

IMPEP – Integrated Multimodal Perception Experimental Platform

A Bayesian Binaural System for 3D Sound-Source Localisation

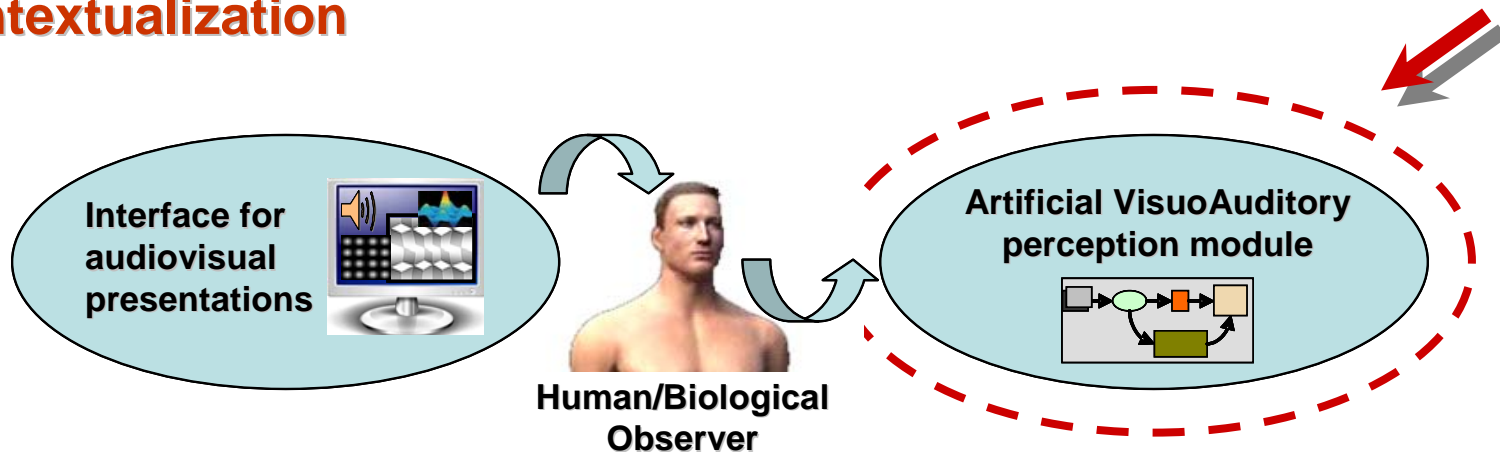


February 2008

Contact Persons:

Cátia Pinho, João Filipe Ferreira, Jorge Dias
Email: {catiap; jorge; jfilipe}@isr.uc.pt

