



**ICT Call 7  
ROBOHOW.COG  
FP7-ICT-288533**

**Deliverable D6.2:**

**Description of the design and realization of the temporal  
projection mechanisms**



**January 31st, 2014**

Project acronym: ROBOHOW.COG  
Project full title: Web-enabled and Experience-based Cognitive Robots that Learn Complex Everyday Manipulation Tasks

Work Package: WP 6  
Document number: D6.2  
Document title: Description of the design and realization of the temporal projection mechanisms  
Version: 1.0

Delivery date: January 31st, 2014  
Nature: Report  
Dissemination level: Restricted (RE)

Authors: Jan Winkler (UNIHB)  
Michael Beetz (UNIHB)

The research leading to these results has received funding from the European Union Seventh Framework Programme FP7/2007-2013 under grant agreement n<sup>o</sup>288533 ROBOHOW.COG.

# Summary

Work package 6 is to design, implement, and empirically analyze temporal projection mechanisms that enable a cognitive agent to predict the outcome of its actions in different grades of accuracy. The predictions will be based on physical simulations of an abstract model of the real world scenario at hand, and will give the agent hints on the feasibility of manipulation actions. Two gradually different approaches are surveyed, which cover a crude, object-based simulation environment, reacting to dynamic changes based on robot manipulation, and a fine grained fluid simulation based approach, showing the result of pouring actions during house work chores such as cooking. Work package 6 investigates and develops prediction techniques based on such simulation environments, taking processing times into account. While a crude, object-based approach leads to on-the-fly results for fast decision making during task execution, a more detailed simulation enables the agent to analyze a more realistic outcome of composed and continuous events, such as fluid pouring, before and after performing the task. Temporal projection of planned actions enables an agent to envision different scenarios in order to decide about the best course of action and action parameterizations. This includes qualitative information such as appropriate put-down positions for objects, locations to stand at in order to perform an action, but also information of more quantitative nature, such as the tilting angle for containers when pouring a liquid pancake mix onto a heater oven.

Based on the results gathered from such predictions, vague action descriptions can be refined, and fundamental flaws during plan execution can be prevented, and general plans can, to a certain degree, disambiguated for the situation at hand.

In the second year, we have focussed on two problem instances in which we employ temporal projection mechanisms for making action decisions:

- First, generalized “fetch and place” plans which include multiple vague location, object, and action descriptions that must be parameterized, can be disambiguated during plan execution. Consisting of very general plan instances such as *“fetch the cup from the table”*, a playthrough using an abstract physics simulation model of the real world allows for trying out strategies to solve such plans. Sequences of actions leading to a desired outcome are then recalculated based on events and action effects in the real world.
- Second, when performing actions that must be parameterized rather precisely on a quantitative level, physics simulations lead to action parameters for achieving a specific desired outcome. We investigated the problem of pouring a certain amount of liquid onto a heater oven without spilling it. Using a dynamic fluid simulation, liquids are represented as graph-based models that can be used to identify parameters fit for the required effect.

**Simulating Action Sequences for Fast Plan Parameterization** Robot plans designed for ambiguous tasks, such as pick and place plans, must be applicable to a very wide range of situations to be useful. Especially parameters like the place from where an object is collected, and where it is put, are strongly situation-dependent. To allow a cognitive agent to decide on a first guess of how to overcome this ambiguity problem, it needs an understanding of the effects of its own actions. The multiple problem dimensions during this seemingly simple task include a place to stand at to perform an action (reaching for an object), but also the effect of manipulative tasks, such as putting an object on a table (i.e. shoving other objects off the table, finding a stable position for the new object). A physical simulation, based on the real world situation, lets the robotic agent try out action sequences to a certain degree, and dismiss possible action routes it assumes to be invalid based on such results. These aspects are explained in terms of a pick and place scenario in [Mösenlechner and Beetz, 2013], as attached to this deliverable.

**Generating Action Parameterizations from Simulation of Atomic Tasks** Performing tasks in the real world requires action parameters. Such parameters depend on the situational context, and especially the task at hand and its desired outcome. In order to evaluate the outcome of such an atomic task, an elaborate physics simulation is employed to identify and dismiss possibly undesired effects before actually performing the task. A notably difficult task is the handling, and pouring, of liquids. In order to show the relevance of a physics based simulation for such tasks, the pouring of a liquid pancake mix onto a heater oven is elaborated on. During this precisely carried out simulation, physical effects on subjects of interest, such as spilling the mix, or pouring the wrong amount, are detected. The design, implementation, and application of such mechanisms is carefully elaborated on in [Klapfer et al., 2012], as attached to this deliverable.

