

Detection, Localization and Tracking of 3D Audio-Visual Objects

Vassil Khalidov, Florence Forbes, Miles Hansard, Elise Arnaud,
Radu Horaud
INRIA Grenoble, France

In collaboration with Heidi Christensen, Sheffield University

Work done under the POP project:
<http://perception.inrialpes.fr/POP/>

Outline

- Introduction
- Audio-Visual Features Integration
- Multimodal Clustering Framework
 - Conjugate Mixture Models
 - Conjugate Particle Tracking
- Experimental Results
- Conclusions and future work

Interpreting the scene...



Several problems to be addressed



- Estimate the number of audio-visual objects
- Localize and track every object
- Select the ones that are involved in auditory activity

Several problems to be addressed



- Estimate the number of audio-visual objects
- Localize and track every object
- Select the ones that are involved in auditory activity

Several problems to be addressed



- Estimate the number of audio-visual objects
- Localize and track every object
- Select the ones that are involved in auditory activity

Several problems to be addressed



- Estimate the number of audio-visual objects
- Localize and track every object
- Select the ones that are involved in auditory activity

Main difficulties

- Steps towards audio-visual perception: How can we recognize an object that is both seen and heard?
- In most natural situations it is difficult to extract unambiguous information from a single modality
- Visual data is dense, light sources are not relevant
- Auditory data is sparse, acoustic sources must be detected in the presence of reverberations
- Integration required !

Main difficulties

- Steps towards audio-visual perception: How can we recognize an object that is both seen and heard?
- In most natural situations it is difficult to extract unambiguous information from a single modality
- Visual data is dense, light sources are not relevant
- Auditory data is sparse, acoustic sources must be detected in the presence of reverberations
- Integration required !

Audio-Visual integration types

- Integration is based on features' coherence
- Two coherence types:
 - static: features related to the same location in ambient space
 - dynamic: correlated evolution through time
- Which features to choose?

Experimental setup



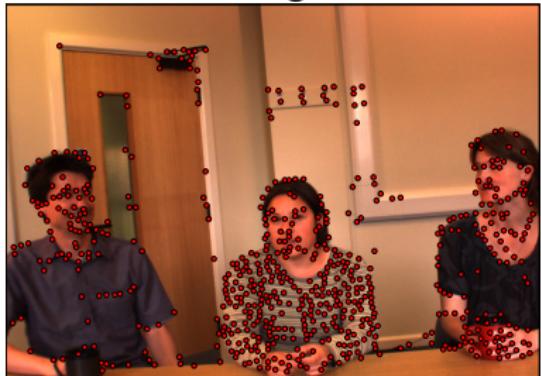
- Stereoscopic camera pair
- Binaural microphone pair

Monocular observations (1)

Right camera image:



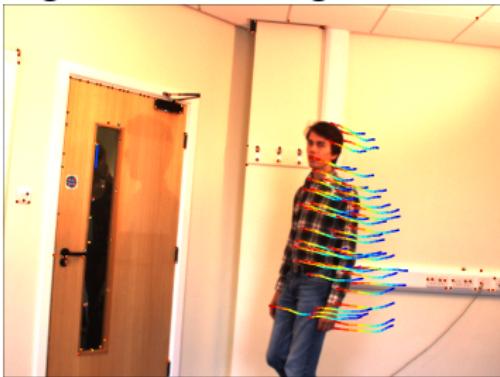
Left camera image:



- Detect 'interest points' in the left and right images (Harris et al., 1988)
- $\mathbf{f}^l = \{\mathbf{f}_1^l, \dots, \mathbf{f}_{m^l}^l, \dots, \mathbf{f}_{M^l}^l\} \in \text{2D}$
- $\mathbf{f}^r = \{\mathbf{f}_1^r, \dots, \mathbf{f}_{m^r}^r, \dots, \mathbf{f}_{M^r}^r\} \in \text{2D}$

Monocular observations (2)

Right camera image:



Left camera image:



- Detect monocular motion in the left and right images (Lucas et al., 1981)
- $\mathbf{w}^l = \{\mathbf{w}_1^l, \dots, \mathbf{w}_{m^l}^l, \dots, \mathbf{w}_{M^l}^l\} \in \text{2D}$
- $\mathbf{w}^r = \{\mathbf{w}_1^r, \dots, \mathbf{w}_{m^r}^r, \dots, \mathbf{w}_{M^r}^r\} \in \text{2D}$

Binocular observations (1)

