

A Developmental Roadmap for Task Learning by Imitation in Humanoid Robots

Baltazar's Story

Manuel Lopes^{†,‡}

Alexandre Bernardino[†]

José Santos-Victor[†]

[†]Instituto de Sistemas e Robótica, Instituto Superior Técnico, Lisboa, Portugal

[‡]Escola Superior de Tecnologia, Instituto Politécnico de Setúbal, Setúbal, Portugal

<http://vislab.isr.ist.utl.pt>

{macl,jasv}@isr.ist.utl.pt

Abstract

We present a strategy whereby a robot follows a developmental pathway to (i) explore its own visuo-motor capabilities, (ii) understand its surrounding environment, (iii) become aware of people acting in the environment and finally imitate observed actions. We describe some results of the different developmental stages, involving perceptual and motor skills, implemented in our humanoid robot, Baltazar. In addition to the overall system, another important contribution is the use of a two-phase, uncalibrated algorithm for object grasping. The last phase is driven under closed-loop vision based control, where the Jacobian is learned online.

1 Introduction

“Friendly” and social interaction between robots and humans is a grand challenge for robotics. Due to the diversity of actions/tasks to be performed and the range of possible interactions with objects and humans, it would be impractical (if not impossible) to explicitly pre-program a robot with such capabilities. Instead, such systems must be able to learn by themselves what tasks to execute and how they should be performed, which requires sophisticated motor, perceptual and cognitive skills.

To address these challenges, we (as well as other researchers) adopt two fundamental metaphors: (i) learning by imitation as a powerful means to teach a complex humanoid-like (social) robot and (ii) a developmental approach that can balance the complexity of the system at the various levels of functional performance.

Learning by imitation is likely to become the primary form of teaching such social, cognitive robots. Let us consider a system able to learn how to solve some tasks by imitation, e.g. by observing a human manipulating a set of objects. This problem of skill transfer has three major difficulties: (i) how to gather task-relevant information? (ii) how to convert the data that are valid for a human for a robot? and (iii) how to infer the important parts of the demonstration (e.g. “understand” the task).

Several approaches have been adopted to gather the



Figure 1: Baltazar. A 14 degrees of freedom humanoid torso.

information for imitation. Schaal et al. (2003) use an exoskeleton to capture kinematic data. Oztop and Arbib (2002) rely on some marks to get visual features for hand detection and grasping, in the context of imitation and modeling of the Mirror neurons. Lopes and Santos-Victor (2003) exploit task-contextual bias to modulate the information extraction process. Imitation and skill transfer between systems with different bodies (kinematics, dynamics and skills) was addressed by Nehaniv and Dautenhahn (1998) using an algebraic formulation (bodies with different skills were considered). For the case of a humanoid robot, Nakaoka et al. (2003) introduce adaptation of the trajectories to be able to guarantee the correct balance during task execution.

Kuniyoshi et al. (1994) proposed one of the first

works in imitation, a system able to learn how to imitate an assembly task by extracting a hierarchical description of the task. Billard et al. (2004) address the problem of inferring the important parts of the task by casting it into an optimization framework. Zöllner and Dillmann (2003) present a system where two hand tasks are imitated, using information about the functionality of each object and handling temporal task restrictions, in a symbolically manner.

Even if imitation can allow a robot to learn an extremely large variety of tasks, it is clear that it requires the robot to have several sophisticated motor, perceptual and cognitive capabilities. Hence, building such complex skills can become an overwhelming task in itself. For learning one particular skill, many other systems may need to be present and their inter-connections properly established.

The developmental perspective, as proposed by e.g. Weng (1998), is a new paradigm aiming at overcoming this complexity problem, of learning and properly integrating many perceptual, motor or cognitive skills.

The robot should “start” with a minimal subset of core capabilities (as newborns do) to bootstrap learning mechanisms that, through self-experimentation and interaction with the environment and other humans, would progressively lead to the acquisition of new skills, adapted to particular contexts, and having the system integrating all the learning methods internally. Metta (1999) used a developmental approach for a robot that successively acquired vergence, saccade and vestibular control, as well as head-arm coordination.

Amongst the capabilities required to interacting with objects, understand their spatial configuration and learning by imitation, perception is perhaps the most fundamental. They allow gathering (task or contextually relevant) information and training samples for all other forms of learning. Then, some motor capabilities need to be in place before the robot can start interacting with the world and providing “calibration” information for other modules (e.g. relating depth information from vergence with arm length).

The development of imitation capabilities requires an appropriate definition of the sequence of learning steps to reach that goal, as well as adequate performance evaluation methods to decide when to switch to higher developmental levels. In other words, it is important to define the overall hierarchy of developmental stages and the skills that must be acquired at each level. Table 1 shows the structure we adopt for the main developmental stages the robot (or a human infant Arsénio (2004); Natale (2004)) will go

through: (i) Learning about the self; (ii) Learning about objects and the world and (iii) Learning about others and imitation.