Contextualization



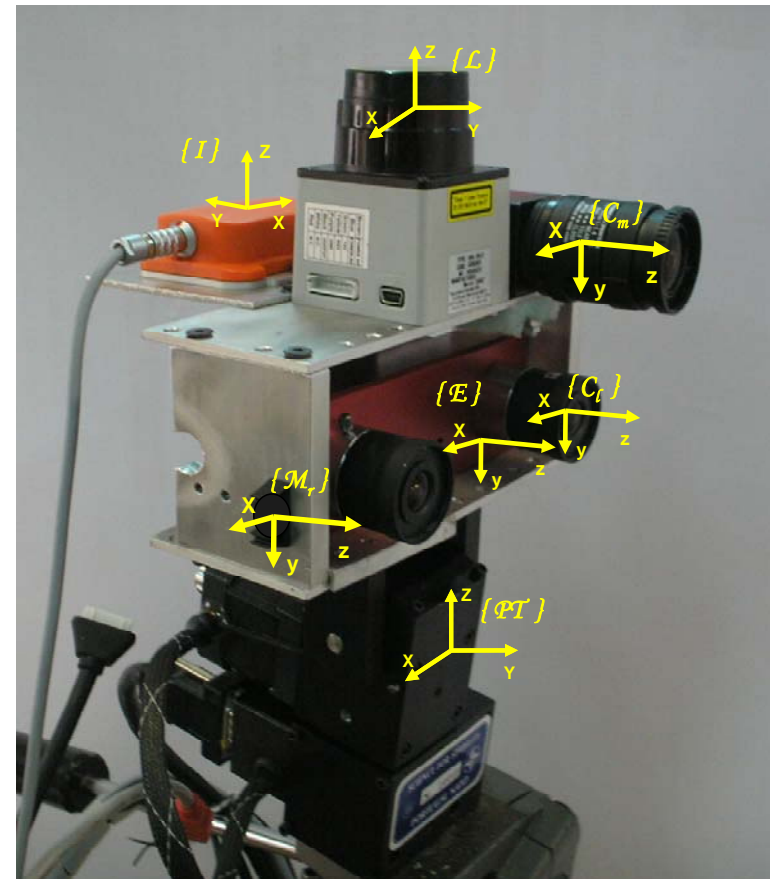
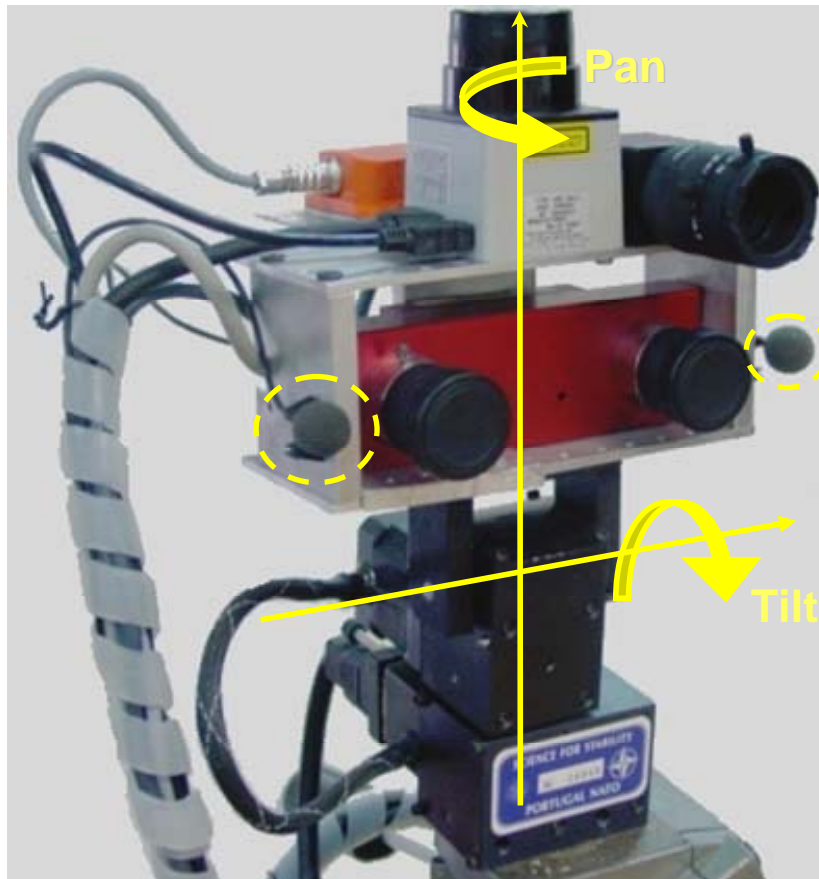
Modelling Approach

- To develop a **visuoauditory perception solution** with the ability to yield **probabilistic estimates** for the complete frame of 3D information:
 - *azimuth* (θ), *elevation* (ϕ), *distance* (ρ)

