

Help me to help you: how to learn intentions, actions and plans

H. Khambhaita and **G-J. Kruijff**

Language Technology Lab,
DFKI GmbH,
D-66123 Saarbruecken, Germany

S.R. Fanello and **M. Gianni** and **P. Papadakis** and **F. Pirri** and **M. Pizzoli** and **A. Rudi**

ALCOR, Cognitive Robotics Lab,
Sapienza, University of Rome, DIS
00186, Rome, Italy

M. Mancas

University of Mons,
F.P.Ms/IT Research Center/TCTS Lab
31, Bd. Dolez, 7000 Mons, Belgium

Abstract

The collaboration between a human and a robot is here understood as a learning process mediated by the instructor prompt behaviours and the apprentice collecting information from them to learn a plan. The instructor wears the Gaze Machine, a wearable device gathering and conveying visual and audio input from the instructor while executing a task.

The robot, on the other hand, is eager to learn both the best sequence of actions, their timing and how they interlace. The cross relation among actions is specified both in terms of time intervals for their execution, and in terms of location in space to cope with the instruction interaction with people and objects in the scene. We outline this process: how to transform the rich information delivered by the Gaze Machine into a plan. Specifically, how to obtain a map of the instructor positions and his gaze position, via visual slam and gaze fixations; further, how to obtain an action map from the running commentaries and the topological maps and, finally, how to obtain a temporal net of the relevant actions that have been extracted. The learned structure is then managed by the flexible time paradigm of flexible planning in the Situation Calculus for execution monitoring and plan generation.

Introduction

In this paper we outline a collaboration model between human-robot in which the final goal is to learn the best actions needed to achieve the required goals (in this case, reporting hazards due to a crash accident in a tunnel, identifying the status of victims and, possibly, rescuing them). The collaboration is here viewed as a learning process involving the extraction of the correct information from the instructor behaviours. The instructor communicate his actions both visually and with the aid of his comments delivered while executing the actions.

In particular, actions and intentions are obtained by elaborating on the instructor path, while inspecting the accident



Figure 1: The instructor Salvo Candela with the Gaze Machine

place, what he¹ looks at, together with his running commentaries recorded via the Gaze Machine (GM). The GM (early described in (MP08) and in (BPC07)), worn by the instructor, is a complex wearable device, illustrated Figure 1 and Figure 3² that allows to gather several perceptual data from a subject executing a task.

The extraordinary vantage point obtained by the Gaze Machine enables an agent to observe, at any time step, not only what effectively the tutor is doing and communicating it but also how the tutor adapts his behaviours, by instantiating with common sense the prescribed laws, that is, those usually regulating his conduct in similar circumstances. It allows to get his intentions, tracking the relationship between saccades and motion towards a direction, namely something interesting in the scene. Finally, by the joint localisation of the instructor's gaze, his current position and his running commentaries and the noise in the scene, it is possible to infer affordances, namely a well defined sequence of the preferred interactions between the instructor and the surroundings.

From these extremely rich source of information an agent is in the condition of learning a well temporised sequence of actions and thus, to generate a suitable plan, in order to correctly operate in a difficult and hazardous environment.

In this paper we describe at a very general level, the following aspects of this learning and generation process:

1. We define two paths, the instructor path in the scene obtained by visually localising the instructor via the gaze machine (Section 1) and the instructor gaze path, obtained via the stereo pairs mounted on the GM and the cameras staring at the eye pupils.
2. From the two paths and a suitable segmentation and clustering of motion directions (of both body and head), both the motion and vision actions are obtained and labelled. On the other hand, as the instructor actively (and benignly) comments his behaviours, all the manipulation actions are identified by the the running commentaries and the association of head motions and body position (Section 2). Indeed, actions are defined as processes with a start and an end action, and with time varying.
3. A plan library of possible activities and affordances, according to the context, is defined a priori with the contribution of the instructor. In particular, the “what to do in such a situation” can be earlier formulated. According to the prior plan library and the effective sequence accomplished, following the instructor behaviours induced by *common sense* a flexible plan of action processes is generated, where the timelines are settled according to the flexible instantiation provided by the difference between coded rules and common sense (Section 3).