For each stage in this “developmental pathway”, we show the set of skills to be acquired, and the time line explaining the restrictions governing the system. We do not distinguish between innate versus learned behavior in biological systems (“the nature versus nurture” question). Instead, we just request all the modules to be present before the system can develop to the next level.

Table 1: Developmental pathway for the Perceptual and Motor capabilities (in *italic* the modules that are learnt by the robot)

Time line	Perceptual/Motor Capabilities
↓ self-awareness	eye vergence random movements <i>Arm-head</i> coordination near-space mapping
↓ world-awareness	near-space mapping <i>visually initiated reaching</i> <i>visual control of grasp</i>
↓ imitation	detection of other’s actions imitation of tasks

In the first developmental level, the robot acquires very simple and yet crucial capabilities: vergence control and object foveation/tracking. Then, by executing random arm movements, in a self exploratory mode, it begins to coordinate head and arm configurations, by creating a arm-head map. This map is accurate enough to allow for reaching and grasping objects in easy positions.

In the second developmental stage, the robot builds a map of the surrounding area (object positions and identification). Driven by attentional cues, the robot engages in more challenging grasping tasks, for which the previously learned arm-head map is not sufficiently accurate. For that reason we propose a novel method for visually controlled grasping, which improves over time and ensures the necessary robustness.

Previous approaches for object grasping were either completely visual controlled Kragic et al. (2002), with problems in guaranteeing the presence of the hand in the visual field, or completely open-loop Natale (2004) with no capability of error correction. Instead, we combine the two modalities, with an open-loop phase moving the effector to the field of view followed by a closed-loop method with the precision necessary to put the effector in contact with the object.

At the final developmental stage, the presence of a demonstrator will elicit a task imitation behavior, that will decompose the actions and then replicate with a given metric. For this purpose, the system must be able to decompose the observed action into the relevant key elementary actions that must be executed for performing a task.

To conclude, our main contributions are two-fold. On one hand we present a developmental strategy for humanoid robots, according to widely accepted stages in developmental psychology. On the other hand, we propose a visually guided grasping process, where learning is driven by the motivation to precisely grasp objects, that continuously adapts over time (open ended learning).

All experiments in the paper were implemented in our humanoid robot, Baltazar, equipped with a 4 *dof* binocular head, a 6 *dof* arm and an 11 *dof* (under-actuated by 4 motors) hand. The robot is shown in Figure 1 and described in detail in Lopes et al. (2004).

In Section 2 we present the development of self awareness. Section 3 deals with the understanding of the world and the interaction with objects. Imitation learning of tasks is presented in Section 4. Conclusions and future work are drawn in Section 5.

2 Self-Awareness

Humans take a long time before becoming self-sufficient. Knowing how to walk, how to recognize objects, understanding how to solve a task, interacting with objects, are very difficult tasks, and only after several years all mechanism necessary are available and reliable. Becoming aware of its own body and then start to coordinate it is the first step to survival. Infants have several mechanism guiding its development.

For the case of the head-eye system, voluntary control appears very early. Some reflexive movements are evident from birth (head-righting reflex Payne and Isaacs (1999)) but voluntary control becomes apparent only at the end of the first month. A five-month old child already shows a good control. This control of the head will enable the tuning of the vision system to start looking at (and understanding) the environment. In van der Meer et al. (1995) there is a discussion about the significance of neonate's, apparently, random arm movements.

Several reflexes allow newborns to look to their hand. The "Asymmetric Tonic Neck Reflex" can be elicited when the baby is prone or supine. When the head is turned to one side or the other, the limbs on the face side extend while the limbs on the opposite

side flex. This reflex is believed to facilitate the development of an awareness of both sides of the body as well as help develop eye-hand coordination.

The interaction between eyes and hand is very important. This interaction will allow the newborn to tune its eyes, distinguish depth and recognize touchable objects. For a baby exploring the hand, how it moves and how it looks will be the most interesting thing in the first few months. Learning to make it do what it wants to do, will be a very complex learning task. In the end the reward will be tremendous: being able to predict hand movements and to touch objects.

In this section we present the capabilities allowing the system to be aware of its own body and to learn how to coordinate it.

2.1 Near-Space perception

The near-space contains the touchable objects and our own body. Being able to understand what happens there is fundamental. Bernardino and Santos-Victor (2002) suggest a method where the disparity between images is used, together with some neuronal-based filters, to segment objects at different depths. The head can be moved to look toward the hand using disparity as a feedback signal to control it. Figure 3 shows a result of verging on an object. This same mechanism will later be used to map object positions.