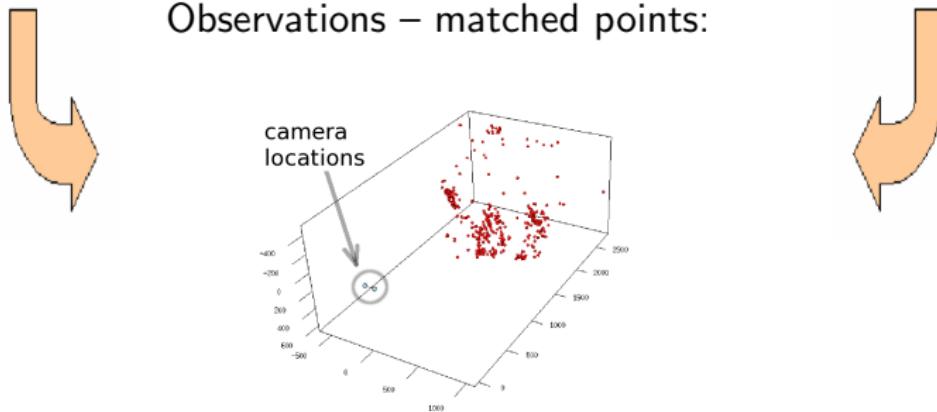
Right camera image:



Left camera image:



Observations – matched points:



Binocular observations (2)

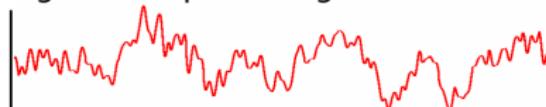
- $\mathbf{f} = \{\mathbf{f}_1, \dots, \mathbf{f}_m, \dots, \mathbf{f}_M\} \in \text{3D}$
- $\mathbf{f}_m = (u, v, d)$: (u, v) = pixel, d = disparity
- $\mathbf{s} = (x, y, z)$: 3D coordinates of an audio-visual object
- 3D visual disparity is linked to 3D coordinates :
$$\mathbf{f}_m = (u, v, d) = \mathcal{F}(\mathbf{s}) = \left(\frac{x}{z}, \frac{y}{z}, \frac{B}{z} \right)$$

Binaural observations (1)

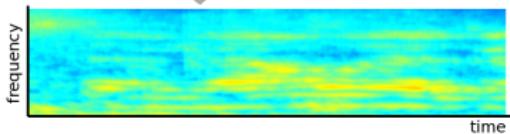
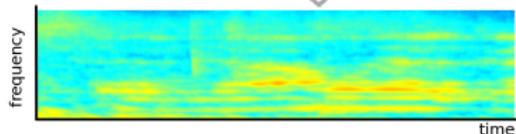
Left microphone signal



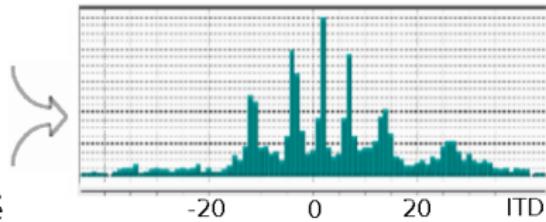
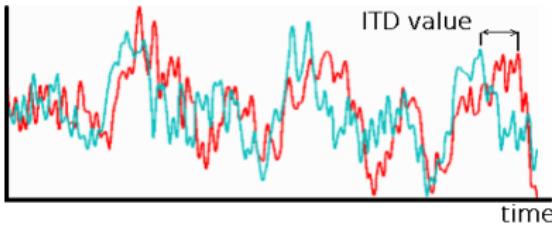
Right microphone signal:



Gammatone filterbank

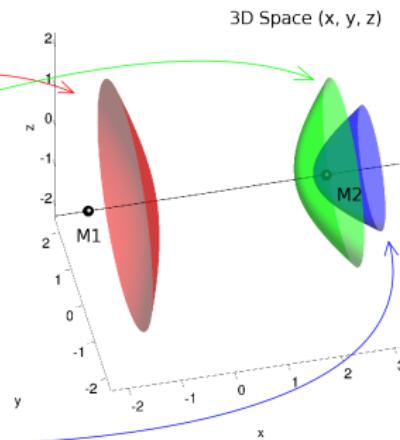
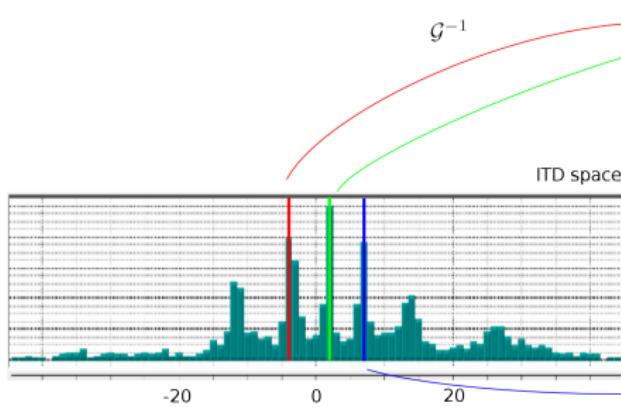