A Bayesian Binaural System for 3D Sound-Source Localisation

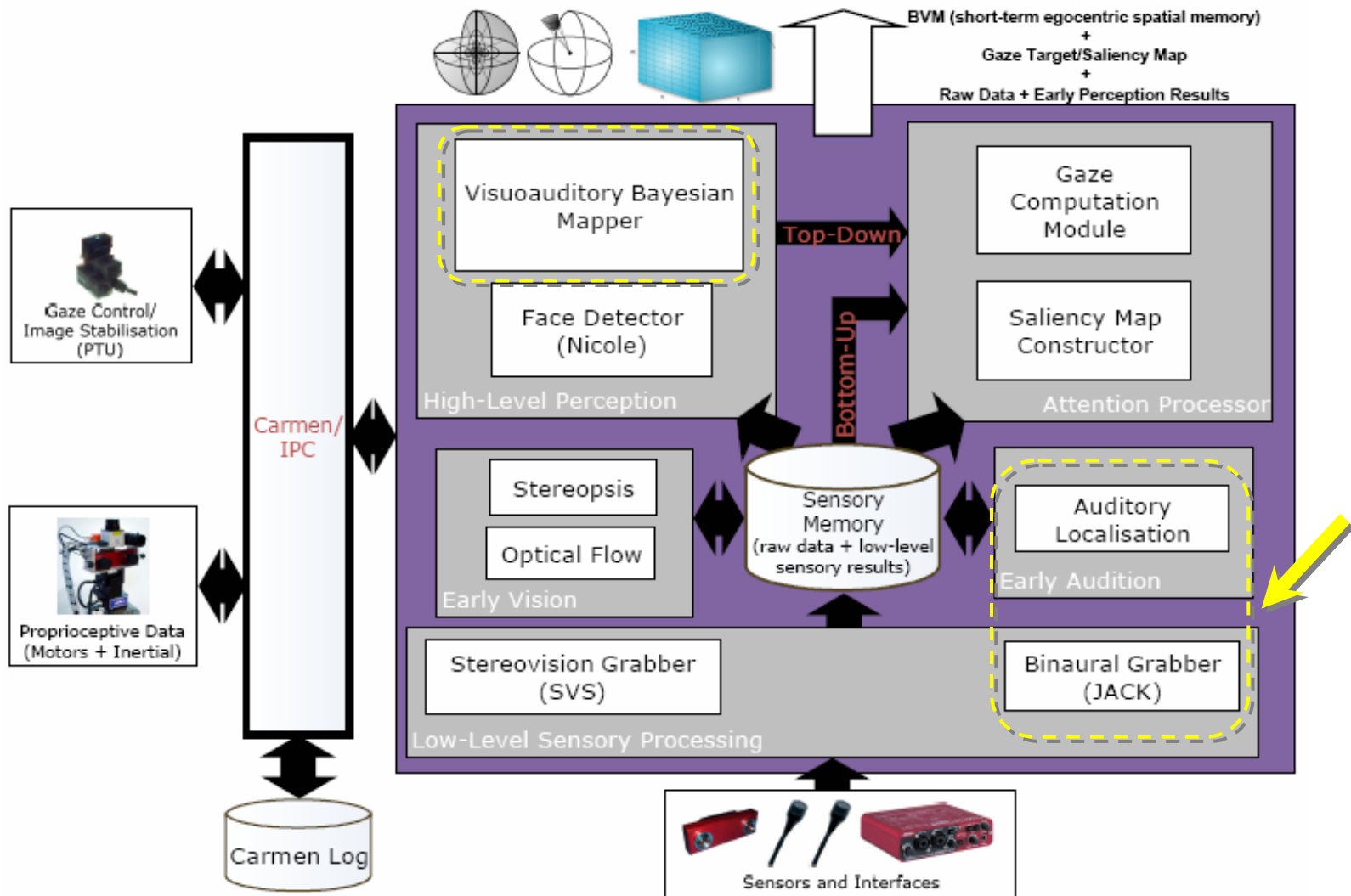
IMPEP 1 (Integrated Multimodal Perception Experimental Platform)