A schema of the model is given in Figure2.

¹A effective fire fighter instructor

²See also <http://www.dis.uniroma1.it/~alcor/site/index.php/research/the-gaze-machine.html>

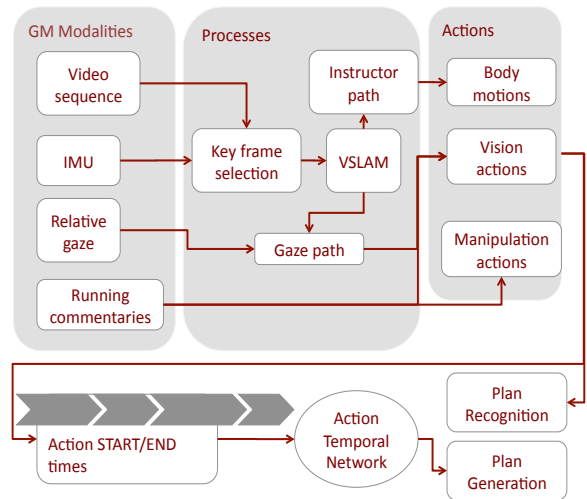


Figure 2: Schema of the flow of information and processing to learn actions from the collaboration instructor-robot, starting from the gaze machine.

The problem of inferring a plan from the observations of actions, in the context of knowledge representation, is called *plan recognition*, and it has been earlier introduced by (SSG78; KA86; Kau87). For a review of the consistency based and probabilistic based approaches to plan recognition see (AA07). Geib (Gei09) introduced a method of plan-recognition where plan-library is first converted to a lexicon similar to that used in combinatory categorical grammar. By this way author is able to introduce concept of headedness, which avoids early commitments to plan and goal hypothesis in the process of plan-recognition, which eventually results in increased speed of the plan-recognition system. On the other hand in the realm of learning and computer vision the analogous concepts of acting based on observations have been specified as *action recognition*, *imitation learning* or *affordances learning*, as mainly motivated by the neurophysiological studies of Rizzolatti and colleagues (PGF⁺92; GFFR96) and by Gibson (Gib77; Gib55). Reviews on action recognition are given in (MHK06; Pop10; AC99) and on learning by imitation in (ACVB09; SIB09).

The two approaches have, however, evolved in completely different directions. Plan recognition assumed actions to be already given and represented, in so being concerned only in the technical problems of generating a plan, taking into account specific preferences and user choices, and possibly interpreting plan recognition in terms of theory of explanations (CG93). On the other hand action recognition and imitation learning has been more and more concerned with the robot ability to capture the real and effective sequence and to adapt it to changing contexts. As noted by Krüger and colleagues in (KKG07) the terms action and intent recognition, in plan recognition, often obscure the real task achieved by these approaches. In fact, as far as plan recognition assumes an already defined set of actions the observation process is purely indexical. On the other hand the difficulties with the learn-



Figure 3: The instructor while is rescuing a victim.

ing by imitation and action recognition approaches is that they lack important concepts such as execution monitoring, intention recognition and plan generation.

Our contribution fosters a more tight integration between the two approaches wherein the actions are segmented via the Gaze Machine and the instructor running commentary and the consequent plan recognition that is based on these actions.

1. Visual Localization

Two paths can be obtained by the instructor running in the disaster theatre. The first concerns the position of his body and the second the position of his gaze, not mentioning the position and direction of his head obtained via the inertial sensor placed on the GM.

For the instructor position we relied on both the extended Kalman filter (EKF)-base visual slam introduced by (DRMS07) and the particle filter ones introduced by (PC06), but suitably extended to cope with the specific head motions and consequent change blindness (SL98), and the advantages of the calibrated stereo pairs. The challenges that we face in this procedure stem from the inherent scene peculiarities of rescue environments as well as the loosely constrained movement of the camera setup which follows the movement of the instructor's head. In detail, the scene characteristics of a rescue environment include a wide range of lighting conditions and a plurality of solid but also non-solid obstacles (such as smoke). The position-orientation of the camera setup is also highly variable as the instructor rushes within the accident area due to looming hazards. We can note that because most of the computational effort is carried out off-line, we can take advantage of the techniques developed in the context of *Structure and Motion* recovery in order to deal with the higher variance in the camera motion and particular lighting conditions. The selection of a stable sequence of frames, that turn out to be key frames not only for the localisation and mapping process but also for action segmentation is a crucial step. This is more deeply discussed in Section, see Figure 4. Furthermore, we rely on well known methods for feature extraction and optical flow to predict the displacement of the tracked features and bundle adjustment between pairs of stereo images for motion estimation. In SAM problems, 3D structure is used to estimate