Interaural time differences (ITDs) are obtained through the analysis of cross-correlogram of the filtered signals



Binaural observations (2)

- $\mathbf{g} = \{g_1, \dots, g_k, \dots, g_K\} \in \text{1D}$
- g_k = interaural time difference (ITD)
- 1D auditory disparity (ITD) is related to 3D coordinates:
$$g_k = \mathcal{G}(\mathbf{s}) = \frac{1}{c} (\|\mathbf{s} - \mathbf{s}_{M_1}\| - \|\mathbf{s} - \mathbf{s}_{M_2}\|)$$



The geometry of 3D audio-visual fusion

- $s = (x, y, z)$: 3D coordinates of an audio-visual object
- 1D auditory disparity (ITD):
$$g_k = \mathcal{G}(s) = \frac{1}{c} (\|s - s_{M_1}\| - \|s - s_{M_2}\|)$$
- 3D visual disparity (projective space):
$$f_m = (u, v, d) = \mathcal{F}(s) = \left(\frac{x}{z}, \frac{y}{z}, \frac{B}{z}\right)$$
- Auditory and visual data are put on an equal footing

Approach

- Exploit temporal and spatial relations between binocular and binaural observations
- Cast the audio-visual localization problem in the framework of maximum likelihood with missing data
- Maximize the expected complete-data log-likelihood:

$$E[\log P(\underbrace{\text{audio, vision}}_{\text{observed data}}, \underbrace{\text{audio-visual objects}}_{\text{missing data}})]$$

Associating observations with objects

- Objects $1, 2, \dots, N$; outlier class $N + 1$
- Two sets of assignment variables (the hidden data):
 $\mathbf{A} = \{A_1, \dots, A_M\}$ for video, and
 $\mathbf{A}' = \{A'_1, \dots, A'_K\}$ for audio
- The notation $A_m = n$, $n = 1, \dots, N, N + 1$ means that the visual observation m is assigned to object n or it is assigned to an outlier class $N + 1$.
- Similarly, $A'_k = n$, $n = 1, \dots, N, N + 1$ means that the audio observation k is assigned to object n or it is assigned to an outlier class $N + 1$.
- Both \mathbf{A} and \mathbf{A}' are missing: the *observation-to-object assignments* are unknowns.

Clustering with conjugate mixture models

- $s = \{s_1, \dots, s_n, \dots, s_N\}$ - object locations in 3D
- Probability of an observation to belong to an object (inlier):
 - $P(f_m | A_m = n) = \mathcal{N}(f_m | \mathcal{F}(s_n), \Sigma_n)$
 - $P(g_k | A'_k = n) = \mathcal{N}(g_k | \mathcal{G}(s_n), \sigma_n^2)$
 - Alternatively, we can also use the t-distribution.
 - The methodology is independent of the normal/t-distribution choice.
- Likelihood of an observation to be an outlier:
 - $P(f_m | A_m = N + 1) = \mathcal{U}(V_{3D}) = \frac{1}{V}$
 - $P(g_k | A'_k = N + 1) = \mathcal{U}(U_{1D}) = \frac{1}{U}$
- The model's parameters:

$$\theta = \{s_1, \dots, s_N, \Sigma_1, \dots, \Sigma_N, \sigma_1, \dots, \sigma_N\}$$

The observed-data log-likelihood

We assume that the variables are independent and identically distributed, then:

$$\log P(\mathbf{f}, \mathbf{g}) = \sum_{n=1}^{N+1} \log \left(\pi_n \sum_{m=1}^M P(\mathbf{f}_m | A_m = n) + \pi'_n \sum_{k=1}^K P(\mathbf{g}_k | A'_k = n) \right)$$

where π_n and π'_n are the prior probabilities.

The maximization of the log-likelihood is not tractable because of the presence of the hidden variables.