Platform description

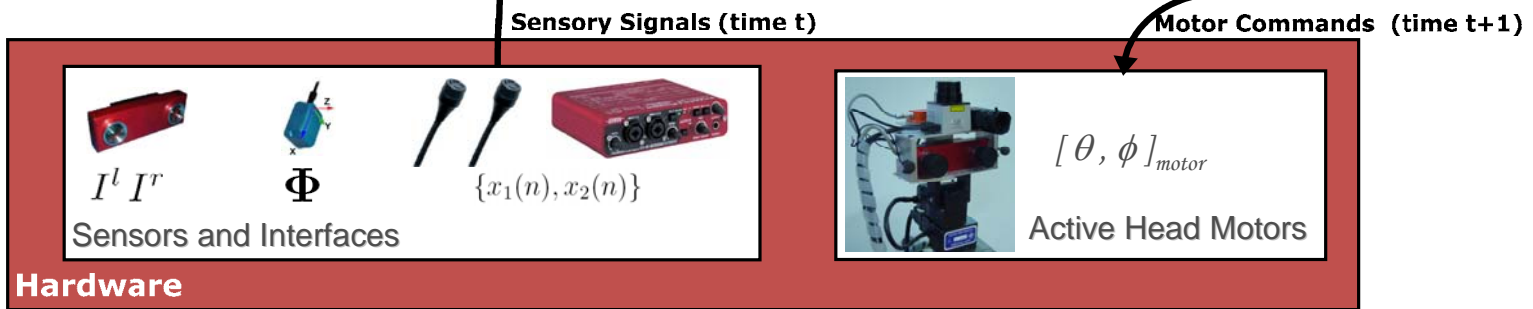
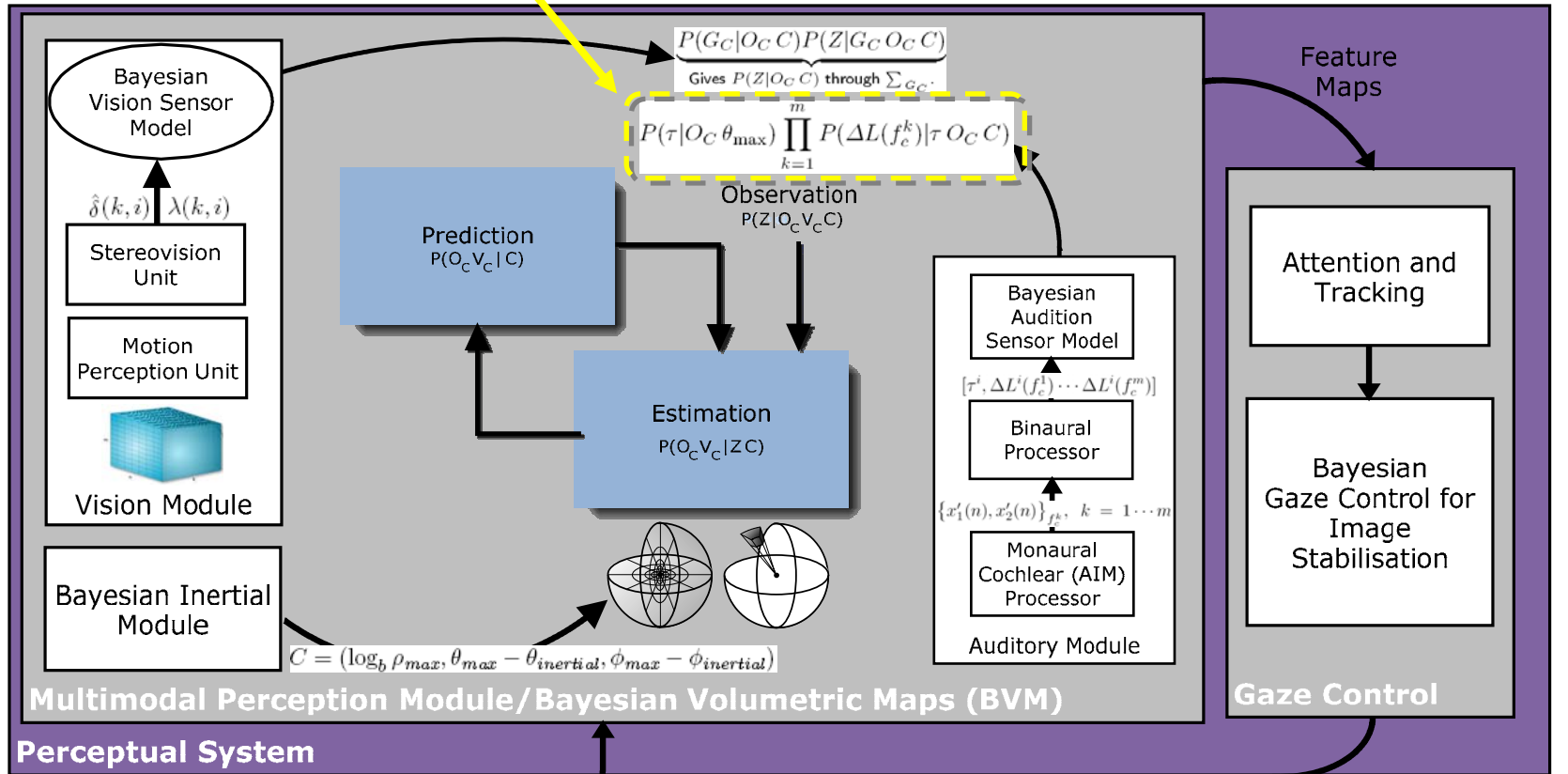


A Bayesian Binaural System for 3D Sound-Source Localisation

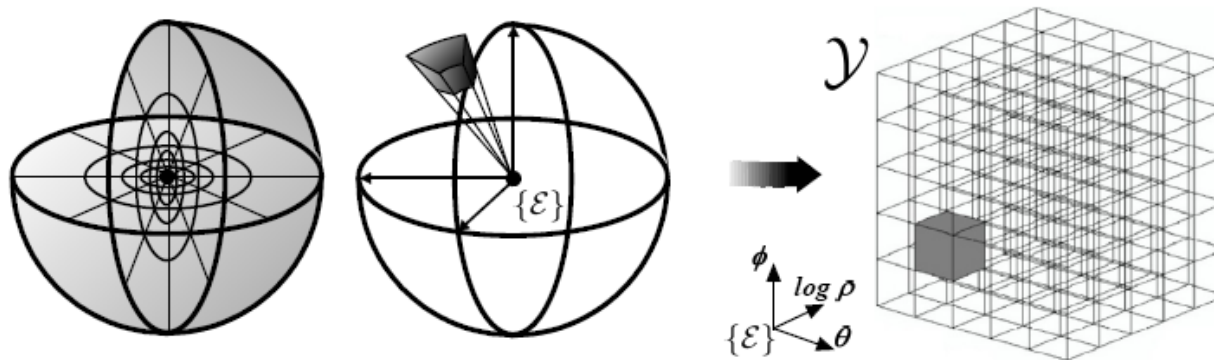
IMPEP - Outputs (to Multimodal Integration System)



