



THE UNIVERSITY *of* EDINBURGH

## Edinburgh Research Explorer

### **Structural Bootstrapping - A Novel, Generative Mechanism for Faster and More Efficient Acquisition of Action-Knowledge**

**Citation for published version:**

Wörgötter, F, Geib, CW, Tamosiunaite, M, Aksoy, EE, Piater, JH, Xiong, H, Ude, A, Nemec, B, Kraft, D, Krüger, N, Wächter, M & Asfour, T 2015, 'Structural Bootstrapping - A Novel, Generative Mechanism for Faster and More Efficient Acquisition of Action-Knowledge', *IEEE Transactions on Autonomous Mental Development*, vol. 7, no. 2, pp. 140-154. <https://doi.org/10.1109/TAMD.2015.2427233>

**Digital Object Identifier (DOI):**

[10.1109/TAMD.2015.2427233](https://doi.org/10.1109/TAMD.2015.2427233)

**Link:**

[Link to publication record in Edinburgh Research Explorer](#)

**Document Version:**

Peer reviewed version

**Published In:**

IEEE Transactions on Autonomous Mental Development

**General rights**

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

**Take down policy**

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact [openaccess@ed.ac.uk](mailto:openaccess@ed.ac.uk) providing details, and we will remove access to the work immediately and investigate your claim.



# Structural bootstrapping - A novel, generative mechanism for faster and more efficient acquisition of action-knowledge

Florentin Wörgötter<sup>a,i</sup>, Chris Geib<sup>b,c,i</sup>, Minija Tamosiunaite<sup>a,d,i</sup>, Eren Erdal Aksoy<sup>a</sup>, Justus Piater<sup>e</sup>, Hanchen Xiong<sup>e</sup>, Ales Ude<sup>f</sup>, Bojan Nemec<sup>f</sup>, Dirk Kraft<sup>g</sup>, Norbert Krüger<sup>g</sup>, Mirko Wächter<sup>h</sup>, Tamim Asfour<sup>h</sup>

<sup>a</sup>*Georg-August-Universität Göttingen, Bernstein Center for Computational Neuroscience, Department for Computational Neuroscience, III Physikalisches Institut - Biophysik, Göttingen, Germany*

<sup>b</sup>*School of Informatics, Edinburgh, United Kingdom*

<sup>c</sup>*College of Computing and Informatics, Drexel University, Philadelphia, USA*

<sup>d</sup>*Department of Informatics, Vytautas Magnus University, Kaunas, Lithuania*

<sup>e</sup>*Institute of Computer Science, University of Innsbruck, Innsbruck, Austria*

<sup>f</sup>*Humanoid and Cognitive Robotics Lab, Dept. of Automatics, Biocybernetics, and Robotics, Jožef Stefan Institute, Ljubljana, Slovenia*

<sup>g</sup>*Cognitive and Applied Robotics Group, University of Southern Denmark, Odense, Denmark*

<sup>h</sup>*Institute for Anthropomatics and Robotics, Karlsruhe Institute of Technology, Karlsruhe, Germany*

<sup>i</sup>*These authors have contributed equally to this work.*

---

## Abstract

Humans, but also robots, learn to improve their behavior. Without existing knowledge, learning either needs to be explorative and, thus, slow or – to be more efficient – it needs to rely on supervision, which may not always be available. However, once some knowledge base exists an agent can make use of it to improve learning efficiency and speed. This happens for our children at the age of around three when they very quickly begin to assimilate new information by making guided guesses how this fits to their prior knowledge. This is a very efficient *generative learning mechanism* in the sense that the existing knowledge is generalized into as-yet unexplored, novel domains. So far generative learning has not been employed for robots and robot learning remains to be a slow and tedious process. The goal of the current study is to devise for the first time a general framework for a generative process that will improve learning and which can be applied at all different levels of the robot's

cognitive architecture. To this end, we introduce the concept of structural bootstrapping – borrowed and modified from child language acquisition – to define a probabilistic process that uses existing knowledge together with new observations to supplement our robot’s data-base with missing information about planning-, object-, as well as action-relevant entities. In a kitchen scenario, we use the example of making batter by pouring and mixing two components and show that the agent can efficiently acquire new knowledge about planning operators, objects as well as required motor pattern for stirring by structural bootstrapping. Some benchmarks are shown, too, that demonstrate how structural bootstrapping improves performance.

*Keywords:* Generative Model, Knowledge Acquisition, Fast Learning

---

## Introduction

It has been a puzzling question how small children at the age of three to four are suddenly able to very quickly acquire the meaning of more and more words in their native language, while at a younger age language acquisition is much slower. Two interrelated processes are being held responsible for this speeding-up. The primary process is semantic bootstrapping where the child associates meaning from observing their world with co-occurring components of sentences. For example, if the word “fill” is consistently uttered in situations where “filling” occurs, then the meaning of the word can be probabilistically guessed from having observed the corresponding action again and again [1, 2]. Once a certain amount of language has been acquired, a second process – named syntactic bootstrapping – can speed this up even more and this is achieved by exploiting structural similarity between linguistic elements. This process can take place entirely within language and happens in a purely symbolic way without influence from the world. For example, if a child knows the meaning of “fill the cup” and then hears the sentence “fill the bowl”, it can infer that a “bowl” denotes a thing that can be filled (rather than a word meaning the same thing as “fill”) without ever having seen one ([1, 3, 4, 5, 6, 7, 8, 9] see [10] for a comparison between semantic and syntactic bootstrapping). Thus, the most probable meaning of a new word is being estimated on the basis of the prior probability established by previously encountered words of the same semantic and syntactic type in similar syntactic and semantic contexts.

These two generalization mechanisms – semantic and syntactic bootstrapping – are very powerful and allow young humans to acquire language without explicit instruction. It is arguable that bootstrapping is what fuels the explosion in language and conceptual development that occurs around the third year of child development [8, 11].

In general “the trick” seems to be that the child possesses at this age already enough well-ordered knowledge (grammar, word & world knowledge) which allows him/her to perform guided inference without too many unknowns. Grammar and word-knowledge are highly structured symbolic representations and can, thus, provide a solid scaffold for the bootstrapping of language. Symbolic representations, however, do not stop short at human language. For robots, planning, planning operators, and planning languages constitute another (non-human) symbolic domain with which they need to operate. Thus, it seems relatively straightforward to transfer the idea of se-



mantic and syntactic bootstrapping to the planning domain for robot actions. The current paper will first address this problem.

The question, however, arises whether related mechanisms might also play a role for the acquisition of other, non-linguistic cognitive concepts, for example the properties of objects and tools. Briefly, if you know how to peel a potato with a knife, would there be a way to infer that a potato peeler can be used for the same purpose? This example belongs to the second set of problems addressed in this study: How can a cognitive agent infer role and use of different objects employing the knowledge of previously seen (and used) objects, how can it infer the use of movement and force patterns, etc.?

The goal of the current study is to address one complex scenario all the way from the planning-level down to sub-symbolic sensorimotor levels and implement (different) bootstrapping processes for the fast acquisition of action knowledge. The only requirement for all these different bootstrapping mechanisms is that there exists a well-structured scaffold as a basis from where on different inference processes can take place. The different scaffolds, thus, form the structures upon which bootstrapping can be built. Hence, we call these processes “structural bootstrapping”.

One can consider structural bootstrapping as a type of semi-supervised probabilistic learning, where an agent uses an internal model (scaffold) to quickly slot novel information (obtained for example by observing a human) into appropriate model categories. This is a *generative process* because existing knowledge is generalized into novel domains, which so far had not been explored. The advantage of such a bootstrapping process is that the agent will be able to very quickly perform these associations and grounding needs only to take place afterwards by experimenting in a guided way with the new piece of knowledge. Evidently, as this is based on probabilistic guesswork, bootstrapping can also lead to wrong results. Still, if the scaffold is solid enough all this can be expected to be much faster and more efficient than the much more unstructured and slow process of bottom-up exploration learning or than full-fledged learning from demonstration. Thus, structural bootstrapping is a way for the generative acquisition and extension of knowledge by which an agent can more efficient redeploy what it currently knows, but where its existing knowledge cannot be *directly* employed. The distinction between syntactic and semantic components is, however, less evident when considering structural (e.g. sensori-motor) elements. It will become clear by the examples below that structural bootstrapping often contains both aspects.

Here we will show that one can implement structural bootstrapping across different levels of our robotics architecture in the humanoid robot ARMAR-III [12, 13] trying to demonstrate that bootstrapping appears in different guises and will, thus, possibly not be limited to the case studies presented in this paper. As a major aspect, this work is meant to advocate structural bootstrapping as a way forward to a more efficient extension of robot-knowledge in the future. Early on we emphasize that the complexity of the here-shown aspects prevents exhaustive analyses. After all we are dealing with very complicated and possibly human-like cognitive generative (inference) processes for which children and adults need years of experience to reach their final efficiency.

The paper is structured in the following way. First we provide an overview of the bootstrapping idea, then we show details on the system, processes, and methods. Next we show six different types of structural bootstrapping at different levels. This will be followed by some benchmarks and a discussion section which also includes the state of the art in robot knowledge acquisition.

## Overview

The goal of this work is to use a humanoid robot (ARMAR III) to demonstrate several ways to perform structural bootstrapping at different levels of its intrinsic cognitive architecture. Thus, we define a traditional 3-layer architecture consisting of a Planning level, a Mid-level, and a Sensorimotor Level [14]. In order to perform a task, the robot first needs to make a (symbolic) plan. The mid-level acts as a symbol-to-signal mediator (explained below) and couples the planning information to the sensorimotor (signal) level. The sensorimotor level then performs execution but also sensing of the situation and the progress and potential errors of the robot’s motor actions. Details of the actual sensorimotor control loops shall be omitted here for the sake of brevity (see e.g. [14] for this).

Every layer uses different syntactic elements; for example the Planning layer uses Planning Operators. But all syntactic elements will always be embedded in their layer-specific scaffold. For the Planning layer its is the Planning Language that defines how to arrange and use the Planning Operators. Hence the Planning Language is the scaffold of this layer. Similar structural relations between syntactic elements and scaffolds are defined for the two other layers.