Maximum likelihood with hidden data

- The expectation-maximizaton (EM) algorithm finds a local maximum of the likelihood
- Maximize the expected complete-data log-likelihood w.r.t. θ :

$$\begin{aligned} E_{\mathbf{A}, \mathbf{A}'} [\log P(\mathbf{f}, \mathbf{g}, \mathbf{a}, \mathbf{a}') | \mathbf{f}, \mathbf{g}] = \\ -\frac{1}{2} \sum_{m=1}^M \sum_{n=1}^N \alpha_{mn} \left((\mathbf{f}_m - \mathcal{F}(\mathbf{s}_n))^{\top} \Sigma_n^{-1} (\mathbf{f}_m - \mathcal{F}(\mathbf{s}_n)) + \log |\Sigma_n| \right) \\ -\frac{1}{2} \sum_{k=1}^K \sum_{n=1}^N \alpha'_{kn} \left((g_k - \mathcal{G}(\mathbf{s}_n))^2 / \sigma_n^2 + \log \sigma_n^2 \right) \\ + \text{ other constant terms} \end{aligned}$$

Posterior probabilities

- α_{mn} and α'_{kn} are the posterior probabilities of a visual observation m and of an auditory observation k to belong to the speaker n (or to be an outlier), given the observations:

$$\alpha_{mn} = P(A_m = n | \mathbf{f}_m) \quad \text{and} \quad \alpha'_{mn} = P(A'_k = n | g_k)$$

- Can be computed explicitly using Bayes' formula

Maximizing the expectation

- Dependency of the target function on θ is expressed through non-linear mappings \mathcal{F} and \mathcal{G}
- Maximization problem does not yield a closed-form solution as in the standard EM algorithm
- Conjugate EM algorithm: method of generations combining local ascent and global search
- Several possibilities to speed-up
 - compute local Lipschitz constants for gradient for local ascent
 - sample manifolds $\mathcal{F}^{-1}(\cdot)$ and $\mathcal{G}^{-1}(\cdot)$ for global search

Interpreting the algorithm's output

- Object locations and their audio-visual characteristics:

$$\boldsymbol{\theta} = \{s_1, \dots, s_N, \Sigma_1, \dots, \Sigma_N, \sigma_1, \dots, \sigma_N\}$$

- Auditory activity estimation:

- assign every audio observation to an object or to an outlier class using the MAP estimate:

g_k was generated by η'_k , where $\eta'_k = \arg \max_{n=1, \dots, N+1} \alpha'_{kn}$

- for each object make a decision on auditory activity based on the number of assigned observations

Implementation

- Verified on CAVA database:
http://perception.inrialpes.fr/CAVA_Dataset
- Each sequence of input data is split into time-intervals, each interval is 1/8 seconds long
- There are 1000 visual observations and 10 auditory observations per time-interval
- Next interval EM iteration is initialized with the results of the previous iteration

Examples

Detection, Localization and Tracking of 3D Audio-Visual Objects

V.Khalidov, F.Forbes, M.Hansard,
E.Arnaud & R.Horaud

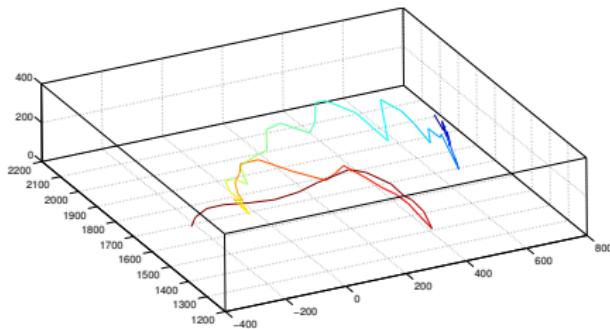
SGSSRF'09, Hirschegg

Qualitative results

- Auditory activity estimation statistics:

	<i>time-int</i>	<i>AV-int</i>	<i>AV-OK</i>	<i>AV-missed</i>	<i>AV-false</i>
M1	166	89	75	0.16	0.14
TTOS1	76	69	60	0.13	0.43