2.2 Arm-Head Coordination

Many tasks need a very fine Arm-head coordination. Object manipulation is only possible with precise visual control of the hand. In order to coordinate Arm positions with Head position, we are going to create an *Arm-Head Map*. This map is bidirectional. If the head position is fixed moving the arm to the mapped position puts the hand in the fovea of the two eyes. If the arm is fixed, we can visually locate it by moving the head to the mapped position.

Several approaches can be used to learn this map. In Lopes and Santos-Victor (2003), a neural network was used to map from arm feature points to joint angles. In D'Souza et al. (2001), a very powerful method is used to learn inverse kinematics of a humanoid. Vijayakumar and Schaal. (2000) created a method, *Locally Weighted Projection Regression*, that will be used for learning the map. This method is linear with the number of samples and every new sample can be added easily. As it is not capable of extrapolating, the working space must be well covered in the training set.

The data set is gathered by self-observation. The arm is moved around in the space, while the hand is

tracked and foveated. Figure 2 shows the hand being moved to the front of the eyes by using the *Head-Arm map*. The quality of this map is good enough to guarantee that the hand is always in the image but not in the fovea. In our experiments, the average error is about 5cm, corresponding to 15% of the image.

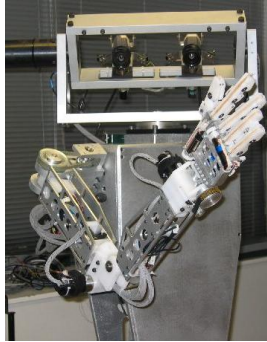


Figure 2: Head-Hand Coordination

This map will enable the system to reach and, in special cases, grasp objects. This will be very motivating and in the next level object grasping will develop further.

2.3 Attention Mechanism

When looking around us some objects attract more our attention than others. This is related to the urgency of each task and to reduce the amount of information to process. Context may influence the attention drawn by some objects, (e.g. food when hungry). In our approach, the attentional process depends on the developmental stage. In the beginning, the hand is the main focus of attention, facilitating the learning of the Head-Arm Map. Also salient objects in the scene attract the system's attention in a bottom-up process. Later on, in the second stage, Baltazar will search around him, pay attention to all objects, one at a time and create a map of the nearby area. In the final stage, attention will be driven toward the person doing the demonstration and the manipulated objects. In these later stages, attention becomes gradually more driven in a top-down, context and task dependent manner.

3 World Awareness

As the robot gains control over its own perceptual and motor capabilities, it gets more and more interested in exploring its surrounding world. This exploratory

motivation will call for the development of more advanced manipulative capabilities as opposed to the rudimentary skills available during phase one.

For object grasping, it is necessary to have several motor programs: the arm must be able to approach the object (reaching) before finally grasping it, the hand must be able to have a stable grasp and pre-shaping can be necessary for faster movements or moving objects. However at this stage all the robot can do is to fixate at salient objects and approach them in a primitive form of grasping. The development path will require the following new skills:

1. detect object's positions in the nearby space and store this information in some sort of representation (near space map).
2. learn how to reach objects in a controlled manner, using visual feedback, and grasp them.

This section describes algorithms that solve all these steps, allowing a robot to move on to the next developmental level, where it gains awareness of others (humans or robots) and the actions they perform. In addition to the reaching step based on the *Head-Arm Map* presented in Section 2, we propose a new algorithm to grasp objects based on visual servoing techniques estimated online.

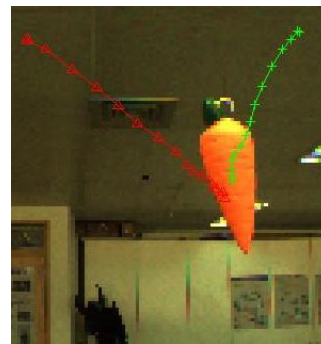


Figure 3: Verging on an object. Left (Δ) and right (+) eye.

3.1 Near-Space (Objects) Mapping

There is neurological evidence of spatial aware neurons that are active when movement or objects are present near the skin Rizzolatti et al. (1977). It is also known in developmental psychology that infants became aware of the near and far space very early. It is very useful to know where an object is and whether it can be grasped or not. After all the time spent interacting with its own hand, the system can already

distinguish objects at different depths and search for the desired one.

By this exploratory behaviour, we create a map of the localization of objects around us - the peripersonal map - through various steps:

1. Find an object in the visual space
2. Foveate on this object
3. Memorize the object position in head (proprioceptive) coordinates (Θ_{Head}).

Through exploration, the robot thus creates a mental image of the surrounding space. The position of the objects are memorized in terms of head (proprioceptive) coordinates. In Figure 4 Baltazar is searching for “fruits” around him where different objects are assumed to have different colours.