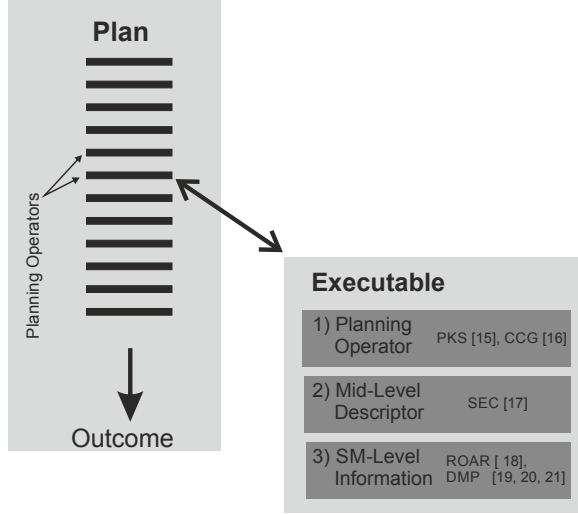


Figure 1: Structure of an Executable and its link to the robotics plan.

The general structural bootstrapping idea is now rather simple: Semantic and/or syntactic similarity at the level of the scaffold is used to infer, which (known) syntactic entities can take the role of which other (unknown, but currently observed) syntactic entities. In other words: Using the appropriate layer-specific scaffold, the robot makes inferences about the role of an observed but “incomprehensible” entity, for which the machine does not have any representation in its knowledge base. Based on these inferences the unknown entity can be replaced with one that is known (one, for which there is an entry existing in the knowledge base). This replacement will allow the machine to continue with its planned operation ideally without any additional information.

### Structures

To allow bootstrapping we need to define the actual data structures, which are used by the robot for execution of a task and which need to be set up in a way to allow for structural bootstrapping, too (Fig. 1).

At the top layer we use a conventional robotics planner [15] to create a plan for a given task. The plan consists of a sequence of Planning Operators. As such these planning operators cannot be executed by a robot. Thus, to achieve this, we define a so-called *Executable*, which consists of several components using methods from the literature:

1. a planning operator, by which the Executable is linked to the Plan [15, 16], together with its
2. mid-level descriptors [17] and
3. all perception and/or control information from the sensorimotor level for executing an action [18, 19, 20, 21].

Hence, during execution the different planning operators are called-up and each one – in turn – calls the belonging Executable, which contains the required control information to actually execute this chunk of the plan.

Some of these aspects are to some degree embodiment specific (most notably the control information), some others are not. Note, the structure of an Executable is related to the older concept of an Object-Action-Complex (OAC, [22, 23]). OACs had been defined in our earlier works as rather abstract entities [23], the Executables – as defined – here extend the OAC concept by now also including planning operators and are finally defined in a very concrete way (to actually allow for execution, which had not yet been the case for the OAC).

Essential to this work is that we use the concept of bootstrapping now in the same way at these three levels. The syntactic representations used to compute aspects of a given level are level-dependent where we have the following syntactic representatives:

1. planning operators,
2. syntactic structure of mid-level descriptors, and
3. perceptual (sensor) and control (motor) variables.

Therefore, we employ different (grammatical) scaffolds for the bootstrapping:

1. planning language,
2. semantic event chains (SECs<sup>1</sup> [24, 17]), and
3. sensorimotor feature/parameter regularity

from where the bootstrapping commences.

---

<sup>1</sup>Semantic Event Chains (SECs) encode for an action the sequence of touching and untouching events that happen until the action concludes. A more detailed description is given in the Methods section below.

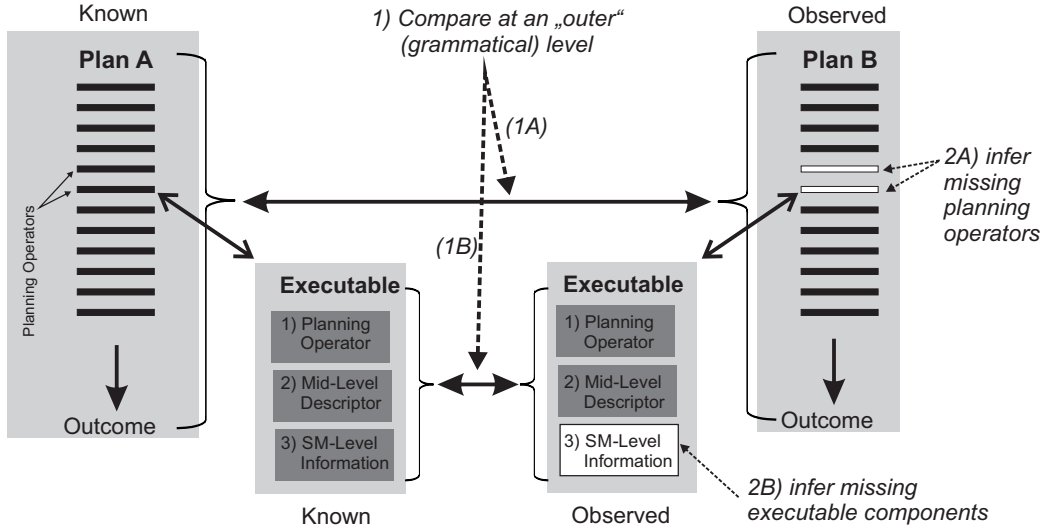


Figure 2: Schematic of structural bootstrapping.

### *Implementing Structural Bootstrapping at different levels*

Figure 2 shows a schematic representation of the bootstrapping processes implemented here. A known plan A (left) consists of a set of planning operators (black bars) and each has attached to it an Executable consisting of the planning operator itself, a mid level descriptor and sensorimotor level information. The plan, being executed, also has a certain “outcome”, which can be considered as the goal of this action sequence. An observed plan B (right) of a similar action (with similar goal), will normally consist of many planning operators which are identical or highly similar to the ones of the known plan and also the outcome will be similar or the same. Still some planning operators may be dissimilar and hence unknown to the agent (white bars). In the same way, individual newly observed Executables (right) may contain unknown components (white). The goal of bootstrapping is to fill in all this missing information. To the end, first (1) the respective entities, Plans (1A) or Executables (1B), will be compared at an “outer”, grammatical level to find matching components. This way, in the second step one can try to infer the respective missing entities, planning operators (2A) or components of the Executables (2B).

Hence, a *central statement* is that structural bootstrapping always “propagates downward”. It uses type-similarities of entities from one level above to

define the missing syntactical elements of the currently queried (lower) level. Plan similarities are used to infer planning operators, Executable similarities to infer Executable parameters such as objects, trajectories, forces, poses, and possibly more.

The main difficulty for implementing structural bootstrapping is to define appropriate scaffolds on which the bootstrapping can be based where – as described – the goal is to create novel information by generative processes which compare existing knowledge with newly observed one, without having to perform an in-depth analysis.

In the following we will now provide the details of the performed experiments, where we will show six different examples of structural bootstrapping for the different layers. These examples should allow the reader to more easily understand the so-far still rather abstract concept of structural bootstrapping.

## **Setup, procedures and specific problem formulation**

### *Scenario (task)*

ARMAR operates in a kitchen scenario. The task for the robot is to pour two ingredients (e.g. flour and water) and mix them together to obtain batter. For this the robot has the required knowledge to do it in one specific way (by using an electric mixer), but will fail whenever it should react flexibly to a changed situation (e.g. lack of the mixer). The goal of this work is to show that bootstrapping will quickly provide the required knowledge to successfully react to such a change. This process is based on observing a human providing an alternative solution (stirring with a spoon) where bootstrapping lead to the “understanding” of the meaning of objects and actions involved.

### *Prior knowledge*

As bootstrapping relies on existing knowledge we have provided the robot with several (pre-programmed) Executables and we assume that the robot knows how to:

- pick up an object;
- put down an object;
- pour an ingredient;

- mix with an electric mixer.

In addition, robot has learned earlier to execute one apparently unrelated action, namely:

- wipe a surface with a sponge [25, 26, 27].

Furthermore the robot has a certain type of object memory where it has stored a set of objects together with their roles, called the *Repository of objects with attributes and roles (ROAR)*. This prior knowledge can be inserted by hand or by prior experience. It allows objects to be retrieved by their attributes, and attributes of novel objects to be inferred, based on proximity in a low-dimensional, Euclidean space in which both, objects and attributes, reside [18].

The following entries exist in the ROAR:

- Sponge, rag, brush = objects-for-wiping with outcome: clean surface
- Mixer tool ends, whisks, sticks = objects for mixing with outcome: batter or dough.

Furthermore we have endowed the machine with a few recognition procedures:

- The robot can generate and analyze the semantic event chain (SEC) structures of observed (and own) actions by monitoring an action sequence using computer vision. Thus, the machine can recognize known actions at the SEC level [24, 17].
- The robot can recognize known objects (tools, ingredients, batter) using computer vision [28, 29, 30].
- The robot can explore unknown object haptically [31] and extract object features such as deformability and softness [32, 33, 25]

#### *Problem definition*

The problem(s) to be solved by structural bootstrapping are defined by several stages as spelt out next:

Normal System Operation: If all required entities are present (mixer, ingredients, containers, etc.) the robot can make a plan of how to make batter

and also execute it.

System Break-Down: Planning and execution will fail as soon as there is no mixer.

Alternative: The robot observes a human making batter by stirring the dough with a spoon.

Goal: The robot should find a way to understand the newly observed action and integrate it into its knowledge base and finally be able to also execute this.

Problem: The robot has no initial understanding of

- the planning requirements,
- the objects involved, and
- the movement patterns seen

in the newly observed stirring action. For example the robot does not know how to parameterize the rhythmic trajectory. Also, it does not know what a spoon is. Furthermore, the robot does not have any planning operator for stirring with a spoon in its plan-library.

Requirement (for the purpose of this study): The process of understanding the new action should happen without in-depth analysis of new actions constituents (hence without employing exploration based processes) but instead by using bootstrapping.

## **Methods - Short Summary**

To not extend this paper unduly, methods are only described to the details necessary to understand the remainder of this paper. References to specific papers are provided where more details can be found.

### *Planning Methods*

In this project, we are using the so-called Combinatory Categorical Grammars (CCGs) [16] to address the planning problem. CCGs are in the family



of *lexicalized* grammars. As such they push all domain specific information into complex categories and have domain independent combinators that allow for the combination of the categories into larger and larger categories. As we have already alluded to, at the planning level, structural bootstrapping is a specialized form of learning new syntactic categories for known actions. A number of different methods have been suggested for this in the language learning literature [34, 35] for this project however we will be applying a variant of the work by Thomford [36]. However, we note that to do the kind of learning that we will propose it will be critical that the same syntactic knowledge, which is used by the system to plan for action, is also used to recognize the plans of other agents when observing their actions. This is not a new idea, however, there are very few AI planning and plan recognition systems that are able to use the exact same knowledge structures for both tasks.

Imagine that, as in our example, the high level reasoner knows about a plan to achieve a particular goal. It knows all of the actions that need to be executed, and for each action has encoded as CCG categories the knowledge necessary to direct its search for the plan. Further we suppose the same knowledge can be used to parse streams of observations of actions in order to recognize the plan being executed by others.