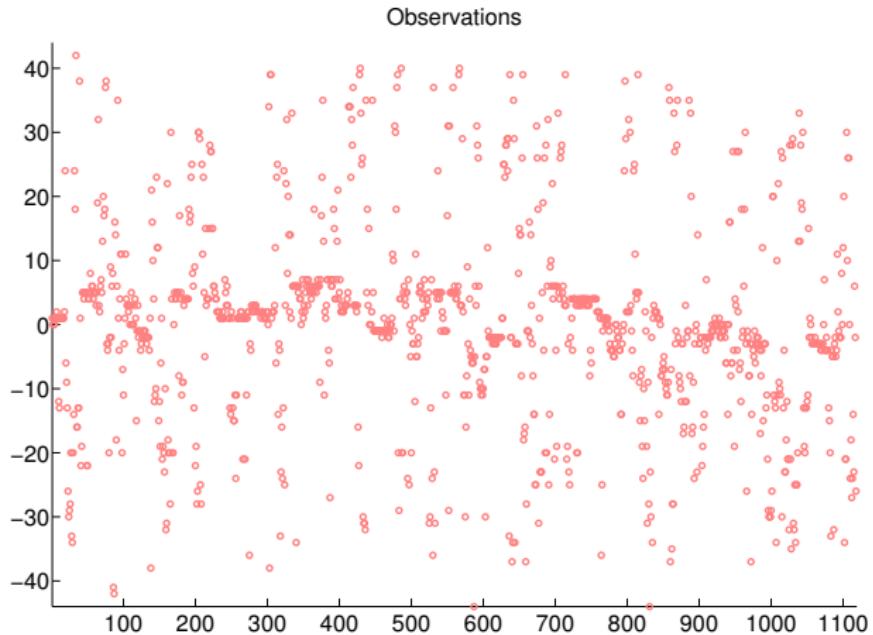
- Localization accuracy is $\pm 10\text{cm}$



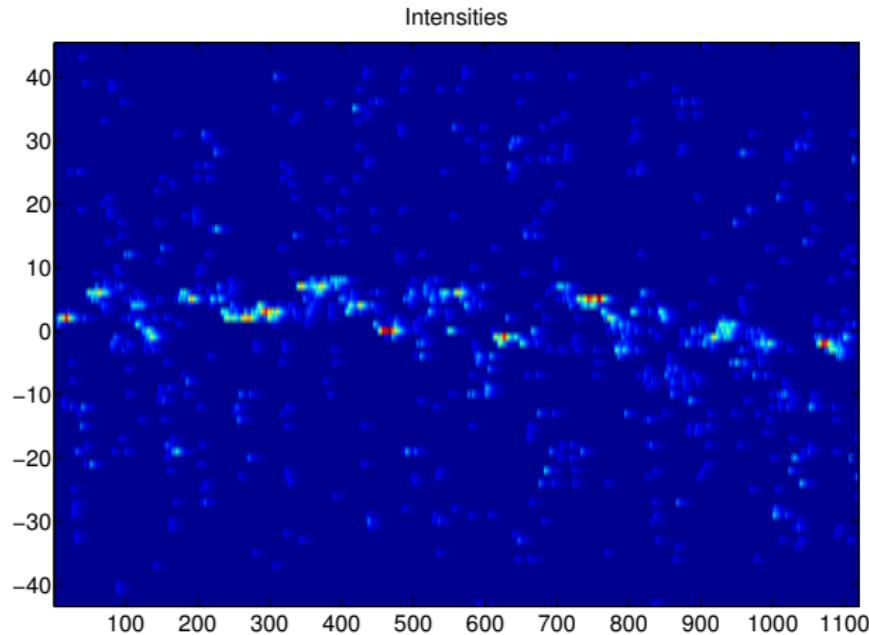
Major drawbacks of the approach

- Dynamic cues (motion vectors) are not included into the model
- Temporal properties of observation streams are ignored
- Possibilities to follow dynamic configuration changes are limited
- Background model is weak, assignments are incorrect