Contextualisation of Auditory Sensor Model within the IMPEP System



BVM - Bayesian Volumetric Map



- An egocentric, log-spherical spatial memory map has been devised as a framework for multimodal sensor fusion, named the Bayesian Volumetric Map (BVM)
- This map stores the independent probabilistic states of occupancy and velocity for each cell in a volumetric grid with log-spherical configuration, defined by the domain \mathcal{Y} :

$$\mathcal{Y} \equiv] \log_b \rho_{\text{Min}} ; \log_b \rho_{\text{Max}}] \times] \theta_{\text{Min}} ; \theta_{\text{Max}}] \times] \phi_{\text{Min}} ; \phi_{\text{Max}}]$$

DASM - Direct Auditory Sensor Model

- Defined as a **normally distributed** BVM operator

$$P(Z | O_C C) = P(\tau | O_C \theta_{\max}) \prod_{k=1}^m P(\Delta L(f_c^k) | \tau O_C C)$$

C – Cell Identifier ($\log_b \rho_{\max}$, θ_{\max} , ϕ_{\max})

Z – Measurement taken by the audio sensor (generic notation)

O_C – Occupancy of cell C ($0 = \text{no sound-source}$; $1 = \text{sound-source}$)

τ – frequency invariant interaural time differences (**ITDs**) - ms

$\Delta L(f_c^k)$ – frequency dependent interaural level differences (**ILDs**) - dB

- The direct sensor model is used to update the BVM by **inverting the direct model** (DASM) using Bayesian inference so as to estimate the probability of occupancy for each cell
- Note:** the angular and range resolution of the auditory sensor space is typically lower than the BVM

IASM - Inverse Auditory Sensor Model

$$P([O_c = 1] | ZC) = \frac{P([O_c = 1] | C)P(Z | [O_c = 1]C)}{\sum_{O_c=0,1} P(O_c | C)P(Z | O_c C)}$$

$$= \frac{P([O_c = 1] | C)P(Z | [O_c = 1]C)}{P([O_c = 0] | C)P(Z | [O_c = 0]C) + P([O_c = 1] | C)P(Z | [O_c = 1]C)}$$

DASM

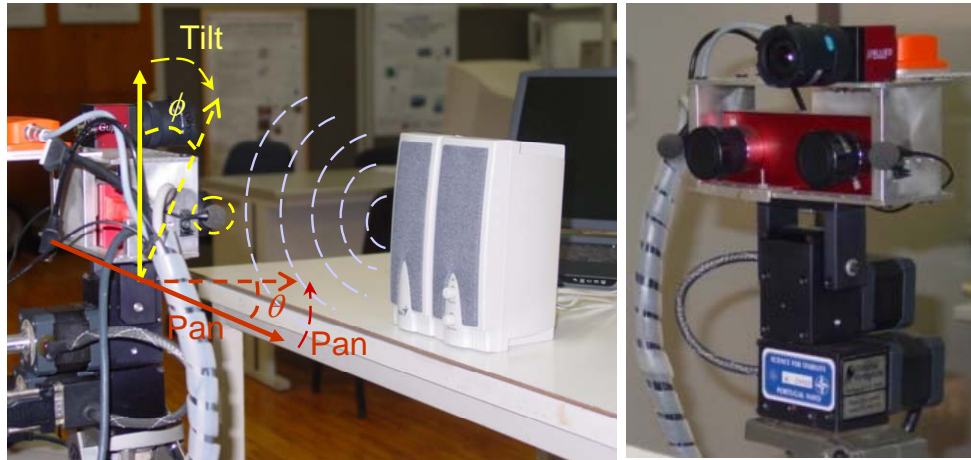
$$P(Z | [O_c = 0,1]C) \begin{cases} O_c = 0 \rightarrow P(\tau | [O_c = 0]C) \prod_{k=1}^m P(\Delta L(f_c^k) | \tau [O_c = 0]C) \\ O_c = 1 \rightarrow P(\tau | [O_c = 1]C) \prod_{k=1}^m P(\Delta L(f_c^k) | \tau [O_c = 1]C) \end{cases}$$

Objective

- The auditory calibration's purpose is to characterise the normal distributions of the **Direct Auditory Sensor Model (DASM)**
- This will allow the full localisation of sound-sources in three-dimensional space:
 - Azimuth (θ)
 - Elevation (ϕ)
 - Distance (ρ)
- **Final aim:** To feed the BVM framework with an accurate direct audition sensor model so as to allow multimodal perception of 3D structure and motion.