Now suppose the agent sees the execution of another plan that achieves the same goal. Let us assume that this new plan differs from the known plan in exactly one action. That is, all of the actions in the new plan are exactly the same as the actions in the known plan except for one action. Since the agent knows that the plan achieved the same goal, and it knows the CCG categories for each action that would be used to recognize the original plan, it is not unreasonable for the agent to assume that the new action should be assigned the same CCG category as its opposite action in the known plan.

If this addition is made to the grammar the agent now knows a new plan to achieve the goal and will immediately know both how to recognize others executing the plan and how to build the new plan for the goal itself (at the higher “abstract” level). The system will have performed structural bootstrapping at the planning level.

In this case, the system will have leveraged knowledge about the outcome of the observed plan being the same as the previously known plan, along with syntactic knowledge about how the previously known plan was constructed to provide new syntactic knowledge about how to construct and recognize the new plan.

### Methods for the Mid-Level: Semantic Event Chains (SECs)

Semantic Event Chains [24, 17] encode in an abstract way the sequence of events that occur during a complex manipulation. They are used for two purposes: (1) Every event provides a specific temporal anchor point, which can be used to guide and temporally constrain the above described scene and motion analysis steps. And (2) the SEC-table itself (see Fig. 3 b), is used to define the mid-level of an Executable.

Fig. 3 shows the corresponding event chains extracted for a *stirring* action. SECs basically make use of image sequences (see Fig. 3 a, top) converted into uniquely trackable segments. The SEC framework first interprets the scene as undirected and unweighted graphs, nodes and edges of which represent image segments and their spatial touching or not-touching relations, respectively (see Fig. 3 a, bottom). Graphs hence become semantic representation of the relations of the segments (i.e. objects, including hand) presented in the scene in the space-time domain. The framework then discretizes the entire graph sequence by extracting only the main graphs, which are those

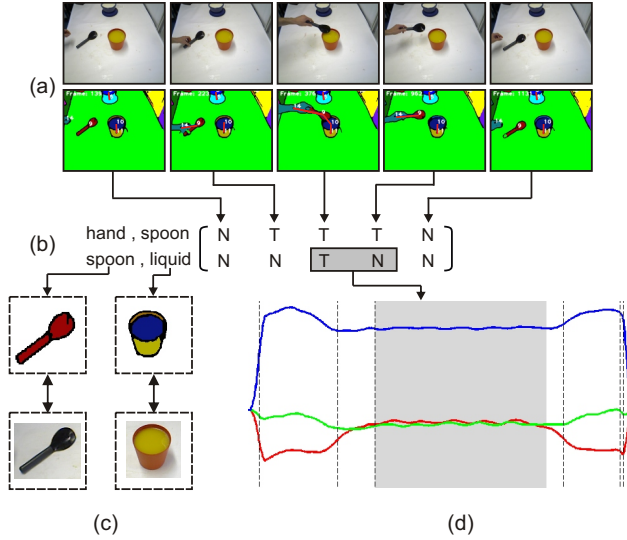


Figure 3: A real action scenario: “*Stirring liquid with a spoon*”. (a) Sample original key frames with respective segments and graphs. (b) Corresponding SEC where each key frame corresponds to one column. Possible spatial relations are N, T, and A standing for “*Not-touching*”, “*Touching*”, and “*Absence*”, respectively (*A* does not happen here.). Shaded box shows a sample relational transition. (c) Object identities derived from segments (d) Complete trajectory information for the hand. Trajectory segment for the time-chunk covered by shaded box in (b) is indicated in gray color.

where a relation has changed (e.g. from not-touching to touching). Each main graph, thus, represents an essential primitive of the manipulation. All extracted main graphs form the core skeleton of the SEC which is a sequence table (the SEC-table), where columns correspond to main graphs and rows to the spatial relations between each object pair in the scene (see Fig. 3 b). SECs consequently extract only the naked spatiotemporal relation-patterns and their sequentiality, which then provides us with the essence of an action, because SECs are invariant to the followed trajectory, manipulation speed, or relative object poses.

Columns of a SEC represent transitions between touching relations. Hence, they correspond to decisive temporal moments of the action and, consequently, they allow now to specifically pay attention “at the right moment when something happens” to additional action relevant information (such as objects, poses, and trajectories). Fig. 3 (c-d)) illustrate syntactic elements of the manipulation. Manipulated objects, e.g. spoon and liquid, are extracted from the rows of event chains, i.e. from the nodes of the main graphs. Temporal anchor points provided by SECs can also be used to segment the measured hand-trajectory into parts for further analysis.

### *Sensorimotor Methods*

Sensory Aspects: Visual scenes are analysed to recognize objects and their attributes, measure movement trajectories, and record object poses.

Basic object and pose recognition is performed in a straight-forward way using pre-defined classes of the different objects which occur during the actions of “stir”, “wipe”, and “mix” and in addition adding some distractor objects (e.g., cups, knives, etc.). Any suitable method can be used for object detection, recognition, and pose estimation; such as edge-based, statistical shape representations [28, 29, 30, 37].

Another important aspect is object recognition for the construction of the repository of objects with attributes and roles (ROAR).

Our primary input for the ROAR consists of a table such as the one shown in Table 1.

Objects and attributes are (discrete) labels; values can be categorical, discrete or continuous. Examples of objects are “bowl” or “knife”; examples of attributes are “cuts”, “food”, “is elongated”, “gripper orientation for grasping”, “fillable”, etc. We then use Homogeneity Analysis to project ob-

	Attribute 1	Attribute 1	Attribute 1
Object A	$Value_{A,1}$	$Value_{A,2}$	$Value_{A,3}$
Object B	$Value_{B,1}$	$Value_{B,2}$	$Value_{B,1}$

Table 1: ROAR encoding

jects and (attribute) values into the same, low-dimensional, Euclidean space (the *ROAR space*) [18]. This projection is set up such that:

- Objects that exhibit similar attribute Values are located close together,
- Objects that exhibit dissimilar attribute Values are located far apart,
- Objects-as-such are close to their attribute Values.

Euclidean neighborhood relations allow us to make the following general types of inference:

- Attribute value prediction: Say, we have an object of which we know some but not all attribute Values. We can predict missing attribute Values by projecting the object into the ROAR and examining nearby attribute Values.
- Object selection: Say, we have a set of required attribute values. We can find suitable objects in the vicinity of these Values in the ROAR.

Note we cannot generally expect that very complex object/attribute relations will be faithfully represented in a low-dimensional Euclidean space. While we are currently working on more powerful representations for such relations, this is a complex research issue [18, 38, 39, 40, 41]. For us the ROAR is at the moment just a viable way forward, which allows us to demonstrate different aspects of structural bootstrapping.

Motor Aspects: Trajectory information is encoded by Dynamic Movement Primitives (DMPs), which were proposed as an efficient way to model goal-directed robot movements [19]. They can be applied to specify both point-to-point (discrete) and rhythmic (periodic) movements. A DMP consists of two parts: a linear second order attractor system that ensures convergence to a unique attractor point and a nonlinear forcing term. The forcing term

Original Plan	Observed Plan
<pre> testName: xpermix; initialState: [ ]; observations: [   pickA( left, beaker, t ),   pourA( left, liquid1, beaker, mixingBowl ),   placeA( left, beaker, t ),   pickA( left, cup2, t ),   pourA( left, liquid2, cup2, mixingBowl ),   placeA( left, cup2, t ),   pickA( right, mixer1, t ),   mixA( mixer1, liquid1, liquid2, mixingBowl ) ]; </pre>	<pre> testName: xpermixnew; initialState: [ ]; observations: [   pickA( left, beaker, t ),   pourA( left, liquid1, beaker, mixingBowl ),   placeA( left, beaker, t ),   pickA( left, cup2, t ),   pourA( left, liquid2, cup2, mixingBowl ),   placeA( left, cup2, t ),   pickA( right, UNKNOBJ, t ),   UNKNACT( UNKNOBJ, liquid1, liquid2, mixingBowl ) ]; </pre>

Figure 4: Comparing known with observed plan. The arrow indicates where there is a novel, unknown planning operator found in the new plan. This is also associated with an, as yet, unknown object (the spoon).

is normally given as a linear combination of basis functions that are defined along the phase of the movement. The basis functions are either periodic or nonzero only on a finite phase interval. The type of basis functions decides whether the DMP defines a discrete or a periodic movement. DMPs have many favorable properties, e.g. they contain open parameters that can be used for learning without affecting the overall stability of the system, they can control timing without requiring an explicit time representation, they are robust against perturbations and they can be modulated to adapt to external sensory feedback [19, 42].

## Concrete Examples of Structural Bootstrapping

### *Structural Bootstrapping at the Planning Level*

The existing plan of making batter with a mixer is compared to the observed sequence of actions during making batter with a spoon. Due to the fact that all sub-actions, but one, are identical between known-action and new-action the agent can infer that the unknown sub-action (stirring with a spoon) is of the same type as its equivalent known sub-action (mixing with a mixer). Hence the grammatical comparison of known with unknown action renders a new (syntactic) planning operator entry for the unknown

sub-action. This process is very similar to syntactic bootstrapping as observed in child language acquisition. A semantic element enters here due to the same outcome of both actions being recognized as batter. We use CCG as our planning language and we employ the PKS planner [15] for the actual planning processes of ARMAR III.

The actual inference process makes use of the similarity of known plan with newly observed plan, where in our example all but one action are identical.

Figure 4 shows the comparison between a known (and executable) plan on the left and an observed new one (right). Structural (grammatical) one-by-one comparison shows that there is just one unknown planning operator present. When the plan recognizer is run on the observed plan it would result in the following explanation of those observations with the highest probability:

```
[ addIngC(left, liquid1, beaker, mixingbowl),
  addIngC(left, liquid2, cup, mixingbowl),
  pickC(left, UNKNOBJ, table),
  UNKNACT(left, UNKNOBJ, liquid1, liquid2, mixingbowl)]
```

Note, the category name for the previously unseen action is simply denoted as UNKNACT. This is a special purpose category used to complete the explanation when we have an action that has never been seen before.