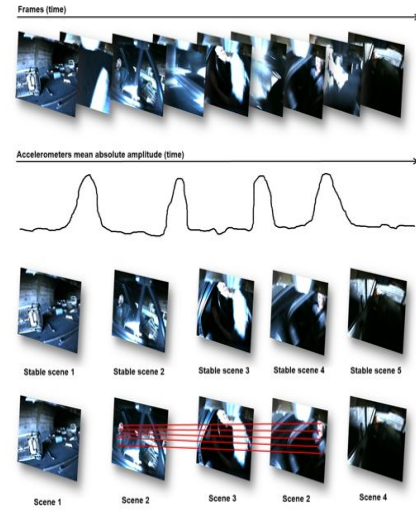


Figure 4: Key frames extraction. First row: acquired frames through time. Second row: acquired accelerometer absolute mean amplitude through time. Third row: frames corresponding to accelerometers peaks (movement from a LOI to another) are discarded. Forth row: features are extracted by the different key frames. If a lot of features match between the key frames of different scenes, this means that those scenes are the same and thus they are grouped together under the same label.

the camera pose by resectioning. Thus, the computation of the motion is also complemented by the usage of dense disparity maps. The process goes through three main steps:

1. build a Viewing Graph (Fig. 5) and compute the Essential Matrices (HZ03) between each pair of views, obtained from the key frames set;
2. given the estimated position at time t , factorise the Essential Matrices to produce observations for an EKF, which provides the current position at time t' ;
3. bundle adjust among the estimated sequence of 3D structure and camera motion.

Steps 1-2 provide a local consistency between different temporal frames. On the other hand they do not take into account sudden movements, which are filtered out in the key frame selection. In order to maintain a global consistency a bundle adjustment step is required where the re-projection error is minimised. Using the above described visual-based SLAM we are able to obtain an estimate of the instructor's path which, in turn, is used to derive the gaze path within the scene. It is interesting to note that due to inhibition of return, typical of the gaze when a salient feature come up hiding previous saliency levels, often a large amount of images are required in order to effectively track features. The viewing graph will tell on which configurations it is possible to rely in so avoiding the constraints induced by a sequence of pairs of images. Given the position of the instructor the localisation of his gaze is immediately obtained by the stereo pair.

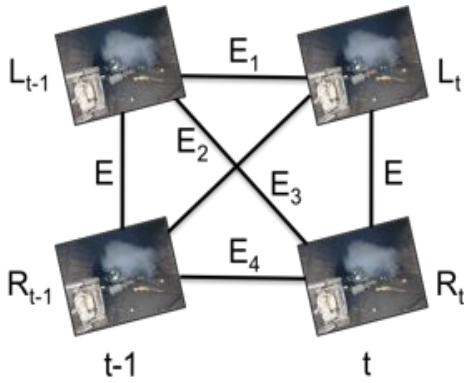


Figure 5: Example of the Viewing Graph used for the visual localisation. L_t and R_t are respectively the left and the right scene cameras at time t . E_i is the Essential Matrix between different views.

In the following section the two paths are going to be segmented according to the recognition of actions from (i) the running commentary and (ii) the video sequence. The recognition of the actions will in turn enable to infer the spatiotemporal information of an action. Indeed, the two paths prove to be essential for action segmentation as they can correctly specify the where and when an action is performed as well as the corresponding spatiotemporal information of the instructor's gaze: what and when the instructor is gazing at, during a particular action.