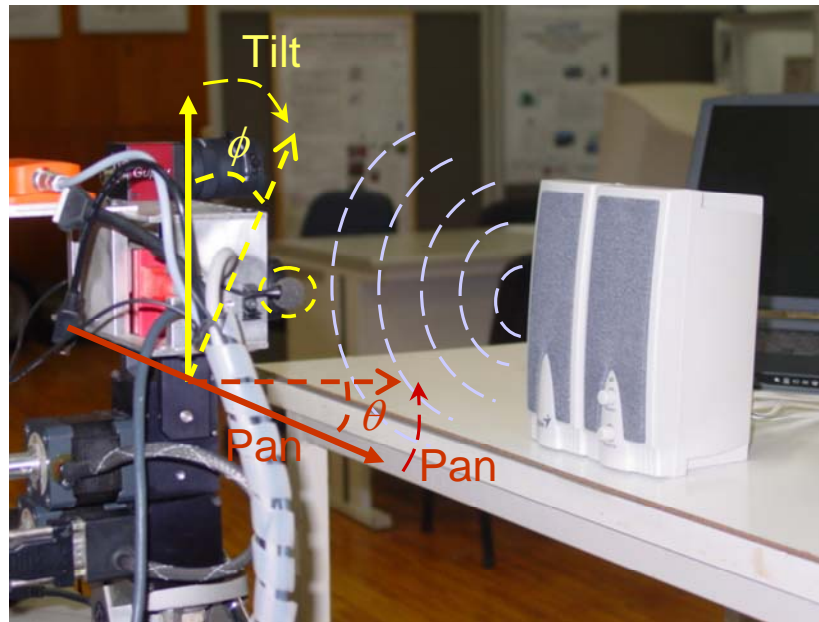
Experimental plan

- **Goal:** to capture binaural readings using the stereo microphones of the IMPEP Active Perception Head for each cell in the auditory sensor space of a broadband noise sound-source (1s)



- To cover the complete auditory sensor space, the sound-source must be positioned at the centre of each grid cell on the BVM
- **Methodology:** sound-source is positioned at a specific distance (ρ) from the IMPEP head, directly facing the front of the Pan & Tilt Unit (PTU), and the corresponding relative rotation is performed by the PTU; The different distances (ρ) between the sound-source and the IMPEP head are obtained manually.
- To avoid redundancies \Rightarrow **only one quadrant** is used due to:
 - Symmetry from **front-back confusion phenomenon**
 - Left-Right anti-symmetry ($ITD = -ITD$ and $ILD = -ILD$)

Experimental plan (example)



Resolution (Pan / Tilt):

- 2 degree of azimuth (θ) – performed by the **Pan** motor – $360^\circ / 2^\circ = 180$ cells (θ)
- 10 degrees of elevation (ϕ) – performed by the **Tilt** motor – $180^\circ / 10^\circ = 18$ cells (ϕ)
- Acquisition for $N_d = 2$ different distances $\Rightarrow d : d_1 \sim 0.32\text{m} ; d_2 \sim 3.2 \text{ m} = 2$ cells (ρ)

Experimental plan (example contd.)

To avoid redundancies \Rightarrow **only one quadrant** is used

Because:

- Symmetry from **front-back confusion phenomenon**
- Left-Right anti-symmetry ($ITD = -ITD$ e $ILD = -ILD$)

- N_d – number of distances
- N_m – numbers of measurements in each cell

Auditory sensor space angular ranges simplify to (including PTU spec limitations for elevation):

$$\text{Azimuth (} \theta \text{) : } 90^\circ / 2^\circ = \mathbf{45 \text{ cells}}$$

$$\text{Elevation (} \phi \text{) : } (-30^\circ \text{ to } 30^\circ) / 10^\circ = \mathbf{6 \text{ cells}}$$

Consequently: Azimuth meas. Elevation meas.

- $[N_d \times (\boxed{45} \times \boxed{6})] \times N_m = 2 \times 270 \times N_m = \mathbf{540} \times N_m$ sets of measurements
(N_m measurements = 20 stimulus) in each place to perform a statistical description