Now the agent has been told (or can observe) that the observed plan is a plan that achieves makeBatterC (making batter), and we will assume that all of the actions in the observed plan are relevant to the plan. The agent's job is to infer a category to replace UNKNACT that allows the completing of the parse. If the agent wants to build a category to assign to the unknown action that will result in a complete plan with the goal of makeBatterC, all it needs to do is walk the explanation from right to left collecting the categories and adding them to the complex category in order. This will result in the unknown action being given the following category:

```
action: UNKNACT(hand, UNKNOBJ, ingredient, ingredient, bowl)
  [ ((( makeBatterC( 2, 3, 4 ))\
    {addIngC( 0, 2, obj(1), 4)}\
    {addIngC( 0, 3, obj(2), 4)}\
    {pickC( 0, 1, table(1)) } ] ];
```

Note the agent also infers the types and co-reference constraints for the basic category’s arguments from the plan instance. In the above definitions we have denoted those arguments to the basic categories by numbers indicating when an argument is bound to the same argument as the action. (i.e. All references to “0” in the category refer to the hand used in the action because it is the zeroeth argument for the action. Likewise all reference to “4” in the category refer to the bowl argument of the action since it is the fourth argument.)

This category would represent the most restrictive hypothesis about the plan structure since it will require both that the actions be executed in the same order (and we know the ingredients can be added to the plan in either order) and that all of the arguments that co-refer in the example plan must co-refer in future instances. In this case, it would require that the same hand be used for all of the ingredient adding and mixing which we know to be overly restrictive.

If we compare the new category to the category for the known mix action (mixA), we can see that the only differences are exactly in these overly restrictive areas:

1. The ordering of the categories for the add ingredient steps. The known category is more general allowing the ingredients to be added in any order while the new learned category has a required order.
2. The co-reference constraints are less restrictive in the known category. (Note the numbers indicating, which hand is to be used in the addIngC, are not the same so the plan would not enforce that the same hand be used.)

At this point, on the basis of the structural information provided by the parse and the action grammar, the agent has inferred that “UNKNACT” is equal to (or at least very similar to) “mixA” and the information can be entered directly into the planning grammar of the agent and forms the top-level of the corresponding new executable. We will, for convenience, from now on name it: “stir”, hence we set:

`UNKNACT:=stir.`

While we have now added a new action to the planning grammar, still there is massive information lacking for designing the complete (new) executable for “stir”, for example there is as yet no understanding existing about the UNKNOBJ (the spoon) and nothing is known about several other mid- and low-level descriptors.

<b>A) Picking up</b>									
Hand, Beaker	0	1	1						
Beaker, Table	1	1	0						

<b>B) Putting down</b>									
Hand, Beaker	1	1	0						
Beaker, Table	0	1	1						

<b>C) Pouring</b>									
Hand, Beaker	1	1	1	1	1				
Beaker, MixBowl	0	1	1	1	0				
Beaker, Liquid2	1	1	1	0	0				
MixBowl, Liquid2	0	0	1	1	1				

<b>E) Mix (with Mixer)</b>									
Hand, Mixer	0	1	1	1	0				
Mixer, Dough	0	0	1	0	0				

<b>D) Wipe (with Sponge)</b>									
Hand, Sponge	0	1	1	1	0				
Sponge, Surface	0	0	1	0	0				

<b>F) Stir (was UNKNACT) with Object*</b>									
<b>Unknown SEC</b>									
Hand, Object	x	x	x	x	x	x			
Object, Dough	x	x	x	x	x	x			

<b>G1) Stir (was UNKNACT) with Object*</b>									
<b>SEC from one observation</b>									
Hand, Object	0	1	1	1	1	0			
Object, Dough	0	0	1	0	1	0			

<b>G2) Stir (was UNKNACT) with Object*</b>									
<b>SEC from two observations</b>									
Hand, Object	0	1	1	1	0				
Object, Dough	0	0	1	0	0				

\*Object = "UNKNOBJ", before object specification  
Object = "spoon" after object specification

Figure 5: Several important SECs, which occur during the different actions. Headlines (bold lettering, like “Picking up”, etc.) denote the type-specifiers of the different SECs. Note, sometimes objects can change. E.g. “Beaker” can be replaced by “Cup2”. **A-E)** error-free archetypical SECs from known actions. **F)** So-far unspecified SEC. **G)** SECs from the unknown action extracted from observation of the human performing it. Hence these SECs might contain errors. **G1)** one observed case, **G2)** two observed cases. (In human terms: G1 corresponds to a case where the spoon had intermittently been pulled out from the dough (grey box), whereas for G2 it always remained inside until the action terminated.)

### Structural Bootstrapping at the Mid-Level

At the mid-level, we need to define the correct SEC for “stir”. Figure 5 shows SECs for several actions where (F) represents the so-far unknown SEC for “stir”. Please, ignore panels (G) for a moment. Note, to be able to treat these tables numerically the intuitive notations from Figure 3 for non-touching “N” and touching “T” are now changed to “0” and “1” in Figure 5.

Structural bootstrapping at the mid-level uses as the “outer”, grammatical scaffold the type-similarity of the planning operators (here “stir” and “mix”) ascertained above. Hence we know that UNKNACT=stir.

Following this conjecture the agent can now with a certain probability



assume that so-far unknown SEC for “stir” ought to be identical (or very similar) to the known one from “mix” and use the “mix”-SEC to define the mid-level (the SEC) for the Executable of “stir”. The arrow indicates that the SEC from panel (E) should just be transferred to fill the unknown SEC in (F) with the same entries.

There is a second line of evidence which supports this. Panels (G1) and (G2) represent the actually observed SECs of the stirring action here from a total of three observations of a human performing this. The SEC in panel (G1) had been observed once and the other twice. By comparing these SECs, the robot can with some certainty infer that the transfer of (E) to (F) was correct, because the more often observed SEC in (G2) corresponds to it, while the SEC from panel (G1) might be noisy as it is a bit different. As shown in an earlier study [24, 17], more frequent observations are likely to confirm this even more, but were not performed with the current setup.

### *Structural Bootstrapping at the Sensorimotor Level*

Bootstrapping at the this level is used by the agent to find out how stirring is actually done (motion patterns), what the meaning of “UNKNOBJ” is, and which other objects might have a similar meaning. Before going into details we can quickly state that at the sensorimotor level several bootstrapping processes can be triggered. We note that bootstrapping is a probabilistic process and things can go wrong, too. One such example is, hence, also included. We find that the following processes are possible:

#### **1. Motion**

- (a) Bootstrapping from SEC-similarities [“wipe” and “stir”] to define the motion patterns for “stir”.

#### **2. Objects**

- (a) Bootstrapping from SEC-similarities [“wipe” and “stir”] into the new action. Here arriving at a false conjecture that “sponges” could be used for mixing.
- (b) Bootstrapping from SEC-similarities [“mix” and “stir”] from the repository of objects with attributes and roles (ROAR) into the new action seeking different objects that could potentially be used for mixing.
- (c) Bootstrapping from SEC-similarities [“mix” and “stir”] from the new action into the ROAR, entering the “spoon” into the category of objects for mixing.

To address the sensorimotor level the agent has to bootstrap from the mid-level downwards. It can do this by comparing the *type-similarities* of the different SECs. For this essentially one calculates a sub-string comparison of the rows and columns between one SEC and any other [24, 17]. We obtain that “stir” and “mix” as well as “stir” and “wipe” are 100% type-similar (compare panels D, E, and G2 in Figure 5), whereas “stir” and “pour” are only 52% similar, etc. Thus, the agent can infer that syntactical elements from “mix” and “wipe” might be used to define missing entities at the sensorimotor level of the Executable.

*1a) Motion: Bootstrapping from SEC-similarities “wipe” and “stir” into the new action for completing motor information*

Here we make use of the fact that the SEC for stir is very similar to the known one from wipe. Figure 6 shows the SECs and the different trajectories recorded from human observation for both actions. Note that for “wipe” the complete motor encoding is known and provided by the respective DMP parameters.

We have in our data-base the following description for “wipe”: Since wiping is essentially a rhythmic movement, we use periodic dynamic movement primitives to specify the required behavior [27]. Periodic DMPs are defined by the following equation system [19]

$$\dot{z} = \Omega \alpha_z (\beta_z (g - y) - z) + f(\phi), \quad (1)$$

$$\dot{y} = \Omega z, \quad (2)$$

In the above equations,  $g$  is the anchor point of the periodic movement. The nonlinear forcing term  $f$  is defined as

$$f(\phi, r) = \frac{\sum_{i=1}^N w_i \Psi_i(\phi)}{\sum_{i=1}^N \Psi_i(\phi)} r, \quad (3)$$

$$\Psi_i(\phi) = \exp(h_i \cos(\phi - c_i) - 1),$$

where the phase  $\phi$  is given by

$$\dot{\phi} = \Omega. \quad (4)$$

Here we assume that a complete parameterization of the DMP for wiping has been learnt from earlier experiences. Given this the DMP can be easily modulated by changing:

- the anchor point  $g$ , which translates the movement,

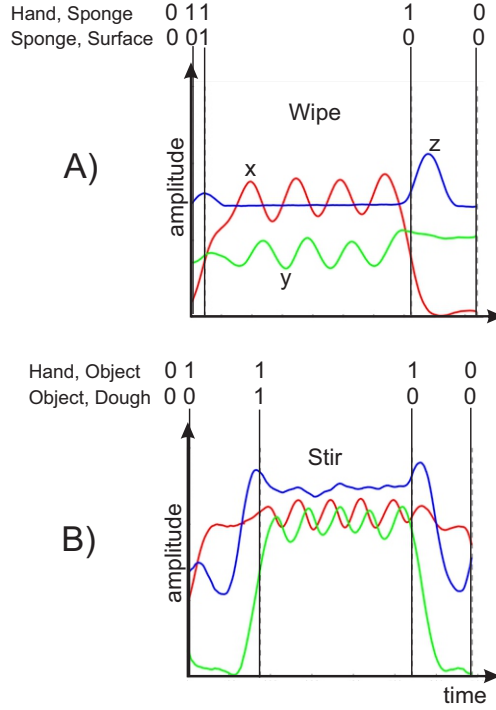


Figure 6: Bootstrapping motor information. SECs (top) and trajectories (bottom) for x, y, and z coordinates in task space are shown for **(A)** wipe and **(B)** stir.

- the amplitude of oscillation  $r$ ,
- the frequency of oscillation  $\Omega$ .

These variables can be used to immediately adapt the movement to sensory feedback.

Bootstrapping progresses by using the concept of temporal anchor points, which are those moments in time when a touching relation changes (from 0 to 1, or vice versa). These anchor points divide the trajectories in a natural way (shown by the vertical lines in the figure.)

Bootstrapping now *just copies* the complete DMP information from “wipe” to the Executable of “stir” between the respective anchor points only leaving the constraint-parameters (e.g. amplitude) open as those are given by the situation (mainly the size of the bowl wherein to stir). Thus, the agent assumes that it can use the motor encoding from “wipe” in an unaltered way

to also perform “stir”. We know from own experience that this largely holds true. Here we can also clearly see the advantages of bootstrapping: we do not need any process that *extracts and generalizes* motor information from the observed example(s) of “stir” (a process which could be more tediously performed by methods from imitation learning [43, 44, 45]). Instead we just copy. Clearly, the agent - like any young child - will have to ground this by trying out the stirring action (see the Discussion section for the “grounding-issues”). It will possibly then have to adjust the force profile, which is likely to be much different for wipe and stir. Still, all this is faster than learning the required motor pattern in any other way. The benchmark experiments below show this clearly.