2. Segmentation and Action Maps

In this section we discuss how we can segment the data acquired using the Gaze Machine to obtain a sequence of performed actions. We shall also discuss the intention recognition via the coup d'oeil, i.e. how it is possible to extract the instructor's intention on the basis of his fixations and the spoken running commentaries.

A library of possible activities and affordances has been compiled in advance with the contribution of the instructor but, due to the high changeability of the scenario, the instructor will not follow a predefined, prioritised sequence of actions. The decision on what to do next is taken on the run, according to the task (i.e. plan the rescue) and the affordances characterising the scenario. The current instructor's intention involves what he is actually able to capture via attention. A saccade that is directed toward a location that is not involved in the current action may indicate a shift in the instructor's attention; depending on the associated saliency, this may or may not fire a head movement. However, also the information provided by the peripheral vision is enough to increase the situation awareness and take decisions. We, thus, introduce the concept of *coup d'oeil* to refer to those time instants in which something in the peripheral view fires a running commentary reporting something relevant in the scene.

The Gaze Machine records the instructor's saccade sequence by tracking his gaze in space. This is accomplished by projecting in the 3D scene the estimated point of regard. Scene structure is recovered via the Gaze Machine stereo rig while both pupils are tracked to extract visual axes (see Figure 1). A first kernel-based segmentation is performed to extract the fixation scan path from the acquired sequence of 3D points of regard. The main problem we address in this step is taking into account the instructor 3D position, as the 3D fixated points changes if the instructor moves.

The segmentation of the image flow acquired from the experienced firefighter is needed as a prior to further analysis of his actions. Key frame selection has been thoroughly investigated in the context of SAM recovery (TFZ98; PGGV⁺04). In this paper we face the problem in the case of wearable cameras and unpredictable human motions. When performing some activity, a person is acquiring information (by gazing) in some important location in order to perform actions and then he moves to another location of interest (LOI). During the movement between two LOIs the acquired images are of little interest as they are most of the time fuzzy and very unstable. Moreover, the extensive visual disruptions due to the firefighter fast motion imply a high probability of change blindness (SL98) which decrease again the usability of the gaze data acquired during those periods. It is thus important to discard the frames which are recorded during the LOI change in order to extract the more stable scenes (Figure 4, third row). Finally, the firefighter can move from one LOI to another and then come back to the first LOI, or he can also be disturbed by some important bottom-up distractor which makes him turn his head and then he can look again to the previous scene. This shows that two stable scenes are not necessary different scenes or LOIs (Figure 4, forth row).

A two-step approach can be used to extract meaningful scenes or key frames from the video flow: first the data from the accelerometer can provide cues on the head stability and then computer vision techniques are able to recognise already seen scenes or novel scenes. Figure 4 illustrates the process. The top row shows the successive frames through time. The shape of the absolute mean amplitude of the accelerometers located in the gaze machine is presented on the second row and shows picks during the firefighter movements between two LOIs and valleys during his stay in the same LOI. By discarding the frames which correspond with the accelerometer peaks, it is possible to keep only the stable scenes. Feature extraction and matching between those different scenes provide information to group together the scenes which are the same. If the features extracted from some key frames of one scene match a number of features above a given threshold on some key frames from another scene, this means that the two scenes are the same as it can be seen on Figure 4, forth row. In that case the two scenes are labelled with the same label.

Along with the instructor changes in position, pose and the running commentaries, 3D fixations are used to detect the starting/ending of an action. We are interested in producing a Map of basic actions, divided in 1) body motion, 2) vision actions and 3) manipulation actions, labelled with the corre-



Figure 6: Fixations from the tunnel sequence labelled by the instructor running commentaries: these are examples of key frames used for action segmentation and to define compatibilities; the third figure above induce the constraint $lookingAt(victim,t)$ **during** $openingDoor(car,t')$.

spondent starting/ending time.

Actions related to body motions are segmented on the basis of the instructor position in the 3D scene. Vision actions involve the generation of a sequence of fixations and are detected by clustering in time and space and recognising special sequences in the running commentaries (i.e. *I see...*). A coupe d’oeil belongs to the vision actions category and is detected making use of the special sequences in the running commentary and, when significant, sudden changes in inertial measurements. Indeed, the coupe d’oeil involves saccades followed by head movements and changes in body directions. For the detection of the manipulation actions we completely rely on the running commentary, as the scene cameras on the GM don’t provide a good point of view for gesture recognition.