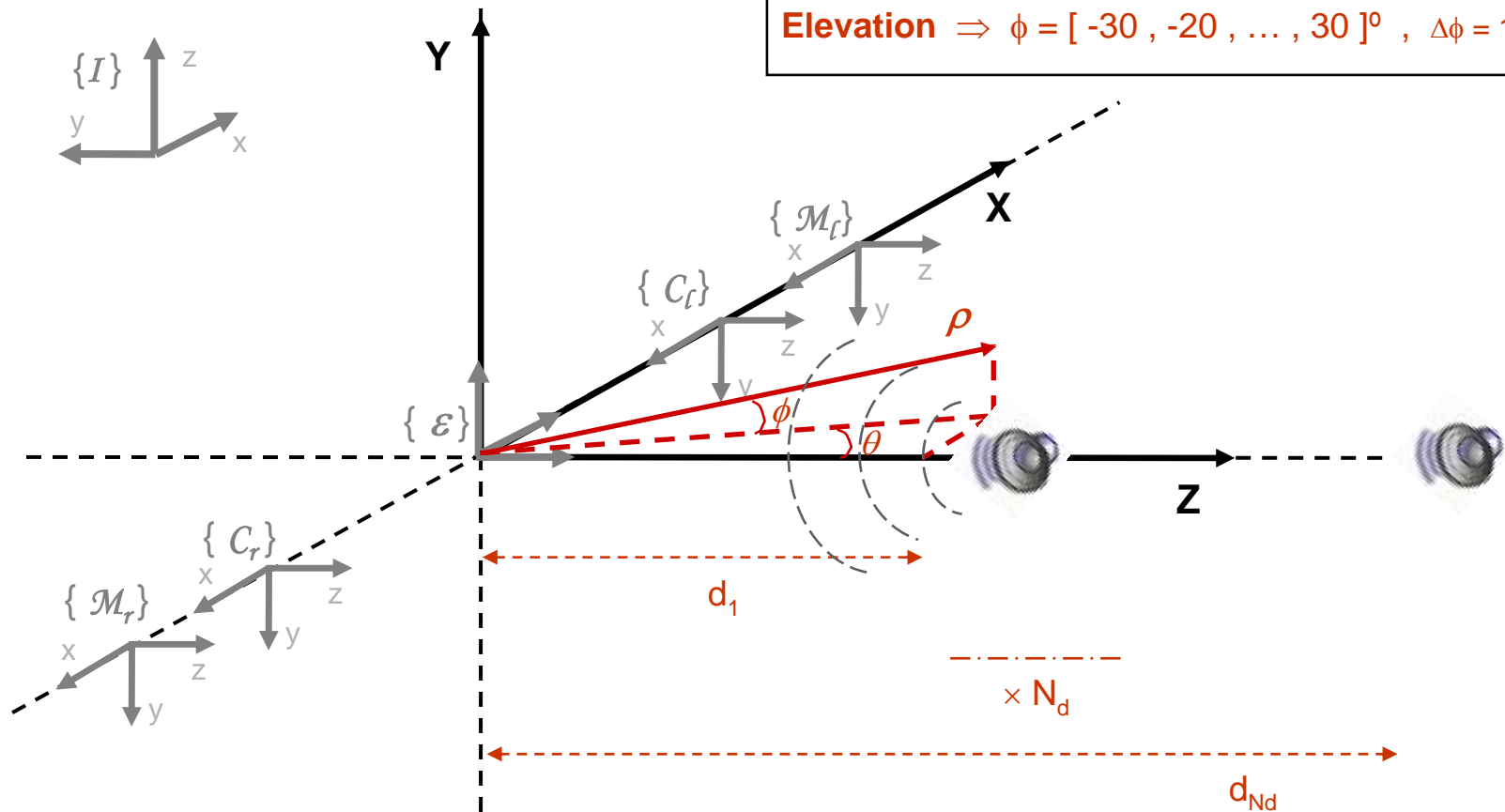
Becoming:

- $540 \times 20 = 10800$
- if each measurement takes 1s plus 1s of pause (play / record), calibration for each distance (i.e. the calibration process is conveniently divided into N_d sessions) will take: $10800 \times 2s / 2 = \mathbf{3 h}$