*2a) Objects: Bootstrapping from SEC-similarities “wipe” and “stir” into the new action for object use*

The SEC-similarities between “wipe” and “stir” allow the agent to also (wrongly!) infer that the object for wiping (sponge) should be suitable for stirring, too. Note this may seem unexpected but can happen during any bootstrapping process due to its probabilistic nature. The use of just one single scaffold (here the SECs) is not strong enough to allow rigorously excluding such false conjectures. For this the agent needs to integrate additional information and, due to the fact that there is a repository of objects with attributes and roles (ROAR), it can indeed obtain evidence that there has been an error.

The agent knows that “stir” and “mix” are at the mid-level (SEC) type-similar action. It finds, however, that sponges are clearly outside the cluster of objects for mixing (Figure 7 A). This lowers the probability substantially that sponges should be used for mixing/stirring actions.

Interestingly, children will many times indeed over-generalize and use “unsuitable” objects for an intended action [46]. It is unknown how the brain represents this, but – clearly – their representation does apparently not yet contain the fine grained-ness of an adult representation.

*2b) Bootstrapping from SEC-similarities “mix” and “stir” from the ROAR to find other suitable objects*

Here the agent falls back (again) on the similarity of the new SECs of “stir” with the known one of “mix”. Due to this similarity, the agent knows that appropriate objects for the novel action might be found in the cluster of “objects for mixing” in the repository of objects with attributes and roles.

## Repository of Objects with Attributes and Roles (ROAR)

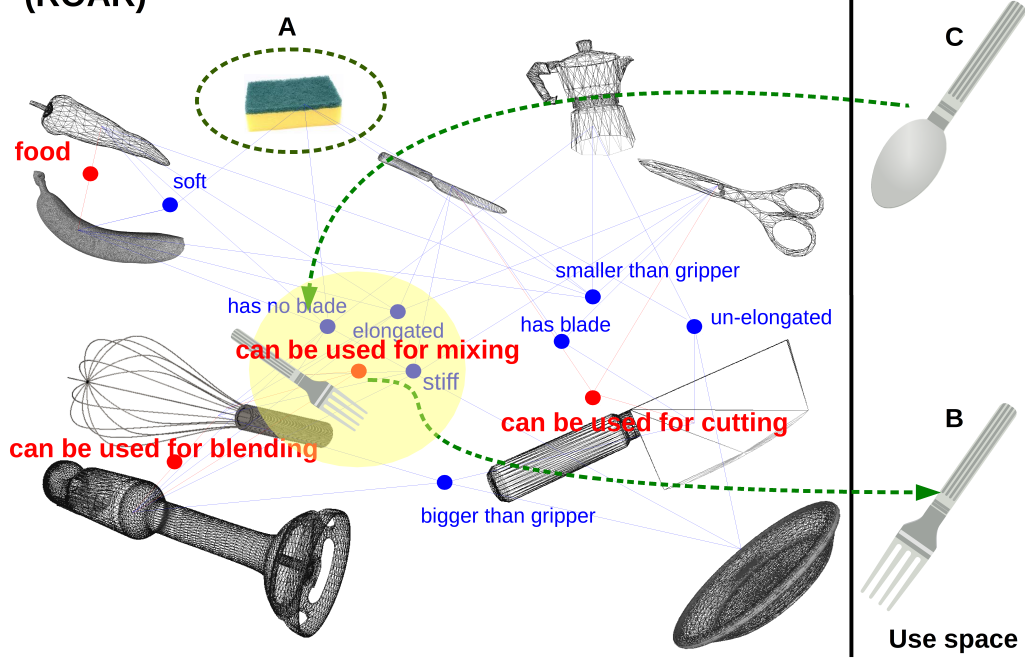


Figure 7: Bootstrapping object information. Graphical rendering of the repository of objects with attributes and roles (ROAR). Depicted are the metric distances between the different objects and the attribute values that describe their respective roles. **A)** The sponge is located far from the attribute value “can be used for mixing”. **B)** Bootstrapping allows inferring that a fork, found close to the “mixing” attribute value, could be used also for “stir”, as “mix” and “stir” are at the SEC-level type-similar. **C)** Following this SEC-similarity, a novel object (spoon) with unknown “mixing” attribute may be hypothesized useful for mixing by the ROAR and also due to other, known attribute values (such as shape, stiffness, and SEC characteristics of known, observed actions).

Hence it can ask the repository for a tool suitable for mixing and maybe locate it somewhere else on the table. Clearly this process will lead to an action relevant result only in those cases where the agent actually find such an object within reach. Then it can try to use this object for stirring, too. Again we can draw similarities to our own behavior. Generally this type of tool-replacement is found for a wide variety of actions where we “define” the tool according to its planned use. Our own generalization properties may here

go far beyond what the ROAR offers to our artificial robotic agent, which is evident from situations where we “abuse” objects for entirely different purposes.

*2c) Bootstrapping from SEC-similarities “mix” and “stir” from the new action into the ROAR to create a new entry*

In the last step, the agent can perform one more bootstrapping procedure to augment the repository of objects with attributes and roles. For this it analyzes the outcomes of the actions realizing that batter is obtained from “mixing” and also from the unknown action of “stirring”.

Thus, the agent can enter the new observed tool (spoon) into the ROAR and can then – by virtue of its resulting position in the ROAR – infer other, unobserved attribute values (uses), which is a bootstrapping effect. This way the repository will be extended by a novel entry following a single-shot experience. This step, however, does require a parametrization of the new object according to the features used for the ROAR.

## **Robotic implementation and benchmark experiments**

Note, the actual bootstrapping processes happen “inside the machine” and any demonstration will, thus, only show that “the robot can do it now”. To go beyond such mere visual inspection, one needs to show quantitative results on performance gain by bootstrapping, which will be presented in the next sections, below.

Still, a complete robotic implementation of these processes is currently being performed using the our robot systems [47]. For brevity, we will here show one central part of this implementation demonstrating the required transfer of human action knowledge (Fig. 8 A) onto the robot. This is the initial step needed to set up action knowledge in the machine before any bootstrapping can happen. The robot acquires here the knowledge to perform mixing with a mixer.

To better be able to extract object relations we have here used a Vicon-based motion capture system from which we immediately get error-free Semantic Event Chains (Fig. 8 B). The complete action relevant information is extracted at the respective key frames and encoded into the required Executables (Fig. 8 C), which can be used by the robot to reproduce this action (Fig. 8 D). The complete experiment is described elsewhere [48].

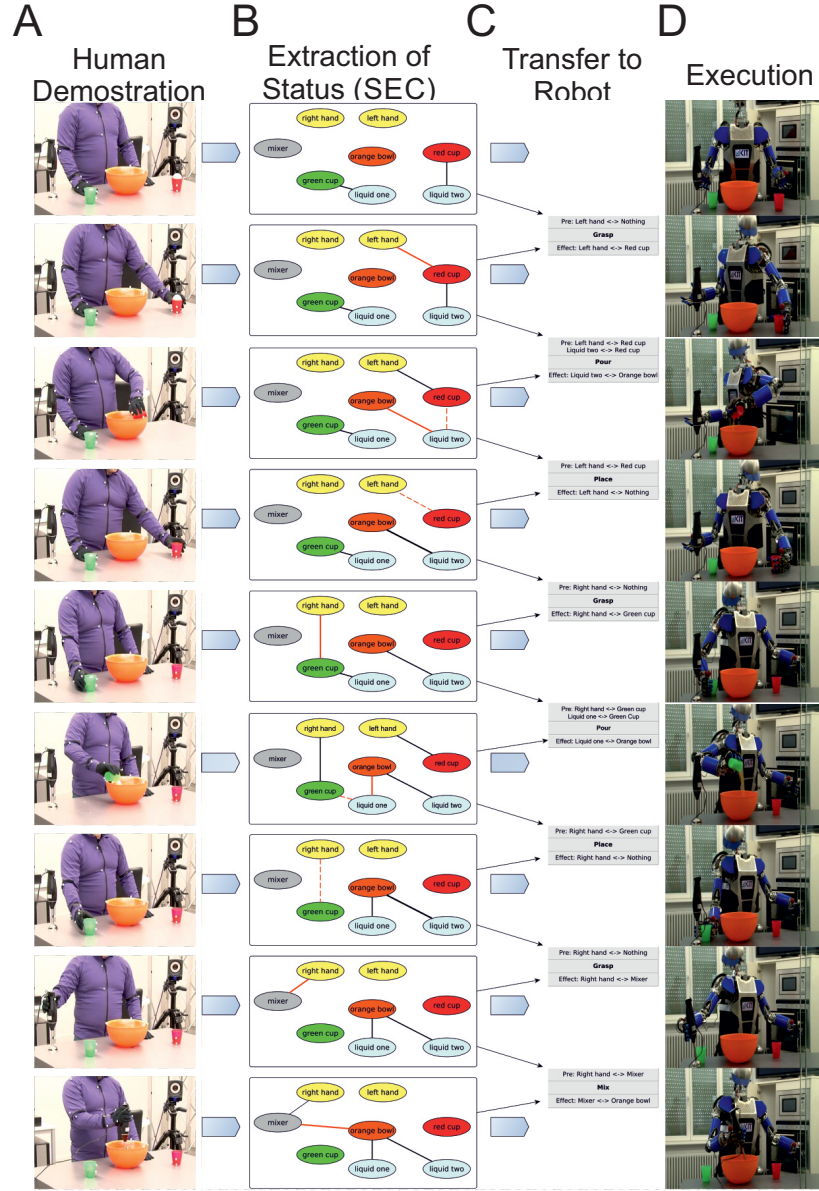


Figure 8: Transfer of action knowledge from human to robot. **A)** Human demonstration, **B)** SEC depicted by ways of its key frame graphs, which show which objects touch which other objects (edges) during human execution. **C)** Abbreviated Executables **D)** Robot execution.