# Contributed Papers

Papers included in this deliverable are:

- Mösenlechner, L., Beetz, M. Fast Temporal Projection Using Accurate Physics-Based Geometric Reasoning. In Proceedings of the IEEE International Conference on Robotics and Automation (ICRA), 2013.
- Kunze, L., Klapfer, R., Beetz, M. Beyond Pick-and-Place — Naive Physics Reasoning for Handling Fluids in Everyday Manipulation. In Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS), 2013.

# Fast Temporal Projection Using Accurate Physics-Based Geometric Reasoning

Lorenz Mösenlechner  
Intelligent Autonomous Systems  
Technische Universität München, Germany  
moesenle@cs.tum.edu

Michael Beetz  
Institute for Artificial Intelligence  
University of Bremen and TZI\*  
beetz@cs.uni-bremen.de

**Abstract**— Temporal projection is the computational problem of predicting what will happen when a robot executes its plan. Temporal projection for everyday manipulation tasks such as table setting and cleaning is a challenging task. Symbolic projection methods developed in Artificial Intelligence are too abstract to reason about how to place objects such that they do not hinder future actions. Simulation-based projection is fine-grained enough but computationally too expensive as it is not able to abstract away from the execution of uninteresting actions (such as navigation). In this paper we propose a novel temporal projection mechanism that combines the strengths of both approaches: it is able to abstract away from the execution of continuous but uninteresting actions and provides the realism and fine grainedness needed to reason about critical situations.

## I. INTRODUCTION

When robots perform everyday manipulation activities such as setting the table or cleaning up, the difficulty of individual pick and place tasks critically depends on where the robot exactly places the objects and on the order in which the robot puts and takes away objects. If the robot places objects too early or at inappropriate places, these objects might hinder subsequent pick and place tasks.

Robot action planning [1], which is the computational task of deciding on the course of action based on predicting and reasoning about the future consequences of the intended plan, aims to avoid such problems. Unfortunately, symbolic action planning methods, mostly based on PDDL (Planning Domain Description Language) derivatives ([2], [3]), are too coarse grained to be informative with respect to whether a reaching or placing action is easy or challenging [4]. Motion planning ([5], [6]) is fine grained enough but very limited in the ways to reason about how scenes can be manipulated in order to simplify actions (but see [7] for a notable exception). Several approaches define more realistic models of world states that consider reachability as the existence of a motion plan [8] or models that are learned from experience [9] but these approaches are limited to navigation actions. Simulation-based projection methods ([10], [11]) are fine grained enough to predict and reason about such effects but they are very resource intensive.

In this paper, we propose a novel robot plan projection method that combines abstract symbolic projection for actions with detailed geometric and simulation-based reasoning for handling the action aspects described above. The system we present allows for inferring plan parameters such as

locations for putting down objects under the constraints defined by the current task and future actions. Instead of projecting a plan by applying a sequence of logical rules that update a symbolic world state, our system uses a geometric representation of the world to calculate predicates such as *Visible* and *Reachable* on demand. The whole system is integrated in CRAM [12] and is based on plans that are specified in the CRAM plan language. Similar to purely symbolic approaches for projection, actions are implemented as symbolic rules for updating the internal representation of the world. However, by using a geometrically accurate 3D representation of the world our approach provides a good compromise between fast but very abstract purely symbolic projection and computationally expensive but very accurate simulation of complete actions. This paper is based on the work presented in [13]. The main contributions are extensions to support temporal reasoning based on time lines, the execution of plans in projection to generate these time lines and the definition of behavior flaws.

This paper is structured as follows. First we give an overview of our system that is capable of generating plan parameters such as locations for placing objects or where to stand not only based on static assumptions of the world but based on the future course of actions. Then we introduce the physics-based reasoning engine that we use, followed by an explanation of reasoning about plans, time line generation and the temporal calculus we use to represent behavior flaws. Finally we demonstrate the expressiveness the system by defining various behavior flaws.