From the Action Map we can define the compatibility conditions for generating a flexible plan. We assume we are given a plan library from the usual instruction on behaviours, and that this plan library includes affordances, given a specific rescue situation. Our aim is to show within the map the common sense raising from the choice of an action according to the urgency of a decision.

3. From Actions to flexible plan and plan recognition

The Actions map is constituted by a timeline indicating the time stamp of each action, the temporal relations among actions and the spatial cluster they belong to. The spatial cluster is obtained by the instructor path (see Section 2). Using the rules specified in the plan library and the Action Map the in-

structor plan execution can be suitably labelled for planning.

For example, according to the plan recognition algorithm of (Gei09), and using the specified plan library, we first generate a combinatory categorical grammar (CCG) type plan-lexicon which maps observations to CCG categories. The algorithm results in a set of *explanations*, mentioning a goal and an ordered sequence of actions. The choice of assigning categories to observations is made according to specified *headedness* value. *Headedness* is a powerful method of controlling the space of possible explanations to be considered during the plan-recognition procedure.

In any case we mainly base the mapping from the Action Map to a possible plans via a temporal network compiling constraints and compatibilities within the Situation Calculus. Temporal relations specify how activities, such as *looking at a victim* and of *opening the car door* are correlate along time. For modelling both temporal constraints and cause-effect relations between activities we adopt, in fact, Temporal Flexible Situation Calculus(FP05), accommodating Allen temporal intervals, multiple timelines among actions and concurrent situations. It intermediates between Situation Calculus formulae and temporal constraint networks. For example, the temporal relations illustrated in Figure 6, first row, last image, can be expressed by the compatibilities

$$T_c = [comp(lookingAt(victim,t), [[(during, openingDoor(car))]])]$$

Here the compatibility states that the activities *look a victim*, involving vision actions, and *opening the car door* have to be performed according to the *during* temporal relation. The temporal network associated with the compatibilities T_c is rep-

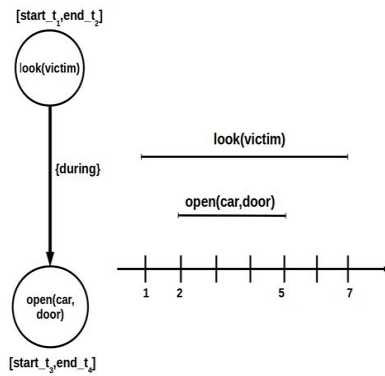


Figure 7: Temporal constraint network represented by T_c along the timelines $do([start_{look}(victim, t_1), end_{look}(victim, t_2), S_0])$ and $do([start_{open}(car, door, t_3), end_{open}(car, door, t_4)], S_0)$

resented in Figure 7. Therefore a way to generate a plan is to exploit the obtained temporal network and the flexible plan in the Situation Calculus.

Conclusion

In this work we have described a new framework for the collaboration between a human and a robot based on a wearable device, the Gaze Machine. This device creates a strong communication between the human, in this case an instructor, and the robot, by allowing the agent to look straightly into the perceptual flux of the companion. We have described how to process this perceptual information in order to obtain an Action Map. The Action Map is a rich labelled graph, starting from which it is possible to use specific methods, such as the transformation from a temporal network to a flexible plan and plan recognition, to generate a plan for the robot to correctly explore the environment.

Acknowledgements

This paper describes research done under the EU-FP7 ICT 247870 *NIFTI* project.

References

Marcelo Gabriel Armentano and Analía Amandi. Plan recognition for interface agents. *Artif. Intell. Rev.*, 28(2):131–162, 2007.

J. K. Aggarwal and Q. Cai. Human motion analysis: A review. *Computer Vision and Image Understanding*, 73:428–440, 1999.

Brenna D. Argall, Sonia Chernova, Manuela Veloso, and Brett Browning. A survey of robot learning from demonstration. *Robot. Auton. Syst.*, 57(5):469–483, 2009.

