (S15) for the visual and auditory priors, and the following expression for the conditional probability

$$p(\theta_{A}|\theta_{V}) = \beta \ p(\theta_{A}) + (1-\beta) \frac{1}{\sqrt{2\pi \ s_{AV}^{2}}} \exp\left(-\frac{d(\theta_{A},\theta_{V})^{2}}{2 \ s_{AV}^{2}}\right)$$

$$p(\theta_{V}|\theta_{A}) = \beta \ p(\theta_{V}) + (1-\beta) \frac{1}{\sqrt{2\pi \ s_{AV}^{2}}} \exp\left(-\frac{d(\theta_{A},\theta_{V})^{2}}{2 \ s_{AV}^{2}}\right)$$
(S17)

In writing Eq. (S17) we assumed that the conditional probability is computed as the weighted sum of the prior unimodal distribution, reflecting the moderate possibility that the two stimuli are independent, and a second term,  $\frac{1}{\sqrt{2\pi s_{AV}^2}} \exp\left(-\frac{d(\theta_A, \theta_V)^2}{2s_{AV}^2}\right)$ , reflecting the probability that the auditory and visual events are originated from the same source.

As in the previous work, we used a value of space constant  $s_{AV} = 1$  deg, assuming a small audio-visual distance when the two stimuli originate from the same source.

## References

- Dacey, D.M. (1993) The mosaic of midget ganglion cells in the human retina. *Journal of Neuroscience*, *13*(12), pp. 5334-5355.
- Middlebrooks, J.C. and Green, D.M. (1991) Sound localization by human listeners. *Annual review of psychology*, *42* (1), pp. 135-159.
- Perrott, D.R. and Saberi, K. (1990) Minimum audible angle thresholds for sources varying in both elevation and azimuth. *The Journal of the Acoustical Society of America*, 87(4), pp. 1728-1731.
- Ursino, M. *et al.* (2017a) Development of a Bayesian Estimator for Audio-Visual Integration: A Neurocomputational Study. *Frontiers in computational neuroscience*, 11, 89.
- Wood, K.C. and Bizley, J.K. (2015) Relative sound localisation abilities in human listeners. *The Journal of the Acoustical Society of America*, *138*(2), pp. 674-686.