## II. RELATED WORK

Related publications in the area of planning using geometry and physics simulation include [14] where the authors integrate a physics engine to calculate state transitions in planning. In [15], the authors combine symbolic and geometric planning by calculating predicate values for a high level planner using a geometric planner. However, the authors in both publications do not integrate high-level concepts such as reachability or visibility. Planning under uncertainty, integrating a geometric representation of the world including free and occluded areas has been shown in [16]. In [17], visibility simulation is integrated into a motion planner and in [18] the authors show similar visibility calculation as presented in this paper in the context of human-robot interaction and perspective taking.

\*The Institute of Artificial Intelligence is part of both the University of Bremen and the Centre for Computing Technologies (TZI).

While all these publications show the implementation of reasoning in geometric domains, none of them provides means for *generating* parameters such as destination poses for objects or poses for the robot to stand. In [19], the authors present capability maps, i.e. pre-calculated maps to quickly check if poses are reachable or where to place a robot to be able to reach a specific pose.

Temporal projection is a well studied field in formal logic. Given a sequence of actions, temporal projection tries to infer the state of the world after executing these actions. In [20] and [21], the authors deal with the problem of uncertain knowledge about the initial state of the world. McDermott introduces the generation of time lines, similar to the work presented in this article. However, purely symbolic approaches are often too abstract to represent geometric properties of the world, for instance occlusions. While simulation based projection as presented in [11] and the authors' previous work in [10] provide similar functionality for geometric reasoning and reasoning about actions as presented in this paper, they suffer from high computational complexity and extremely long run times.

### III. SYSTEM OVERVIEW

The purpose of the system described in this paper is to generate poses for placing the robot's base while manipulating objects and for placing objects, for instance on a counter. The system not only uses the current state of the world to generate such poses but also takes into account future actions of the current plan.

In CRAM plans, locations, objects and actions are specified using designators, instances that are built from conjunctions of symbolic constraints specified as key-value-pairs. For instance, we specify a location for a plate on the counter that is visible and reachable for the robot as follows:

(a location (for plate) (on counter) (reachable-for robot) (visible-for robot))

To resolve these designators, they are compiled into Prolog programs and solutions for them are generated by executing them. To improve readability, instead of writing actual Prolog programs, in this paper we will use a first-order representation that is equivalent to the corresponding Prolog code. The above designator is compiled to the following logical expression:

$$\begin{aligned} & \text{PosesOn}(\text{Counter}, \text{Cup}, ?\text{Poses}) \wedge \text{Member}(?\text{Pose} ?\text{Poses}) \\ & \wedge \text{AssertPose}(\text{Cup}, ?\text{Pose}) \wedge \text{Reachable}(\text{Cup}, \text{Robot}) \\ & \wedge \text{Visible}(\text{Cup}, \text{Robot}) \end{aligned}$$

The *PosesOn* predicate generates candidate poses by drawing samples from a probability distribution that is built from the designator. If all predicates hold for such a candidate pose, it is a valid pose under the constraints defined in the designator and it is considered a solution. In addition to explicitly defined constraints, as shown in the example above, a number of implicit constraints need to hold and are added to the Prolog program. For instance, an object of interest must stand stable at a put-down location, i.e. it must not be in collision with any other objects. In contrast to classical Prolog, the predicates used in designator resolution are computed by

using an accurate simulated 3D representation of the world as a Prolog database. For instance, if the Prolog inference process needs to prove if the *Stable* predicate holds for a specific object instance, a physics simulation is used to check if the object is moving at its current location. If not, it is considered stable and the stable predicate holds.

While proving predicates in a static world database finds locally valid solutions, it does not take into account future actions and thus cannot infer if a specific solution negatively influences the performance of future actions. However, by projecting a plan, and analyzing the generated time line, we can evaluate if specific designator solutions will cause problems in future actions.

Projection is the execution of a plan in a lightweight simulation environment. It generates a time line that allows for reconstructing different world states along the course of actions that can be reasoned about. By matching behavior flaws, i.e. logical definitions of conditions in the world that should be avoided on these time lines, we can define an objective function and evaluate the quality of different solutions of a specific designator.