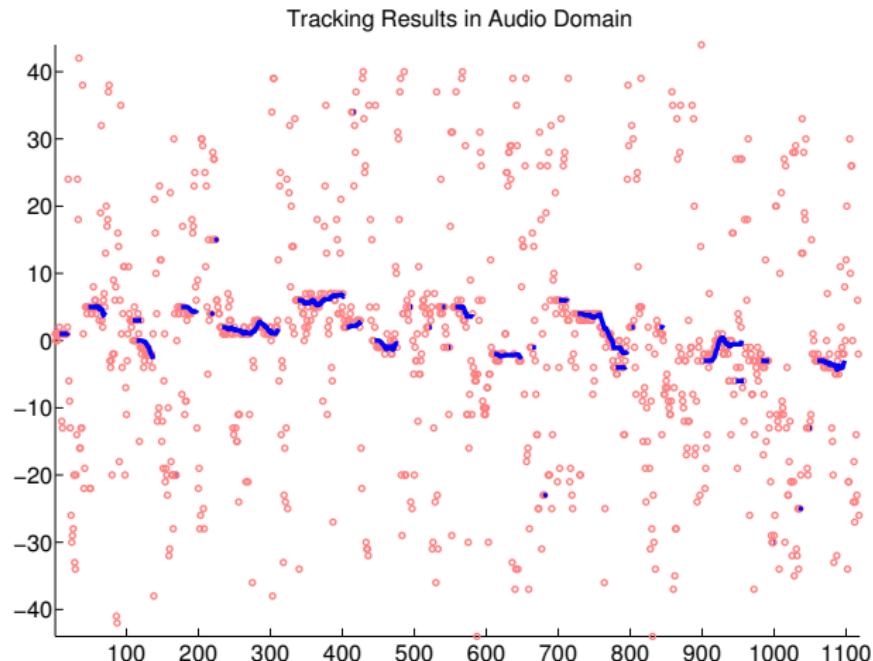
Audio stream intensity



Audio stream intensity



Audio stream intensity



Video observations density

Right camera image:

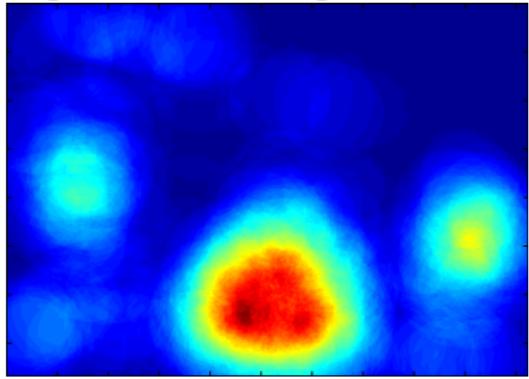


Left camera image:

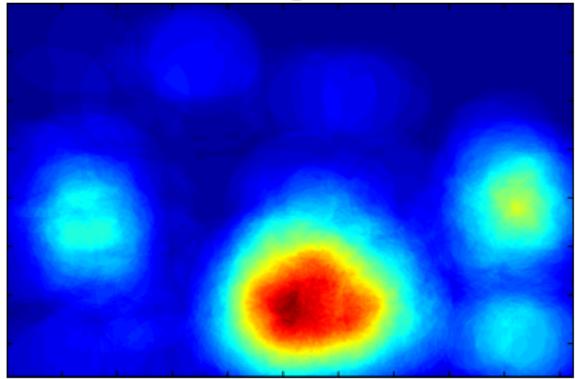


Video observations density

Right camera image:



Left camera image:



Ideas to improve the model

- Audio observations as a point process in ITD space
- Monocular video observations as point processes in image spaces
- Audio-visual object dynamics: diffusion in the parameter space
- Background is stable, motion vector is **0**
- Learning background depth map

Object dynamics

- State evolution

$$\begin{aligned} ds_t &= \mathbf{w}(t)dt + \sigma_s(t)dB_t \\ d\mathbf{w}_t &= \sigma_w(t)dB_t \end{aligned}$$

- Binocular observations at time t

$$f_m = \begin{cases} \mathcal{F}(s_m) + \sigma_f W, & \text{if computed correctly} \\ \mathcal{U}(\mathbb{F}), & \text{otherwise} \end{cases}$$

- Audio observations at time t

$$g_k = \begin{cases} \mathcal{G}(s) + \sigma_g W, & \text{if computed correctly} \\ \mathcal{U}(\mathbb{G}), & \text{otherwise} \end{cases}$$

Image space processes

- Image points should follow projected dynamics equations

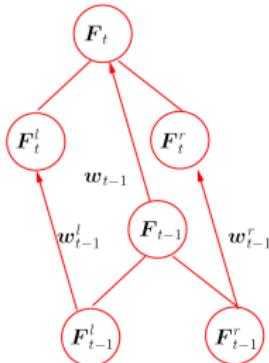
$$\begin{aligned} d\mathbf{f}_t^l &= \mathbf{w}^l(t)dt + \sigma_{\text{f}^l}(t)dB_t \\ d\mathbf{w}_t^l &= \sigma_{\text{w}^l}(t)dB_t \end{aligned}$$

- Points configuration is defined by a pairwise interaction Markov point process:

$$P(\mathbf{f}^l | \mathbf{a}^l; \nu) \propto \mu^{M^l} \prod_{\substack{i \sim j \\ a^l}} h(\|\mathbf{f}_i^l - \mathbf{f}_j^l\|)$$
$$i \underset{a^l}{\sim} j \iff a_i^l = a_j^l \text{ and } \|\mathbf{f}_i^l - \mathbf{f}_j^l\| \leq r$$