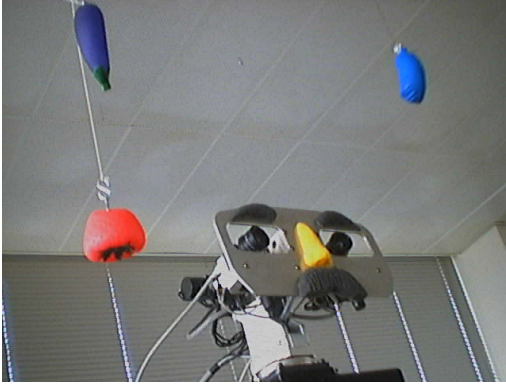


Figure 4: Mapping object positions in head coordinates.

3.2 Object Grasping - a two step approach

Infants start reaching objects without any visual feedback. The movement is only initiated with vision but not guided throughout the entire action. In case of failure, the movement restarts from the beginning.

At the first stage of development, the estimated *Arm-Head map* allows the system to (crudely) move the hand towards an object. Hence, if a simple trajectory is followed, the hand may well succeed in touching the object. The problem with this (open-loop) approach is the absence of a mechanism for error correction. This is the reason why babies in this phase restart the grasp quite often, instead of correcting it Payne and Isaacs (1999).

The second stage of object reaching relies on visual feedback, coping with the problem of error correction. The *Head-Arm map* is used to move the hand

to the objects vicinity. Then, accurate positioning is achieved by visual guidance in closed loop. With this phase, it is possible to grasp objects in a reflex type manner, the hand closing after touch.

The method presented in D’Souza et al. (2001) could be used here. Their approach consists in mapping motor positions and velocities to image velocities, using a very strong statistical learning approach, yielding good results. The disadvantages arise from the lack of extrapolation capabilities and by not having an explicit Jacobian estimation, thus needing more time to gather the information, and preventing the use of well studied visual servoing control algorithms.

We adopted a visual servoing perspective, described by e.g. Hutchinson et al. (1996). However, although it is possible to solve the problem with an algebraic formulation, we adopted a model-less way, as it allows the system to learn and develop from its own experience. A particularly useful method for on-line estimation of visual motor relations is presented by Jaegersand (1996). The image *Jacobian* (J) relating image changes ($\Delta\mathbf{y}$) caused by motor movements ($\Delta\theta$), can be interactively estimated by:

$$\hat{J}(t+1) = \hat{J}(t) + \alpha \frac{(\Delta\mathbf{y} - \hat{J}(t)\Delta\theta) \Delta\theta^T}{\Delta\theta^T \Delta\theta}$$

where α denotes the Jacobian update rate. To move the system to the desired image position y^* , we apply the following control law:

$$\Delta\theta = g(J^+(\mathbf{y}^* - \mathbf{y}))$$

where J^+ represents the pseudo-inverse of J and the function $g(\cdot)$ can be chosen to have an exponential, linear or any other type of convergence.

In order to deal with a larger workspace and to incorporate some open-loop movements, we had to improve the existing algorithm. More details can be found in the Appendix A. Figure 5 shows the resulting behavior of the system while grasping objects. The hand is closed after sensing the contact with the object. The capability of pre-shaping the hand will only develop at a later stage. For small grasping velocities, this type of movements can be sufficient, but bigger velocities will require learning some form of pre-shaping and predicting the time of contact with the object.

At this point in development, the system can not only control its own body and perceptual abilities but also perform relatively complex manipulation tasks, memorize objects spatial configurations, search for objects, etc. It is then ready to start looking at humans or other robots and the tasks they perform.