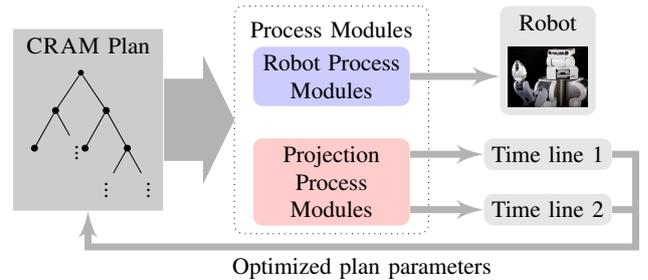


Fig. 1. System overview of plan projection for designator resolution.

Let us consider a simple example: the robot's task is to grasp a cup on the counter and put it on the table and grasp a plate that is on the table and place it on the counter. Apart from being on the counter or table respectively, the put-down locations are unspecified. However, good solutions will increase plan performance. For instance, if the put down location for the first object is close to the other object, the robot does not need to move to grasp the second object. However, the put down location of the cup should not be too close or in front of the plate since it would block the pick up action because the plate cannot be detected by perception or it cannot be grasped because the grasp position needed by the pick-up action is blocked by the cup.

Figure 1 gives an overview of the system design. When the system starts to execute a plan, a number of projections are started in parallel. Each projection spawns off from the current execution and world state, i.e. if the cup has already been grasped, all subsequent projections start with the cup in the robot's gripper. Each projection generates a time line. The time line contains events that indicate changes in the simulated world and for each event, we store a copy of the simulated world state for later reasoning and analysis. In addition to this time line, a separate task tree is stored for each projection that contains information about plan

execution such as which task has been executed when and why, what the task’s status at a given time was and what the outcome of the task was, including information about failures. When one projection thread finishes, the time line is evaluated by matching predefined flaw specifications. Based on the result, a performance value for all location designator solutions is generated. If the solutions in one episode are better than the previous ones, they are cached and used by the actual plan that is executed on the robot. Examples for behavior flaws include the distance the robot has to drive between actions, if objects that are to be manipulated are visible and if objects are blocking other objects.

#### IV. PHYSICS-BASED REASONING

Projection and inference of flaws is based on a physics based reasoning engine as explained in [13]. Instead of proving all predicates based on a purely symbolic fact base, the truth values and bindings of certain predicates are calculated by using the Bullet physics engine, OpenGL off-screen rendering and inverse kinematics calculation and simplified inverse reachability map calculation. We use physics based reasoning instead of a purely symbolic knowledge base because in our domain of a human household, a symbolic representation would be too abstract. Many problems could not be solved. For instance, if an object is visible or not for the robot not only depends on the current location of the robot and its distance from the object but also on all other objects and their shape. In addition, since we need to find solutions for location designators, we need to have a generative model that yields poses that are valid solutions for the respective designator. We solve this by first drawing random samples based on a probability distribution that is generated from the designator’s constraints and then using the physics based reasoning engine to prove the validity of such a sample, considering all explicit and implicit constraints of the designator.

For reasoning, we use a three-dimensional representation of the environment the robot operates. This representation contains all information the robot has about the environment. All static objects such as cupboards, the refrigerator and counter tops are provided by a semantic environment map as described in [22]. The robot’s perception routines assert objects in the database when they are seen. The world database is kept consistent by removing objects that should have been visible but could not be seen anymore. The robot is equipped with a huge database of 3D meshes for objects including the PR2’s household objects database <sup>1</sup>.

To explain the inference process, let us consider a simple location designator for placing a cup on the counter top at a location where it is visible for the robot from its current location. The following first-order logic expression shows how that designator can be resolved:

$$\text{PosesOn}(\text{Counter}, \text{Cup}, ?\text{Poses}) \wedge \text{Member}(?\text{Pose}, ?\text{Poses}) \\ \wedge \text{AssertPose}(\text{Cup}, ?\text{Pose}) \wedge \text{Visible}(\text{Cup}, \text{Robot})$$

<sup>1</sup>[http://www.ros.org/wiki/household\\_objects\\_database](http://www.ros.org/wiki/household_objects_database)

The inference engine processes the different terms sequentially and tries to find bindings for yet unbound variables. If a predicate fails, the engine backtracks and retries with a different solution for an unbound variable until a solution has been found or fails if all possible bindings for a variable have been tried. The first step is to generate a sequence of pose candidates that are possible solutions for the designator. The predicate *PosesOn* generates a probability distribution as shown in Figure 2. Since all poses on the counter top might be valid solutions, each pose on the counter has equal probability. The *PosesOn* predicate binds the (virtual) list of all samples to the variable *?Poses*.