Anna Belardinelli, Fiora Pirri, and Andrea Carbone. Bottom-up gaze shifts and fixations learning by imitation. *IEEE Transactions on Systems, Man, and Cybernetics, Part B*, 37(2):256–271, 2007.

Eugene Charniak and Robert P. Goldman. A bayesian model of plan recognition. *Artif. Intell.*, 64(1):53–79, 1993.

Andrew J. Davison, Ian D. Reid, Nicholas D. Molton, and Olivier Stasse. Monoslam: Real-time single camera slam. *IEEE Trans. Pattern Analysis and Machine Intelligence*, 29:2007, 2007.

Alberto Finzi and Fiora Pirri. Representing flexible temporal behaviors in the situation calculus. In *Proceedings of IJCAI-2005*, pages 436–441, 2005.

C. Geib. Delaying commitment in plan recognition using combinatorial categorial grammars. In *Proceedings of the International Joint Conference on Artificial Intelligence (IJCAI) 2009*, pages 1702–1707, 2009.

V. Gallese, L. Fadiga, L. Fogassi, and G. Rizzolatti. Action recognition in the premotor cortex. *Brain*, 119:593–609, 1996.

J.J. Gibson. Perceptual learning: differentiation or enrichment? *Psych. Rev.*, 62:32–41, 1955.

J.J. Gibson. The theory of affordances. In R. Shaw and J. Bransford, editors, *Perceiving, Acting, and Knowing: Toward an Ecological Psychology*, pages 67–82. Hillsdale, NJ: Lawrence Erlbaum, 1977.

Richard Hartley and Andrew Zisserman. *Multiple View Geometry in Computer Vision*. Cambridge University Press, New York, NY, USA, 2003.

Henry A. Kautz and James F. Allen. Generalized plan recognition. In *AAAI*, pages 32–37, 1986.

Henry A. Kautz. *A formal theory of plan recognition*. PhD thesis, Department of Computer Science, University of Rochester, 1987.

Volker Krüger, Danica Kragic, and Christopher Geib. The meaning of action a review on action recognition and mapping. *Advanced Robotics*, 21:1473–1501, 2007.

Thomas B. Moeslund, Adrian Hilton, and Volker Krüger. A survey of advances in vision-based human motion capture and analysis. *Computer Vision and Image Understanding*, 104(2-3):90–126, 2006.

Stefano Marra and Fiora Pirri. Eyes and cameras calibration for 3d world gaze detection. In *ICVS*, pages 216–227, 2008.

M. Pupilli and A. Calway. Real-time visual slam with resilience to erratic motion. *Computer Vision and Pattern Recognition, 2006 IEEE Computer Society Conference on DOI - 10.1109/CVPR.2006.240*, 1:1244–1249, 2006.

G. Di Pellegrino, V. Gallese, L. Fadiga, L. Fogassi, and G. Rizzolatti. Understanding motor events: a neurophysiological study. *Exp. Brain Research*, 91:176–180, 1992.

Ronald Poppe. A survey on vision-based human action recognition. *Image and Vision Computing*, 28:976–990, 2010.

Marc Pollefeys, Luc Van Gool, Maarten Vergauwen, Frank Verbiest, Kurt Cornelis, Jan Tops, and Reinhard Koch. Visual modeling with a hand-held camera. *Int. J. Comput. Vision*, 59(3):207–232, 2004.

Stefan Schaal, Auke Ijspeert, and Aude Billard. Computational approaches to motor learning by imitation. *Philosophical Trans. of the Royal Soc. B: Biological Sciences*, 358(1431):537–547, 2009.

Daniel J. Simons and Daniel T. Levin. Failure to detect changes to people during a real-world interaction. *Psychonomic Bulletin and Review*, 5:644–649, 1998.

Charles F. Schmidt, N. S. Sridharan, and John L. Goodson. The plan recognition problem: An intersection of psychology and artificial intelligence. *Artif. Intell.*, 11(1-2):45–83, 1978.

P. Torr, A.W. Fitzgibbon, and A. Zisserman. Maintaining multiple motion model hypotheses over many views to recover matching and structure. pages 485–491, jan. 1998.