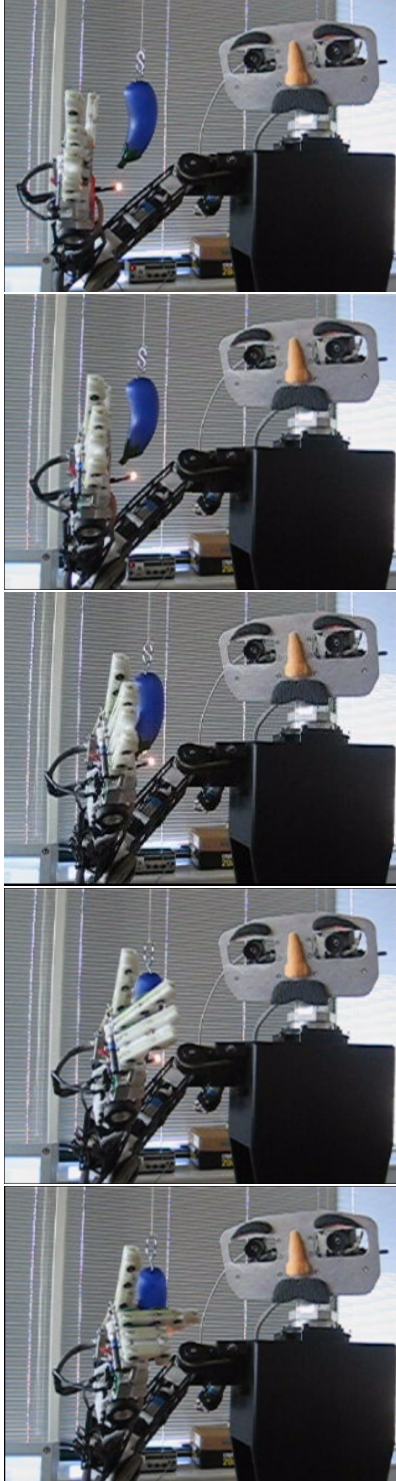


Figure 5: Several frames in the sequence from the initial position resulting from the *Head-Arm Map*, then the visual guided part and finally the object grasping.

4 Imitation

Figure 6 shows an example of a task being executed. It consists of picking up some objects and moving them around. To imitate this task, the robot will first need to understand the spatial relations of objects around the demonstrator (understand the far space). Then, understanding the near space becomes fundamental to establish correspondence between the demonstrator perspective its own (self) viewpoint (i.e. the blue object is on the left of the demonstrator, but it is in front of me). After the observation of the demonstration movements, the important task moments must be extracted and segmented. Finally the task is repeated by the robot, using the task description and all the modules previously learned. The following sections will provide details on the different modules developed at this stage.

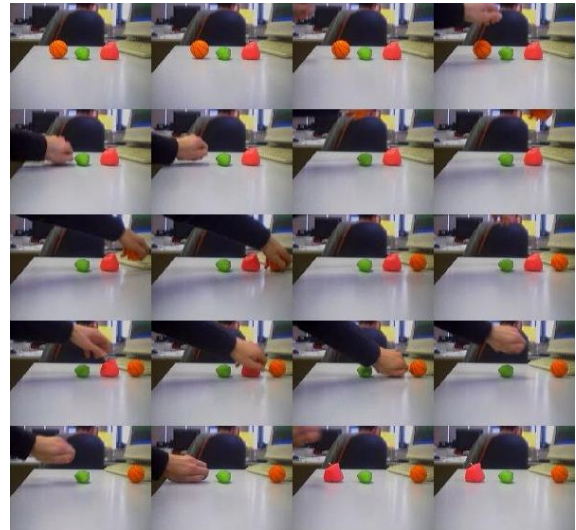


Figure 6: Several frames of the task demonstration.

4.1 Far-Space Interpretation

Understanding events and object's localizations at far distances (i.e. more than the arm can reach) is different from mapping our surrounding space. The frame of reference will no longer be our own body, instead we describe object's positions relative to another person, this is specially useful for imitation learning. Object's position will be codified in terms of allo-coordinates. Some simplifying assumptions can be made about depth in order to reduce the complexity of scene reconstruction.

4.2 Task Segmentation

The actions and movements of the demonstrator must be segmented and codified in a way useful for imitation. We developed a method consisting in a multiple object tracking and a task point detector. When doing manipulation our hand will occlude objects very frequently. Grasping and releasing can be very difficult to detect. Being the hand the only actuator enables the usage of information to deal with occlusions. Every object can have three movement models: static, moving and being moved. When an object is moving its velocity profile can be predicted with Newtonian dynamics, when being moved it has the same velocity as the hand. The algorithm will mark every point in the trajectories of the objects that satisfy the following constraints: all objects are static, the hand is not moving and the hand is not occluding any object.

The task is then codified by having objects with their physical properties (shape and color) and their spatial relations (A between B and C; A right of B or A left of B).

The complete sequence shown in Figure 6 has 234 frames, this sequence was processed online and the task points, shown in Figure 7, were automatically extracted. We can see that the system succeeds in detecting what frames are important to describe the task.

4.3 Imitation

As mentioned in Gergely et al. (2002), imitation goals are not always very clear. In our case the imitation task will proceed in order to have the same spatial relations. In case the demonstrator has made a movement and there is no difference in the ordering of objects (Figure 7), the robot will mimic the absolute spatial positions. We can see that all the modules developed until this point were essential to be able to replicate the task at hand.

5 Conclusions/Future Work

We have presented a developmental route for a humanoid robot¹ to acquire increasingly more complex skills.

The robot first learns about its own body and surrounding environment. All information is gathered by self-exploration. The quality of the Arm-head coordination achieved in this phase is sufficiently good to ensure that the hand always remains in the image and

¹see <http://vislab.isr.ist.utl.pt/baltazar> for videos showing the experiments in this work

that objects can be grasped in simple cases. In a second phase, motivated to further interact with objects, the system develops a closed-loop control behavior capable of precise grasping. It also creates a map of the interesting objects in the surrounding space. In the final developmental phase, people acting in the environment are the major source of information. The observed tasks are segmented in special points in order to finally imitate the task.