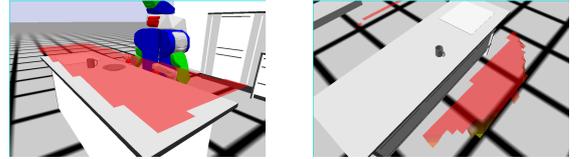


Fig. 2. Probability distributions used for sampling poses that are on the counter top (left) and for the robot to stand to reach the cup (right). The left distribution has an equal probability for all poses on the counter top while the right distribution is based on an inverse reachability map for the robot.

The next step is to draw one sample from the sequence of poses and to put the cup at the corresponding pose. Finally, the *Visible* predicate calculates if the object is visible for the robot at its current pose by rendering the scene from the position of the robot’s camera and counting how many pixels of the object are visible.

For generating locations for objects and the robot and for inferring flaws on time lines, we define predicates for *collisions*, *stability* reasoning, *visibility* reasoning and *occlusions* and *reachability* reasoning and inference of objects that might be *blocking* a grasp action. The following list gives an overview of the predicates we implemented.

- *Contact*(?*W*, ?*O*<sub>1</sub>, ?*O*<sub>2</sub>) holds if the two objects ?*O*<sub>1</sub> and ?*O*<sub>2</sub> are in contact in the world *W*.
- *Stable*(?*W*, ?*O*) holds if all forces on object ?*O* are canceled out, i.e. it does not move in the corresponding world ?*W*.
- *Visible*(?*W*, ?*P*, ?*O*) holds if the object ?*O* is visible from pose ?*P* in the world ?*W*.
- *Occluding*(?*W*, ?*P*, ?*O*<sub>1</sub>, ?*O*<sub>2</sub>) holds if object ?*O*<sub>2</sub> is occluding the object ?*O*<sub>1</sub> when the camera is at pose ?*P*.
- *Reachable*(?*W*, ?*R*, ?*O*) holds if the robot ?*R* can reach the object ?*O* in the world configuration ?*W*.
- *Blocking*(?*W*, ?*R*, ?*O*, ?*B*) unifies ?*B* with the list of objects that might be blocking a grasp for object ?*O*.

As can be seen, all predicates require a world database as first variable. By letting this parameter unbound, a default database is used. In this article we can only briefly explain the implementation of these predicates. Details can be found in [13]. The predicates *Contact* and *Stable* are implemented using the Bullet physics engine<sup>2</sup>. Contacts are inferred by using the Bullet’s collision engine and for inferring stability, we simulate for a short period of time and compare the poses of objects to check if they moved. If an object changed its location, it is not stable. The predicates *Visible* and *Occluding* are inferred using OpenGL. Each object is rendered in a

<sup>2</sup><http://bulletphysics.org/>

different color and by counting pixels, the system can infer how much of an object is visible and by which other objects it is occluded. Reachability is calculated using a standard solver for inverse kinematics and blocking objects are objects the robot is in collision with when reaching for an object. Blocking objects are objects that potentially hinder grasping actions and make motion planning harder.

## V. TEMPORAL PROJECTION OF CRAM PLANS

The reasoning system explained in Section IV has no notion of time. However, projection requires the integration of reasoning about sequences of actions over time and their consequences and we need to extend the system by integrating a temporal calculus.

In CRAM, we consider robot control programs that have annotations in a first-order logic as plans because the annotations allow us to reason about their purpose. Projecting a plan means executing it in projection mode. In projection, instead of interacting with the actual robot's hardware, plans interact only with the world database that we use to implement our physics-based reasoning engine as described in the previous section. Each action generates a sequence of events where each event indicates a change in the world database at a specific point in time. By storing the sequence of events and copies of the corresponding world database at the time the event happened, we construct a time line. By adding predicates for temporal reasoning on these time lines, including predicates for relating events to parts in the plan, we can specify flaws in the robot's behavior as logical expressions.

In this section we show how plans must be constructed to allow for projecting them, how time lines are generated and how reasoning about them is implemented.