Model consistency conditions (1)

- Redundancy in process definition



- Set of *error flags* $\{e, e^l, e^r\}$ to define whether to use 3D motion model, projected motion model or the point process
- One more set of hidden variables to estimate

Model consistency conditions (2)

- Joint distributions of observation assignments $\{a, a^l, a^r\}$ and error flags $\{e, e^l, e^r\}$ are highly important
- Examples of rules:

$$P(e_m | a_{\tau^l(m)}^l, a_{\tau^r(m)}^r) = \begin{cases} p^{\text{MA}} 1_{\{e_m=1\}} + (1 - p^{\text{MA}}) 1_{\{e_m=0\}}, & \text{if } a_{\tau^l(m)}^l = a_{\tau^r(m)}^r \\ 1_{\{e_m=1\}}, & \text{otherwise.} \end{cases}$$

$$P(A_m = n | e_m, a_{\tau^l(m)}^l, a_{\tau^r(m)}^r) = \begin{cases} 1_{\{a_{\tau^l(m)}^l = n\}}, & \text{if } a_{\tau^l(m)}^l = a_{\tau^r(m)}^r \text{ and } e_m = 0 \\ \frac{1}{N+1}, & \text{otherwise.} \end{cases}$$

Inference for conjugate particle tracking model

- Observations

$$\mathfrak{o} = \{\mathbf{f}, \mathbf{f}^l, \mathbf{f}^r, \mathbf{w}, \mathbf{w}^l, \mathbf{w}^r, \mathbf{g}\}$$

- Model parameters

$$\theta = \{\text{positions, motion vectors, (co)variances, priors}\}$$

- Hidden variables

$$\mathcal{A} = \{\mathbf{A}, \mathbf{A}^l, \mathbf{A}^r, \mathbf{A}'\} \text{ and } \mathcal{E} = \{\mathbf{E}, \mathbf{E}^l, \mathbf{E}^r\}$$

- Cannot apply the EM algorithm directly, neither of the steps has a closed form solution

Variational EM approximation

- Consider restricted class of distributions $\tilde{\mathcal{D}}$:

$$q_{A^l}(\mathbf{a}^l | \mathbf{f}^l) = \prod_{m^l=1}^{M^l} q_{A_{m^l}^l}(a_{m^l}^l | \mathbf{f}^l) \quad \text{and} \quad q_{A^r}(\mathbf{a}^r | \mathbf{f}^r) = \prod_{m^r=1}^{M^r} q_{A_{m^r}^r}(a_{m^r}^r | \mathbf{f}^r)$$

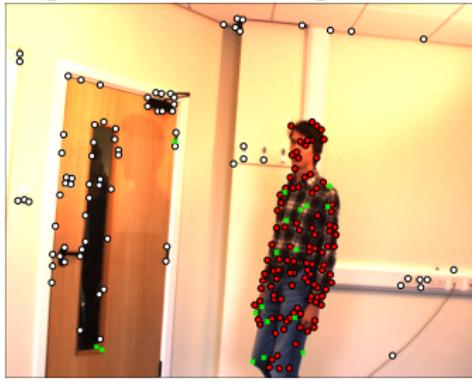
- VE-step: find distribution $q_{\mathcal{A}, \mathcal{E}}^{\eta+1}$ such that

$$q_{\mathcal{A}, \mathcal{E}}^{\eta+1} = \arg \max_{q_{\mathcal{A}, \mathcal{E}} \in \tilde{\mathcal{D}}} \mathbb{E}_{q_{\mathcal{A}, \mathcal{E}}} \{ \log P(\mathbf{o}, \mathcal{A}, \mathcal{E}; \boldsymbol{\theta}^\eta) \} + I[q_{\mathcal{A}, \mathcal{E}}],$$

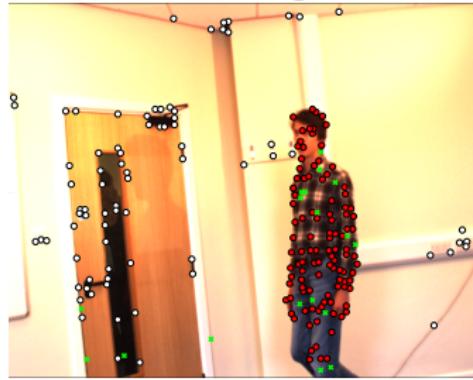
- M-step: use the method of generations to find position parameters s_1, \dots, s_N