The developmental pathway allows the robot to acquire new skills on top of the existing (learned) capabilities. We described results of the various developmental stages of the system: the vergence and object tracking system, the learning of the Arm-head map, the visually initiated object grasping system and a new solution to visually guide grasping. The method consists in two phases: an open-loop controller putting the hand close to the object, and a closed-loop vision-based controller for precisely touching the object. This method does not need calibration and can be learned on-line in a very efficient way. In the future, we will focus our efforts on the aspects of learning the interaction between people and objects.

A Visual Grasp

In this section we present a generalization of the method suggested by Jaegersand (1996), to be used to visually control the arm. The image Jacobian (J) relating image changes ($\Delta\mathbf{y}$) caused by motor movements ($\Delta\theta$), can be iteratively estimated by:

$$\hat{J}(t+1) = \hat{J}(t) + \alpha \frac{(\Delta\mathbf{y} - \hat{J}(t)\Delta\theta) \Delta\theta^T}{\Delta\theta^T \Delta\theta}$$

where α is the Jacobian update rate. To move the system to the desired image position y^* , we can apply the following control law:

$$\Delta\theta = g(J^+(y^* - y))$$

where J^+ represents the pseudo-inverse of J and $g(\cdot)$ can be chosen to have an exponential, linear or any other kind of convergence.

When the working volume is very large the Jacobian can no longer be accurately estimated with only one linear model. To solve this we propose a new method. With only one linear model the update mechanism must be fast enough to have an accurate model for each region. In the case of open-loop movements the system can no longer update the model and a specific model for the new region must

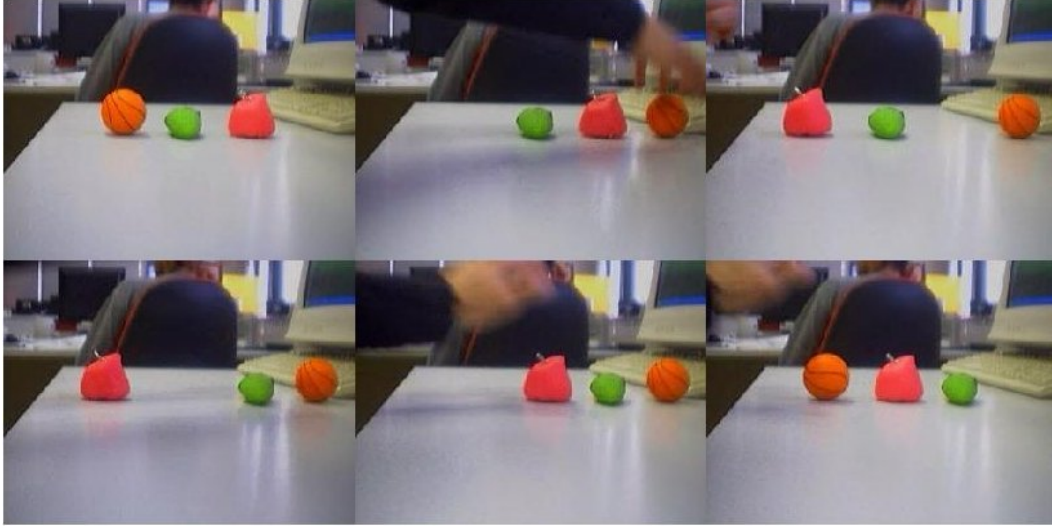


Figure 7: Segmentation of a task. Notice that from the third to the fourth image there is no difference in the ordering of the object, just their absolute distances. These relevant points were extract online from a video sequence with 234 frames.

already be present. The workspace should be partitioned in several regions, R_i , $i = 1 \dots N$. At each instant the distance c is measure between the current position and all the regions, the selected Jacobian J is the one corresponding to the nearest area R_i . We use a Mahalanobis distance with covariance D . The covariance can be updated online to reduce the number of regions and to better adjust the linear model to the non-linear system. Trying to update the regions center creates problems by overlapping regions and with region transitions.

The Jacobian update rate (α) should be larger when the model is inaccurate and then reduced to improve convergence. One measure to access the model quality (mq) can be:

$$mq(t) = mq(t-1) + \gamma \langle \Delta y, J_k \Delta \theta \rangle$$

γ is a decaying factor and $\langle \cdot \rangle$ represents internal product. mq is positive when the observed movements has a direction error less than 90 degrees.

The regions centers x_i may correspond to motor features $x = \theta$, visual features $x = y$ or a combination of them. With visual features there is the possibility of doing planning in visual space but there are different motor positions that give the same visual features and should have different linearizations.

Table 2 presents the complete algorithm for doing the visual controlled grasp.