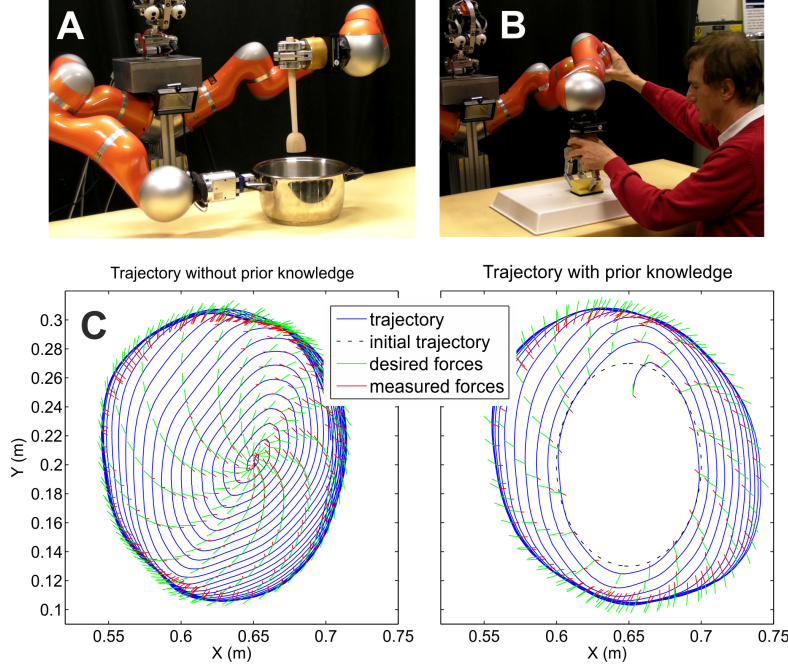


Figure 9: Benchmark experiment demonstrating the gain of learning speed when bootstrapping motion trajectories. **A)** Experimental setup and **B)** demonstration of wiping. **C)** Learning of stirring behavior without prior knowledge and **D)** adaptation of wiping to stirring. The desired and actual forces are shown with red and green vectors.

In the following we will show some experiments from our scenario demonstrating the power of structural bootstrapping for example the speed-up as compared to conventional, exploration based learning methods but also the accuracy of the object attribution methods used in the bootstrapping process.

#### *Bootstrapping Motion - Measuring Learning Speed*

Our setup for learning of stirring behavior is shown in Fig. 9 A. It is composed of two KUKA LWR robots, both equipped with Barred hands. The task is to learn how to stir in a metal pad of diameter of 21 cm using wooden spoon. The position, size and shape of the pad are not known in advance. To define a criterion function for motion learning, we specify the force  $\mathbf{F}_d$  with which the robot should move along the edge of the pot.

We considered two cases: 1) learning without any prior knowledge about



the stirring trajectory and 2) learning where the adaptation process is initialized with wiping trajectory. The wiping trajectory is obtained by imitation learning (Fig. 9 B). We used periodic DMPs to represent the movement [49] and apply a Repetitive Control (RC) algorithm [50, 51]. The RC algorithm iteratively adapts the current estimate of the stirring behavior to achieve the desired outcome as defined by the desired contact force. Task performance is improved with each repetition and eventually, the required behavior is achieved regardless of the initial estimate of the stirring trajectory.

Fig. 9 C,D show the progress of learning in  $x$ - $y$  plane for both cases. The robot learned the policy in approximately 15 cycles without any prior knowledge about the trajectory and in approximately 7 cycles with prior knowledge taken from wiping motion. This demonstrates in a practical example that low-level sensorimotor learning can significantly benefit from the initialization provided by the semantic understanding of the task.

Note that in the specified scenario, the direction of adaptation is provided by the information about the desired contact force. We can expect that the difference between the two approaches would be even bigger if model-free methods such as reinforcement learning were used.

#### *Bootstrapping Objects - Measuring Success*

Trivially immediate object replacement, using the ROAR as suggested for cases 2a, b, and c above, will always be faster than finding an appropriate object by exploration, but will the ROAR find the correct object?

We evaluate the capacity of the ROAR to predict the suitability of given objects for mixing, similar to the scenario above. To this end, we created a database of 10 objects as listed in Fig. 10. Each object is characterized by 10 binary attributes describing its properties (such as shape and stiffness) and usage categories (such as “container” or “mixing tool”), some of which may be unknown. The ROAR ranks objects according to their estimated suitability for mixing. Fig. 10 shows the suitability of the 10 objects as a function of the proportion of missing attribute values. Each column of the graph represents results averaged over 100 runs.

For each run and for each proportion of missing attributes, the designated proportion of object attribute values is randomly chosen from the complete database; these are set to **unknown**. On the resulting copy of the database with missing attribute values, homogeneity analysis is performed (see “Methods”), producing a ROAR. Objects are ranked by the ratio of their Euclidean

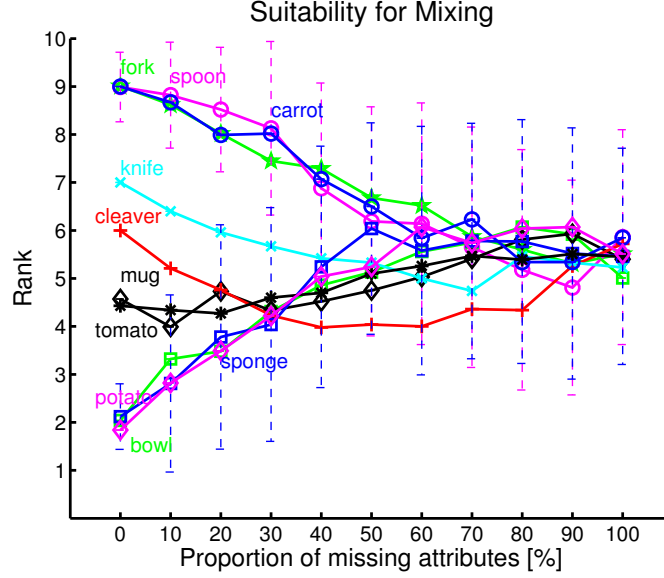


Figure 10: Estimated suitability of 10 different objects for mixing, ranked by the ROAR. With fully known attributes, the ROAR consistently considers the spoon, the fork and the carrot as useful mixing tools, and the sponge, potato and bowl as useless. This consistency degrades gracefully with increasing percentage of missing attributes. Each column of the graph represents ranks averaged over 100 runs. Error bars give standard deviations for our two objects of interest (sponge and spoon); those of the other objects are similar but not shown to reduce clutter.

distances to the “mixing tool”=`true` vs. “mixing tool”=`false` attribute values [18].

With fully known attributes, the ROAR consistently ranks the spoon, the fork and the carrot as the most useful mixing tools, while the sponge, potato and bowl rank last. This consistency degrades gracefully with increasing percentage of missing attributes.

## Discussion

A central issue for the development of autonomous robots is how to quickly acquire new concepts to be used for planning and acting, for example learning to stir in a pot, which is a relatively complex manipulation sequence. Reinforcement learning or association-based learning methods have been applied to this and related problems but are usually too slow to achieve this

efficiently. Thus, one often combines them with supervised learning. Here, especially the paradigm of learning from demonstration has often been successfully employed [52, 43, 53, 54, 55, 45] also because we (humans) are rather good at this. Still none of these methods is generative in the sense that it would take existing knowledge to generalize it into novel unexplored domains. At best one finds in-domain generalization, such a generalizing across different trajectories of the same action-type [20, 56, 57, 58].

This may not make us wonder, though. After all, generative learning is clearly an advanced cognitive trait and the gap between human performance and the current capabilities of machines is exceedingly wide. The central problem seems to be that – on the one hand – one has clear evidence that such processes do indeed happen during human (infant) learning [1, 8, 10, 11, 59], but – on the other hand – no one knows how; let alone, no one seems to have convincing ideas of how to do this with artificial agents either.

This was also the main challenge which we faced in this study: How can one develop a set of generative processes that use an “outer”, grammatical representation to bootstrap missing “inner”, syntactic elements, preferably at different levels of a cognitive architecture (planning, mid-level, and sensorimotor level). Furthermore our goal was to define such processes in a rigorous, algorithmically implementable way, to actually allow a robot to do this.

Language development did offer us a useful analogue on which we could build in this study. Semantic and syntactic bootstrapping [1, 2, 3, 4, 5, 6, 7, 8, 9], by which a child infers the meaning of unknown words using prior knowledge both rely on a general principle which we also used here: Grammar provides a solid scaffold for the probabilistic reasoning required for such inferences. While this was a helpful notion, still it remained unclear what the grammatical elements of an action sequence are (see [60] for a set of articles related to action-verb learning in children).

### *Bootstrapping at the planning level*

Planning languages and planning operators can be rather directly linked to the “language of action”. Since the earliest days of AI research on symbolic planning, the ideas of abstraction and hierarchy and the decomposition of high level plans into lower level plans has been seen as central to efficiently building plans [61, 62]. Many current researchers view knowledge of such plan hierarchies as “domain specific control knowledge”, that is knowledge of how to construct plans that is specific to individual domains. This kind

of knowledge has traditionally been encoded in Hierarchical Task Networks (HTNs) [63]. A formal relationship has been shown between HTNs and other similar plan structures and Context Free Grammars (CFGs) that are used extensively in natural language processing, formal grammar theory and theory of computation [63]. Here essentially we were representing our search control knowledge as a grammar and thereby it becomes quite clear how to extend the idea of syntactic&semantic bootstrapping to the symbolic planning domain. In this case, our objective was to learn the “syntactic knowledge” that encodes how to effectively build a new from an old plan.

Thus, for us it was relatively straight forward to implement structural bootstrapping at the planning level. The similarities of two plans allows inferring missing planning operator information (Fig. 4). But this addresses only the highest, the symbolic, level of an action sequence. It is for robotics totally useless to utter commands like “pour liquid”, without also providing the required, complex sub-symbolic information of how to actually do this.

#### *The problem of mid-level scaffolds*

Hence, more was needed to bridge the gap from symbols all the way down to the control signals of the robot motors. In some earlier studies we had introduced the Semantic Event Chain (SEC) as a possible mid-level descriptor for manipulation actions [24, 17, 64]. The SEC framework analyzes the sequence of changes of the relations between the objects that are being manipulated by a human or a robot. Consequently, SECs are invariant to the particular objects used, the precise object poses observed, the actual trajectories followed, or the resulting interaction forces between objects. All these aspects are allowed to change and still the same SEC is observed which, thus, captures the essence of the action as demonstrated in several action classification tests performed by us [17, 64]. In fact, SECs can be used to form an ontology of manipulation actions, where one finds that there are about 30 manipulation action types existing, which can be captured by the SEC framework [65].

It turned out that SECs offer two important aspects which make them good scaffolds for the bootstrapping of lower-level sensorimotor information.

1. SECs provide temporal anchor points, annotating in an action when “something decisive” has happened. This allows the chunking of an action and thereby provides the agent with a means to perform motor-pattern replacement (here wipe for stir), because “it knows” when to do the replacement.