Experimental results (1) - TTOS1

Right camera image:

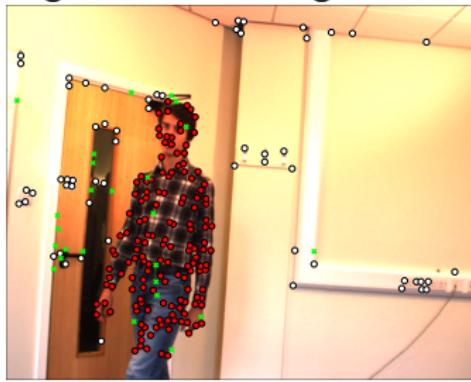


Left camera image:

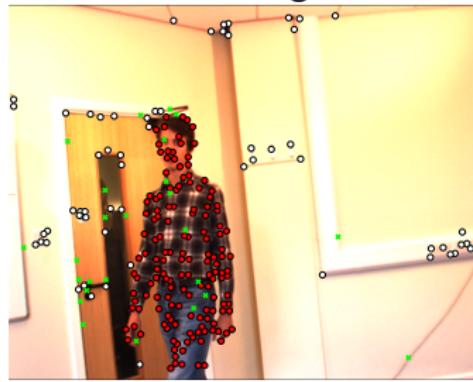


Experimental results (1) - TTOS1

Right camera image:

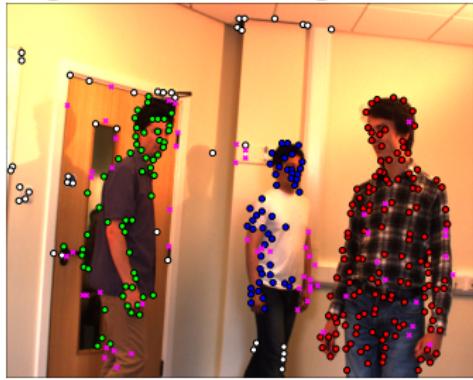


Left camera image:

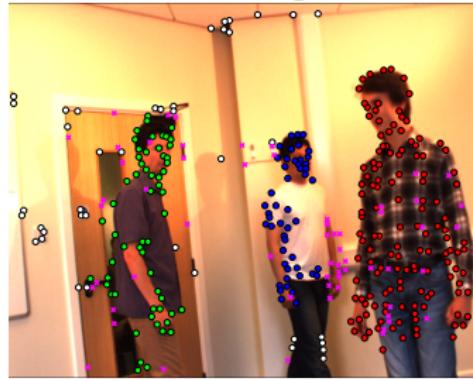


Experimental results (2) - CTMS3

Right camera image:

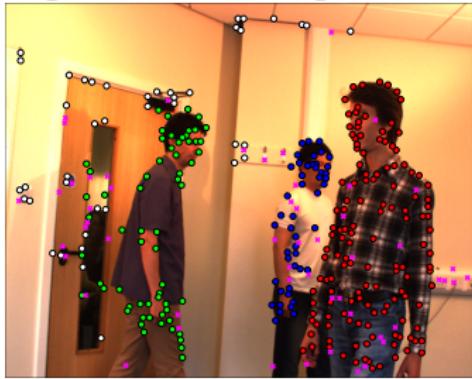


Left camera image:

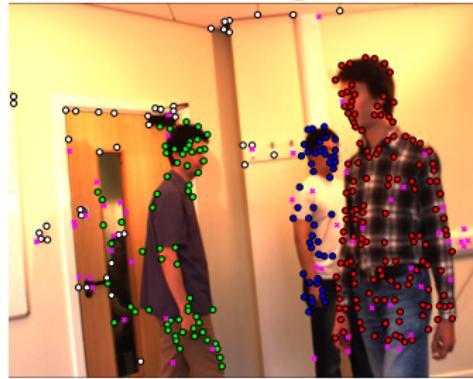


Experimental results (2) - CTMS3

Right camera image:



Left camera image:



Conclusions and future work

- We introduced a framework for audio-visual fusion using spatial and temporal coherence.
- The method is based on unsupervised clustering using the variational EM algorithm.
- Unified approach: the same model for detection, localization and auditory activity estimation.
- Possibility to automatically correct erroneous observations.
- A more sophisticated model should include eye and head motions, i.e., active perception and attention.

Thank you!

Vassil Khalidov <vasil.khalidov@inrialpes.fr>