J^+ must be carefully implemented. As some directions are not observed, the Jacobian inversion will be very unstable. To solve this problem the pseudo

Table 2: Uncalibrated Visual Servoing Algorithm

To move the system to the desired image position y^*

1. Choose the region R_i corresponding to the actual state x :

$$c_i = (x - x_i)^T D_i (x - x_i)$$

$$R_i : \min_i c_i$$

if $\max c_i < C$ create a new area l with $x_l = x$, $D_l = D$ and $J_l = J_i$. Choose $R_i = R_l$.

2. apply the control law:

$$\Delta \theta = K_i \frac{J_i^+ (y^* - y)}{\|J_i^+ (y^* - y)\|}$$

3. observe image changes Δy
4. make the update to the model i corresponding to position x with:

$$\hat{J}_i = \hat{J}_i + \alpha_i \frac{(\Delta y - \hat{J}_i \Delta \theta) \Delta \theta^T}{\Delta \theta^T \Delta \theta}$$

5. if $|y^* - y| > E$ goto 1

inverse is implemented with a SVD method and any singular values less than 10% of the larger are treated as zero.

Chaumette (1998) show some problems present in Visual servoing methods. Our method solves the problem of the Jacobian derivation and the calibration of the robot and cameras. In general these methods are sensitive to initial positions, being prone to fall in local minima but, in our approach, the system always starts near the final position due to the *Head-Arm map*, thus making convergence easier.

We made several experiments to assess the quality of the resulting algorithm. Our system measures a specific dot in the hand with two cameras giving an image position of the hand (u_l, v_l) for the left eye and (u_r, v_r) for the right eye. The features are calculated as follows:

$$\mathbf{y} = \begin{bmatrix} \frac{u_l + u_r}{2} \\ \frac{v_l + v_r}{2} \\ u_l - u_r \end{bmatrix}$$

This gives position and distance information estimation of the hand related to the head. The head was maintained fixed and four arm joints were used. The distance between the central point of each zone was 10 *degrees*. The Jacobian update rate was equal in all regions and chosen as $\alpha = 0.1$ while $mq < 0$ and $\alpha = 0.01$ while $mq > 0$.

Figure 5 shows some quantitative results of the grasp sequence shown in Figure 5 using our proposed algorithm. The hand was positioned near the object using the *Head-Arm map*. The resulting error corresponds to about 8 *cm*. The associated image error is corrected in the final phase (visually controlled) with a linear convergence rate.

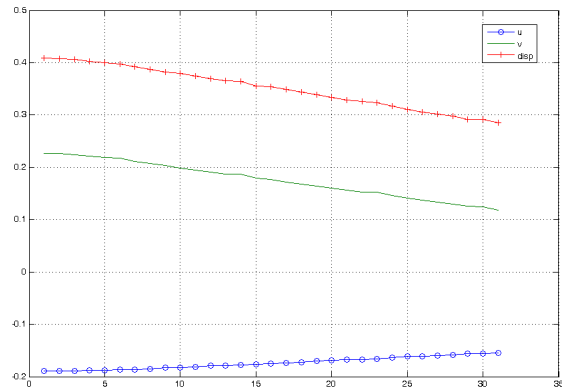
Acknowledgments

Work partially supported by: EU Proj. (IST-004370) RobotCub <http://www.robotcub.net> and by the FCT POSI-Program in the frame of QCA III.

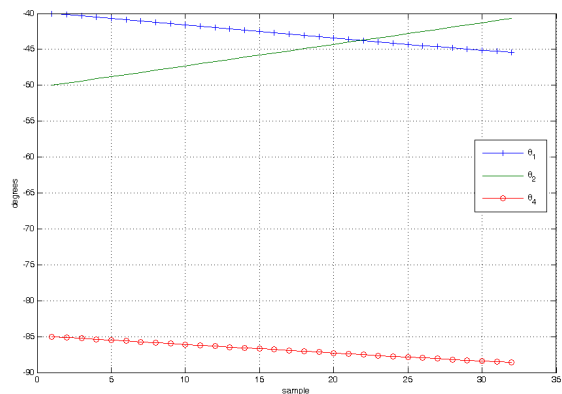
References

Artur Arsénio. *Cognitive-Developmental Learning for a Humanoid Robot: A Caregiver's Gift*. PhD thesis, Massachusetts Institute of Technology, Boston, USA, Sept 2004.