2. Above we stated that SECs are invariant to the particular objects used. This is also essential for the bootstrapping. Only through this, object replacement is immediately permitted as the scaffold (the SEC) is not bound to particular objects as long as the chosen-one performs the same role (performs the same NT, TN transitions).

Ideas to utilize (spatial) relations to approach the semantics of actions first appeared in 1975. Badler [66] used directed scene graphs where each node identifies one object. Edges represent spatial information (e.g., LEFT-OF, IN-FRONT-OF, etc.) between the objects. Based on the object’s motion patterns, events are defined. Taken together this then represents an action. This, approach came to a stand-still, though because only now powerful enough image processing methods are available to provide the required information.

Even by now there are still only a few approaches towards semantic action understanding [67, 68, 69], often based on very complex activity graphs [67]. In [68], segmented hand poses and velocities are used to classify manipulations based on a histogram representation and using support vector machine classifiers for categorization of the manipulated objects. Others [69] introduced a visual semantic graph to recognize the action consequences based on changes in the topological structure of the manipulated objects.

In the context of the current study, potentially all these different approaches could be used as mid-level scaffolds, because they are based on the fact that the human action space is rather limited [65] and we are in fact not restricted by the here used SECs.

#### *Bootstrapping low-level information*

Any of these mid-level scaffolds could thus be used to guide bootstrapping at the control level, where we had shown 4 different examples (bootstrapping headlines 1a, 2a-c, see above). Here mainly visual information is used. This is done by linking shape similarities to action affordances into categories. These categories create the links in the repository of objects with attributes and roles.

The learning of perception and action categories requires quite some time during human development because large scale statistics on perceptual data need to be acquired and evaluated to sufficiently ground the categories. This learning process is working along two tracks. On a behavioral track, a rather small set of archetypical behaviors (as outlined in [46]) ensures the early association of objects with actions. The general execution of an action generates

the required low-level sensorimotor experience later to be used for structural bootstrapping and facilitates a model building by creating internal world knowledge. This – in turn – can be used by older children and adults to perform mental simulations of potential action contingencies thereby creating the second track.

The fundamental problem of these processes is the dimensionality of the potential sensorimotor contingencies (e.g. think of the visual input space, [70]) leading to a level of complexity that generates a very difficult learning/simulation task. To handle this complexity, an appropriate representation of sensorimotor information is required. Analysis of the visual representation in the vertebrate brain suggest that this takes place in form of a deep hierarchy which potentially allows for providing search spaces with different degree of granularity, different order of feature combinations and different levels of semantic abstraction [71]. This may lead to the required complexity reduction and could lead to the emergence of new structures in the internal world model of the agent further speeding up structural bootstrapping.

#### *Acquiring basic experience and the grounding issue*

Finally we would like to return to the key claims of this study: Does structural bootstrapping really represent a concept that (a) will lead to much *faster* knowledge acquisition and (b) is a way for the generative acquisition and extension of knowledge by which an agent can *more efficient* redeploy what it currently knows?

Central to both claims is that there exists already a solid and rich-enough knowledge base on which structural bootstrapping can operate. This knowledge should exist for all layers of the cognitive architecture and here is currently still a big bottleneck in robotics. The currently existing robot systems can only very slowly and tediously acquire sensorimotor experience either by exploration-based learning or from learning-by-demonstration. Thus, there are no rich sensorimotor knowledge bases existing anywhere in the robotics world. This may partly be due to the fact that this layer is usually very much embodiment-dependent and it is hard to define a generic sensorimotor knowledge base. Attempts to achieve this are currently being made in the European ACAT project [72]. Things are a bit better for symbolic knowledge and robotic-relevant knowledge-bases begin to emerge thanks to several large scale efforts like ROBOEARTH [73] or ROBOHOW [74] that use – for example – internet information to shape their knowledge bases and strongly focus on declarative (symbolic) knowledge.

Thus, the here shown experiments can currently only in a very point-wise manner support those two claims. We have indeed now implemented a more complex scenario (“making a salad”), where structural bootstrapping happens on-line (during the execution of the task) providing additional support to the here presented concepts [47]. Both claims, however, will probably get more and more substantiated the richer the individual knowledge bases at the different layers will become in the future. Hence, exploration-based and other, similar bottom-up learning methods will continue to play an important role for achieving this.

In a similar way, exploration will remain important also for the grounding of knowledge inferred by bootstrapping. One example above showed that such an inference process can also go wrong. Hence, very much like humans (especially children), robots also need to try out whether or not the newly acquired entity will actually work and also *how* it works. Inferring knowledge about “cutting”, “skiing”, “playing tennis” will never tell you the actual skills (force profiles, etc.). Humans and robots need to learn this by trial and error. Still, guided by some solid knowledge that came from a bootstrapping process this subsequent grounding process will be much faster than trying to learn any such skill from scratch.

Thus, structural bootstrapping may indeed begin to play an increasingly important role for robotic knowledge acquisition in the near future because it seems indeed that both of the above claims hold and that this set of methods will supersede other learning methods (but will have to still use grounding).

## Acknowledgements

The research leading to these results has received funding from the European Community’s Seventh Framework Programme FP7/2007-2013 (Specific Programme Cooperation, Theme 3, Information and Communication Technologies) under grant agreement no. 270273, Xperience.

## References

- [1] S. Pinker, Language Learnability and Language Development, Cambridge University Press, Cambridge, 1984.
- [2] J. Snedeker, Cross-Situational Observation and the Semantic Bootstrapping Hypothesis, in: E. Clark (ed.), Proc. 13th Ann. Child Language

Research Forum. Stanford, CA: Center for the Study of Language and Information, New York: John Wiley & Sons, 2002, pp. 445–496.

- [3] L. Gleitman, The structural sources of verb meanings, *Language Acquisition* 1 (1990) 3–55.
- [4] C. Fisher, Structural limits on verb mapping: the role of analogy in childrens interpretation of sentences, *Cogn Psychol* 31 (1996) 41–81.
- [5] L. Naigles, The use of multiple frames in verb learning via syntactic bootstrapping, *Cognition* 58 (1996) 221–251.
- [6] C. Fisher, L. Gleitman, Language acquisition, in: Pashler HF and Gallistel CR (eds.), *Steven’s Handbook of Experimental Psychology, Vol 3: Learning and Motivation*, New York: John Wiley & Sons, 2002, pp. 445–496.
- [7] L. Gleitman, Hard words, *Language Learning and Language Development* 1 (2005) 23–64.
- [8] J. Trueswell, L. Gleitman, Learning to parse and its implications for language acquisition, in: *Oxford Handbook of Psycholinguistics*, Oxford, 2007, pp. 635–656.
- [9] C. Fisher, Y. Gertner, R. M. Scott, S. Yuan, Syntactic bootstrapping, *WIREs Cognitive Science* 1 (2010) 143–149.
- [10] G. Chierchia, Syntactic bootstrapping and the acquisition of noun meanings: the mass-count issue, in: B. Lust, J. W. MArgarita Suner (Eds.), *Heads, Projections and Learnability Volume 1*, Hillsdale, New jersey, 1994, pp. 301–318.
- [11] F. Tracy, The language of childhood, *Am. J. Psychol.* 6 (1) (1893) 107–138.
- [12] T. Asfour, K. Regenstein, P. Azad, J. Schröder, N. Vahrenkamp, R. Dillmann, ARMAR-III: An Integrated Humanoid Platform for Sensory-Motor Control, in: *Humanoids*, Genova, Italy, 2006, pp. 169–175.
- [13] T. Asfour, N. Vahrenkamp, D. Schiebener, M. Do, M. Przybylski, K. Welke, J. Schill, R. Dillmann, ARMAR-III: Advances in Humanoid



Grasping and Manipulation, *Journal of the Robotics Society of Japan* 31 (4) (2013) 341–346.

- [14] T. Asfour, P. Azad, N. Vahrenkamp, K. Regenstein, A. Bierbaum, K. Welke, J. Schröder, R. Dillmann, Toward Humanoid Manipulation in Human-Centred Environments, *Robotics and Autonomous Systems* 56 (2008) 54–65.
- [15] R. Petrick, F. Bacchus, A knowledge-based approach to planning with incomplete information and sensing, in: *International Conference on Artificial Intelligence Planning and Scheduling (AIPS)*, 2002, pp. 212–221.
- [16] M. Steedman, *The Syntactic Process*, MIT Press, 2000.
- [17] E. E. Aksoy, A. Abramov, J. Dörr, K. Ning, B. Dellen, F. Wörgötter, Learning the semantics of object-action relations by observation, *The International Journal of Robotics Research* 30 (10) (2011) 1229–1249.
- [18] H. Xiong, S. Szedmak, J. Piater, Homogeneity Analysis for Object-Action Relation Reasoning in Kitchen Scenarios, in: *2nd Workshop on Machine Learning for Interactive Systems*, ACM, 2013, pp. 37–44, workshop at IJCAI. doi:10.1145/2493525.2493532.
- [19] A. J. Ijspeert, J. Nakanishi, H. Hoffmann, P. Pastor, S. Schaal, Dynamical movement primitives: Learning attractor models for motor behaviors, *Neural Computations* 25 (2) (2013) 328–373.
- [20] A. Ude, A. Gams, T. Asfour, J. Morimoto, Task-specific generalization of discrete and periodic dynamic movement primitives, *IEEE Trans. Robot.* 26 (5) (2010) 800–815.
- [21] T. Kulvicius, K. J. Ning, M. Tamosiunaite, F. Wörgötter, Joining movement sequences: Modified dynamic movement primitives for robotics applications exemplified on handwriting, *IEEE Transactions on Robotics* 28 (1) (2011) 145–157.
- [22] F. Wörgötter, A. Agostini, N. Krüger, N. Shylo, B. Porr, Cognitive agents – A procedural perspective relying on predictability of object-action complexes (OACs), *Robotics and Autonomous Systems* 57 (4) (2009) 420–432.