Schematic of the experimental acquisition

Azimuth $\Rightarrow \theta = [0, 2, \dots, 90]^\circ$, $\Delta\theta = 2$

Elevation $\Rightarrow \phi = [-30, -20, \dots, 30]^\circ$, $\Delta\phi = 10$

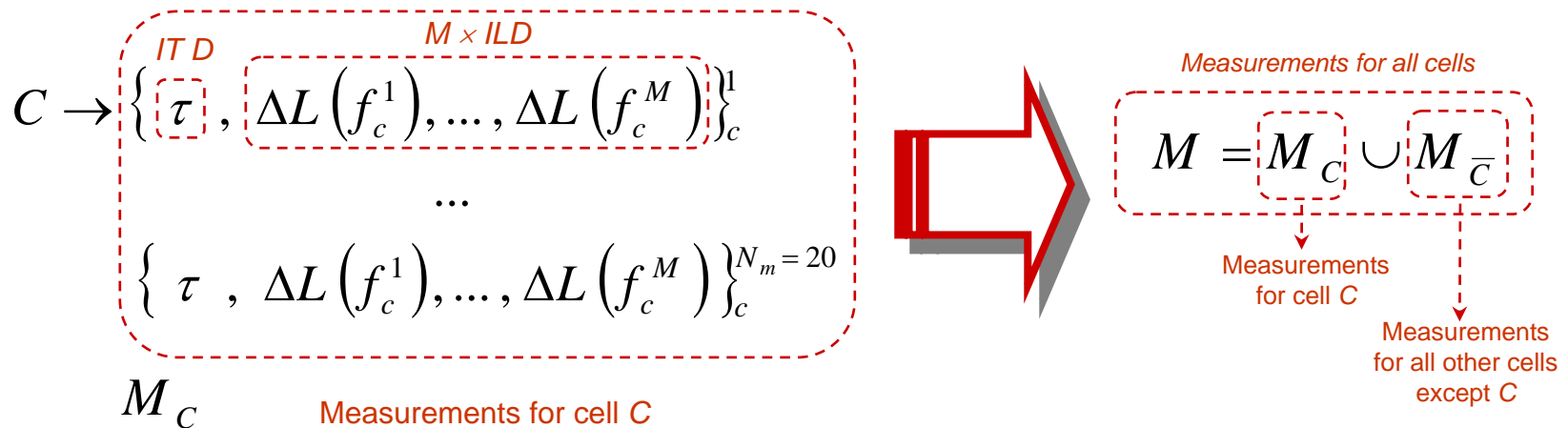


Processing of the data given by Auditory Calibration

- Mathematical formulation - *Direct Auditory Sensor Model*

$$P(\tau | O_C \theta_{\max}) \prod_{k=1}^m P(\Delta L(f_c^k) | \tau O_C C)$$

- Measurement definitions (after applying AIM and binaural processing to binaural readings, [1,2,3,4]):



Processing of the data given by Auditory Calibration

- Normal probability distribution statistical characterisation process:

$$P(\tau | O_C \theta_{\max}) = \begin{cases} P(\tau | [O_C = 1] \theta_{\max}) \Rightarrow \mu_{M_C}(\tau), \sigma_{M_C}(\tau) \\ \text{Cell occupied} \quad \text{Measurements average} \quad \text{Standard deviation} \\ P(\tau | [O_C = 0] \theta_{\max}) \Rightarrow \mu_{M_{\bar{C}}}(\tau), \sigma_{M_{\bar{C}}}(\tau) \\ \text{Cell not occupied} \end{cases}$$

$$P(\Delta L | \tau O_C C) \approx \begin{cases} P(\Delta L | [O_C = 1] C) \Rightarrow \mu_{M_C}(\Delta L), \sigma_{M_C}(\Delta L) \\ \text{Cell occupied} \quad \text{Measurements average} \quad \text{Standard deviation} \\ P(\Delta L | [O_C = 0] C) \Rightarrow \mu_{M_{\bar{C}}}(\Delta L), \sigma_{M_{\bar{C}}}(\Delta L) \\ \text{Cell not occupied} \end{cases}$$

Conclusions

- The auditory calibration's purpose is of characterising the normal distributions of the **DASM** (Direct Auditory Sensor Model) is thus solved efficiently.
- This will allow the full localisation of sound-sources in three-dimensional space:
 - Azimuth (θ)
 - Elevation (ϕ)
 - Distance (ρ)
- within the BVM framework

References:

- [1] - Faller, C. and Merimaa, J. Source localization in complex listening situations: Selection of binaural cues based on interaural coherence. *J. Acoust. Soc. Am.*, 116:30753089, 2004.
- [2] - Akeroyd, M. A. A Binaural Cross-correlogram Toolbox for MATLAB. February 2001.
- [3] - Johannesma, P. I. M. The pre-response stimulus ensemble of neurons in the cochlear nucleus. In *Symposium on Hearing Theory*, pages 5869. IPO, Eindhoven, The Netherlands, 1972.
- [4] - Slaney, M. Auditory Toolbox: A MATLAB Toolbox for Auditory Modeling Work. Technical report, Interval Research Corporation, 1998.
- [5] - C.J. Lee, S.D. Wang, A Gabor filter-based approach to fingerprint recognition, in: *IEEE Workshop on Signal Processing System, SiPS 99*, pp. 371-378, 1999.
- [6] - Dankers, A., Barnes, N., and Zelinsky, A. Primate structures in synthetic dynamic active visual saliency. 6th International Conference on Epigenetic Robotics (Epirob 2006), Modeling Cognitive Development in Robotic Systems, Paris, France 2006.
- [7] - Itti, L., Koch, C., and Niebur, E. A Model of Saliency-Based Visual Attention for Rapid Scene Analysis. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 20(11): 1254–1259, November 1998.
- [8] - Shic F, Scassellati B. A Behavioral Analysis of Computational Models of Visual Attention. *International Journal of Computer Vision*, 73(2):159-177, Jun. 2007.