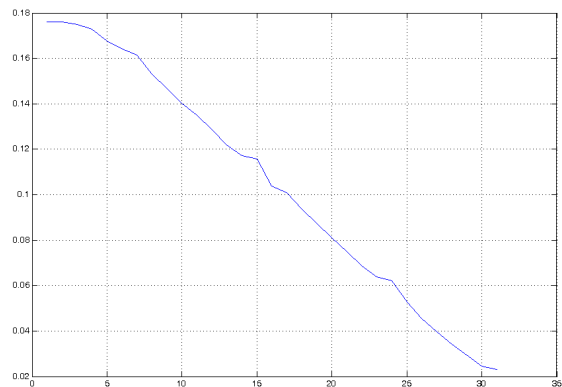
Alexandre Bernardino and José Santos-Victor. A binocular stereo algorithm for log-polar foveated systems. In *Biological Motivated Computer Vision*, Tuebingen, Germany, 2002.



(a) Visual Features (normalized pixel values)



(b) Joint Angles (degrees)



(c) Error

Figure 8: Servoing results for object grasping.

- A. Billard, Y. Epars, S. Calinon, G. Cheng, and S. Schaal. Discovering optimal imitation strategies. *Robotics and Autonomous Systems*, 47:2-3, 2004.
- F. Chaumette. Potential problems of stability and convergence in image-based and position-based visual servoing. In D. Kriegman, G. Hager, and A.S. Morse, editors, *The Confluence of Vision and Control*, pages 66–78. LNCIS Series, No 237, Springer-Verlag, 1998.
- A. D’Souza, S. Vijayakumar, and S. Schaal. Learning inverse kinematics. In *IEEE International Conference on Intelligent Robots and Systems*, Maui, USA, 2001.
- György Gergely, Harold Bekkering, and Ildikó Király. Rational imitation in preverbal infants. *Nature*, 415:755, Feb 2002.
- S. Hutchinson, G. Hager, and P. Corke. A tutorial on visual servo control. *IEEE Transactions on Robotics and Automation*, 12(5):651–670, October 1996.
- M. Jaegersand. Visual servoing using trust region methods and estimation of the full coupled visual-motor jacobian. In *IASTED Applications of Control and Robotics*, Orlando, EUA, 1996.
- D. Kragic, L. Petersson, and H. I. Christensen. Visually guided manipulation tasks. *Robotics and Autonomous Systems*, 40(2-3):193–203, August 2002.
- Y. Kuniyoshi, M. Inaba, and H. Inoue. Learning by watching: Extracting reusable task knowledge from visual observation of human performance. *Trans. on Robotics and Automation*, 10(6):799–822, 1994.
- M. Lopes, R. Beira, M. Praça, and J. Santos-Victor. An anthropomorphic robot torso for imitation: design and experiments. In *International Conference on Intelligent Robots and Systems*, Sendai, Japan, 2004.
- M. Lopes and J. Santos-Victor. Visual transformations in gesture imitation: What you see is what you do. In *International Conference on Robotics and Automation*, Taipei, Taiwan, 2003.
- Giorgio Metta. *Babybot: A Study on Sensori-motor Development*. PhD thesis, University of Genova, 1999.
- S. Nakaoka, A. Nakazawa, K. Yokoi, H. Hirukawa, and K. Ikeuchi. Generating whole body motions for a biped humanoid robot from captured human dances. In *ICRA*, Taipei, Taiwan, 2003.
- Lorenzo Natale. *Linking Action to Perception in a Humanoid robot: a Developmental Approach to Grasping*. PhD thesis, University of Genova, 2004.
- C. Nehaniv and K. Dautenhahn. Mapping between dissimilar bodies: Affordances and the algebraic foundations of imitation. In *European Workshop on Learning Robots*, Edinburgh, Scotland, 1998.
- Erhan Oztop and Michael A. Arbib. Schema design and implementation of the grasp-related mirror neuron system. *Biological Cybernetics*, 87: 116–140, 2002.
- V. Gregory Payne and Larry D. Isaacs. *Human Motor Development: a Lifespan Approach*. Mayfield Publishing Company, California, USA, 4th edition, 1999.
- G. Rizzolatti, L. Fadiga, L. Fogassi, and V. Gallese. The space around us. *Science*, 277:190–191, 1977.
- S. Schaal, A. Ijspeert, and A. Billard. Computational approaches to motor learning by imitation. *Phil. Trans. of the Royal Society of London: Series B, Biological Sciences*, 358:537–547, 2003.
- A. van der Meer, F. van der Weel, and D. Lee. The functional significance of arm movements in neonates. *Science*, 267(5198):693–695, Feb 1995.
- S. Vijayakumar and S. Schaal. Locally weighted projection regression: An $o(n)$ algorithm for incremental real time learning in high dimensional spaces. In *ICML*, Stanford, USA, 2000.
- Juyang Weng. The developmental approach to intelligent robots. In *AAAI Spring Symposium Series, Integrating Robotic Research: Taking The Next Leap*, Stanford, USA, Mar 1998.
- R. Zöllner and R. Dillmann. Using multiple probabilistic hypothesis for programming one and two hand manipulation by demonstration. In *IROS*, Las Vegas, USA, 2003.