- [23] N. Krüger, C. Geib, J. Piater, R. Petrick, M. Steedman, F. Wörgötter, A. Ude, T. Asfour, D. Kraft, D. Omrčen, A. Agostini, R. Dillmann, Object-action complexes: Grounded abstractions of sensorimotor processes, *Robotics and Autonomous Systems* 59 (2011) 740–757.
- [24] E. E. Aksoy, A. Abramov, F. Wörgötter, B. Dellen, Categorizing object-action relations from semantic scene graphs, in: *IEEE International Conference on Robotics and Automation (ICRA)*, 2010, pp. 398–405.
- [25] M. Do, J. Schill, J. Ernesti, T. Asfour, Learn to wipe: A case study of structural bootstrapping from sensorimotor experience, in: *IEEE International Conference on Robotics and Automation (ICRA)*, 2014.
- [26] J. Ernesti, L. Righetti, M. Do, T. Asfour, S. Schaal, Encoding of periodic and their transient motions by a single dynamic movement primitive, in: *2012 12th IEEE-RAS International Conference on Humanoid Robots (Humanoids)*, Osaka, Japan, 2012, pp. 57–64.
- [27] A. Gams, M. Do, A. Ude, T. Asfour, R. Dillmann, On-line periodic movement and force-profile learning for adaptation to new surfaces, in: *2010 10th IEEE-RAS International Conference on Humanoid Robots (Humanoids)*, Nashville, TN, 2010, pp. 560–565.
- [28] P. Azad, T. Asfour, R. Dillmann, Combining harris interest points and the sift descriptor for fast scale-invariant object recognition, in: *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2009, pp. 4275–4280.
- [29] P. Azad, T. Asfour, R. Dillmann, Accurate shape-based 6-dof pose estimation of single-colored objects, in: *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2009, pp. 2690–2695.
- [30] P. Azad, D. Münch, T. Asfour, R. Dillmann, 6-dof model-based tracking of arbitrarily shaped 3d objects, in: *IEEE International Conference on Robotics and Automation (ICRA)*, 2011, pp. 0–0.
- [31] A. Bierbaum, M. Rambow, T. Asfour, R. Dillmann, Grasp Affordances from Multi-Fingered Tactile Exploration using Dynamic Potential Fields, in: *IEEE/RAS International Conference on Humanoid Robots (Humanoids)*, Paris, France, 2009, pp. 168 – 174.

- [32] S. Navarro, N. Gorges, H. Wörn, J. Schill, T. Asfour, R. Dillmann, Haptic object recognition for multi-fingered robot hands, in: IEEE Haptics Symposium, 2012, pp. 497–502.
- [33] J. Schill, J. Laaksonen, M. Przybylski, V. Kyrki, T. Asfour, R. Dillmann, Learning continuous grasp stability for a humanoid robot hand based on tactile sensing, in: IEEE RAS EMBS International Conference on Biomedical Robotics and Biomechatronics (BioRob), Rome, Italy, 2012, pp. 1901–1906. doi:10.1109/BioRob.2012.6290749.
- [34] J. Hockenmaier, M. Steedman, CCGbank: A Corpus of CCG Derivations and Dependency Structures Extracted from the Penn Treebank, *Computational Linguistics* 33 (3) (2007) 355–396.
- [35] T. Kwiatkowski, S. Goldwater, L. S. Zettlemoyer, M. Steedman, A probabilistic model of syntactic and semantic acquisition from child-directed utterances and their meanings, in: EACL, 2012, pp. 234–244.
- [36] E. Thomforde, M. Steedman, Semi-supervised CCG lexicon extension, in: EMNLP, 2011, pp. 1246–1256.
- [37] D. Teney, J. Piater, Continuous Pose Estimation in 2D Images at Instance and Category Levels, in: Tenth Conference on Computer and Robot Vision, IEEE, 2013, pp. 121–127. doi:10.1109/CRV.2013.34.
- [38] M. Ghazanfar, A. Prügel-Bennett, S. Szedmak, Kernel-Mapping Recommender System Algorithms, *Information Sciences* 208 (2012) 81–104. doi:10.1016/j.ins.2012.04.012.
- [39] L. Montesano, M. Lopes, A. Bernardino, J. Santos-Victor, Learning Object Affordances: From Sensory Motor Maps to Imitation, *IEEE Transactions on Robotics* 24 (1) (2008) 15–26. doi:10.1109/TRO.2007.914848.
- [40] W. Mustafa, N. Pugeault, N. Krüger, Multi-view object recognition using view-point invariant shape relations and appearance information, in: IEEE International Conference on Robotics and Automation (ICRA), 2013.
- [41] M. Thomsen, L. Bodenhagen, N. Krüger, Statistical identification of composed visual features indicating high-likelihood of grasp success., in: Workshop 'Bootstrapping Structural Knowledge from Sensory-motor

- Experience. IEEE International Conference on Robotics and Automation (ICRA), 2013.
- [42] S. Schaal, P. Mohajerian, A. Ijspeert, Dynamics systems vs. optimal control – a unifying view, *Progress in Brain Research* 165 (6) (2007) 425–445.
  - [43] A. Billard, S. Calinon, F. Guenter, Discriminative and adaptive imitation in uni-manual and bi-manual tasks, *Robot. Auton. Syst.* 54 (2006) 370–384.
  - [44] B. D. Argall, S. Chernova, M. Veloso, B. Browning, A survey of robot learning from demonstration, *Robotics and Autonomous Systems* 57 (5) (2009) 469–483.
  - [45] R. Dillmann, T. Asfour, M. Do, R. Jäkel, A. Kasper, P. Azad, A. Ude, S. R. Schmidt-Rohr, M. Lösch, Advances in robot programming by demonstration, *KI - Künstliche Intelligenz* 24 (4) (2010) 295–303.
  - [46] F. Guerin, N. Krüger, D. Kraft, A survey of the ontogeny of tool use: from sensorimotor experience to planning, *IEEE TAMD* 5 18 – 45.
  - [47] A. Agostini, M. J. Aein, S. Szedmak, E. E. Aksoy, J. Piater, F. Wörgötter, Using structural bootstrapping for object substitution in robotic executions of human-like manipulation tasks, in: *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, Hamburg, Germany, 2015, p. submitted.
  - [48] M. Wächter, S. Schulz, T. Asfour, E. Aksoy, F. Wörgötter, R. Dillmann, Action Sequence Reproduction based on Automatic Segmentation and Object-Action Complexes, Atlanta, USA, 2013.
  - [49] A. Gams, A. Ijspeert, S. Schaal, J. Lenarčič, On-line learning and modulation of periodic movements with nonlinear dynamical systems, *Autonomous Robots* 27 (1) (2009) 3–23.
  - [50] L. Cuiyan, Z. Dongchun, Z. Xianyi, A survey of repetitive control, in: *IEEE/RSJ International Conference on Robots and Systems (IROS)*, Sendai, Japan, 2004, pp. 1160–1166.

- [51] A. Gams, J. van den Kieboom, M. Vespignani, L. Guyot, A. Ude, A. Ijspeert, Rich periodic motor skills on humanoid robots: Riding the pedal racer, in: IEEE International Conference on Robotics and Automation (ICRA), Hong Kong, China, 2014.
- [52] S. Schaal, Is imitation learning the route to humanoid robots?, *Trends in Cognitive Sciences* 3 (1999) 233–242.
- [53] M. Pardowitz, S. Knoop, R. Dillmann, R. D. Zöllner, Incremental Learning of Tasks From User Demonstrations, Past Experiences, and Vocal Comments, *IEEE Transactions on Systems, Man and Cybernetics – Part B: Cybernetics* 37 (2) (2007) 322–332.
- [54] S. Ekvall, D. Kragic, Robot learning from demonstration: a task-level planning approach, *International Journal of Advanced Robotic Systems* 5 (3) (2008) 223–234.
- [55] R. Cubek, W. Ertel, Learning and Execution of High-Level Concepts with Conceptual Spaces and PDDL, in: 3rd Workshop on Learning and Planning, ICAPS (21st International Conference on Automated Planning and Scheduling), 2011.
- [56] B. Nemec, R. Vuga, A. Ude, Exploiting previous experience to constrain robot sensorimotor learning, in: Proc. 11th IEEE-RAS Int. Conf. Humanoid Robots, 2011, pp. 727–732.
- [57] K. Kronander, M. Khansari-Zadeh, A. Billard, Learning to control planar hitting motions in a minigolf-like task, in: Proc. 2011 IEEE/RSJ Int. Conf. Intelligent Robots and Systems, 2011, pp. 710–717.
- [58] J. Kober, A. Wilhelm, E. Oztop, J. Peters, Reinforcement learning to adjust parametrized motor primitives to new situations, *Auton. Robots* 33 (4) (2012) 361–379.
- [59] J. Piaget, *The Origins of Intelligence in the Child*, Routledge, London, New York, 1953.
- [60] K. Hirsh-Pasek, R. M. Golinkoff (Eds.), *Action Meets World: Now Children Learn Verbs*, Oxford University Press, 2006.
- [61] A. Tate, Generating project networks, in: IJCAI, 1977, pp. 888–893.

- [62] E. D. Sacerdoti, Planning in a hierarchy of abstraction spaces, *Artif. Intell.* 5 (2) (1974) 115–135.
- [63] K. Erol, J. A. Hendler, D. S. Nau, HTN planning: Complexity and expressivity, in: *AAAI*, 1994, pp. 1123–1128.
- [64] E. E. Aksoy, M. Tamosiunaite, R. Vuga, A. Ude, C. Geib, M. Steedman, F. Wörgötter, Structural bootstrapping at the sensorimotor level for the fast acquisition of action knowledge for cognitive robots, in: *IEEE International Conference on Development and Learning and Epigenetic Robotics (ICDL-EPIROB)*, 2013.
- [65] F. Wörgötter, E. E. Aksoy, N. Krüger, J. Piater, A. Ude, M. Tamosiunaite, A simple ontology of manipulation actions based on hand-object relations, *IEEE Transactions on Autonomous Mental Development* 5 (2) (2013) 117–134.
- [66] N. Badler, Temporal scene analysis: Conceptual descriptions of object movements, Ph.D. thesis, University of Toronto, Canada (1975).
- [67] M. Sridhar, G. A. Cohn, D. Hogg, Learning functional object-categories from a relational spatio-temporal representation, in: *Proc. 18th European Conference on Artificial Intelligence*, 2008, pp. 606–610.
- [68] H. Kjellström, J. Romero, D. Kragić, Visual object-action recognition: Inferring object affordances from human demonstration, *Comput. Vis. Image Underst.* 115 (1) (2011) 81–90.
- [69] Y. Yang, C. Fermüller, Y. Aloimonos, Detection of manipulation action consequences (mac), in: *International Conference on Computer Vision and Pattern Recognition (CVPR)*, (In press), 2013.
- [70] G. Granlund, The complexity of vision, *Signal Processing* 74.
- [71] N. Krüger, P. Janssen, S. Kalkan, M. Lappe, A. Leonardis, J. Piater, A. J. Rodríguez-Sánchez, L. Wiskott, Deep hierarchies in the primate visual cortex: What can we learn for computer vision?, *IEEE PAMI* 35 (8) (2013) 1847–1871.
- [72] Acat project webpage (2014).  
URL <http://www.acat-project.eu/>

- [73] Roboearth project webpage (2014).  
URL <http://roboearth.org/>
- [74] Robohow project webpage (2015).  
URL <http://robohow.eu/>