### A. Reasoning About Plan Execution

In CRAM, plans are robot control programs that cannot only be executed but also reasoned about. Programmatically inferring what an arbitrary control program does is intractable at best if not impossible and it gets even more difficult in highly dynamic changing environments such as a human household. However, by providing annotation describing the semantics of a certain part of a plan, reasoning about the semantics of a complete plan becomes possible.

Execution of a CRAM plan generates a plan tree, i.e. a tree that contains all plans and sub-plans that have been executed, their status at any point in time, when they started, when they terminated and their result or failure description. The semantics of a sub-plan are made transparent for a reasoning engine by providing semantic annotations. For instance, a plan that is supposed to place an object *?obj* at a location *?loc* is called *Achieve(Loc(?obj, ?loc))*. The implication of a terminal status *succeeded* is that the robot believes that the object is at the destination location. The term *Loc(?obj, ?loc)* is called an occasion and is defined in the system's reasoning engine. In other words, using this naming scheme, we define the purpose of plans using predicates of an underlying reasoning system which allows for relating plans to the corresponding world states that are stored on the time line.

To reason about plan execution, we define predicates that query the plan tree. For instance, we define the predicate *Task(?tsk)* to unify a variable with a task object. The predicate *TaskGoal(?tsk, ?goal)* can be used to query a task object for its goal where the goal is for instance an achieve expression as introduced before. In addition, the system provides predicates to access the task's status, result, errors, position in the tree including sub-tasks and start and times. A complete list and a more detailed explanation of the system can be found in [23].

### B. Generation of Timelines

CRAM plans need to be written in a general way to allow for executing the same code on the actual robot hardware as well as in projection. This is achieved by a clear, well-defined, minimal and most importantly purely symbolic interface between high-level plans and the low-level components of the robot or projection that execute the actual actions. All actions that are to be performed by the robot are described by action designators, i.e. key-value-pairs specifying the action's parameters in order to execute it. For instance, the following action designator is used to grasp a cup:

(an action (type navigation) (goal (location (to see) (obj Cup))))

The input for all process modules as shown in Figure 3 are action designators. Each process module is activated at the beginning of plan execution and deactivated after a plan finished. The status can be monitored using special signal slots and process modules update the world database using events.

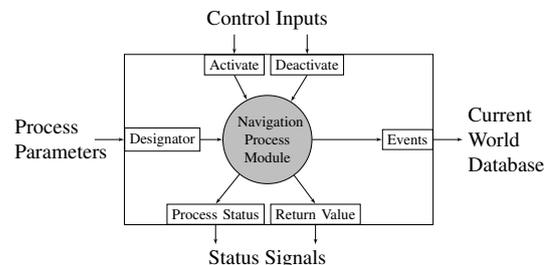


Fig. 3. Process module encapsulating a navigation control process. The input is an action designator specifying that the action is a navigation action and the containing the goal location as represented by a location designator, e.g. *(a location (to see) (obj cup))*.

From the point of view of a high-level plan, process modules can be seen as black boxes. That allows us to keep the high-level plan general enough to work with process modules developed for different robot platforms, simulation or projection. Specifically, when projecting a plan, we implement a process module that just makes assertions and retractions in a projection-local copy of the world database of the reasoning engine instead of sending commands to the actual robot hardware. Depending on the application and the robot, the number of process modules may vary. On the PR2, we currently implement four process modules: manipulation, navigation, perception and a process module for moving the cameras mounted on the robot's head.

To understand how projection and the creation of a time line works, let us have a closer look at the manipulation projection process module as an example for a rather complex process module. In the current implementation it supports the actions *grasp*, *lift*, *put down*, *open* and *close*. When the process module receives a respective action designator it makes assertions in the world database, emits events and increments the simulation time. The following example shows an action designator for grasping a cup:

(an action (to grasp) (obj (object (type cup))))

When the process module executes it, it generates the events *RobotStateChanged()* and *ObjectAttached(Cup, "r\_gripper\_wrist\_link")*. Since projection is supposed to be fast, we do not perform a complete simulation of a trajectory which would also require motion and grasp planning. Instead, we define *key points* for trajectories. For grasping, we define only one key point, the final pose the robot needs to reach in order to grasp the object. For simplicity, we use one out of the following grasps: a side grasp, a front grasp or a top grasp. Which grasp is used depends on the type of the object that should be grasped. For instance, cups are grasped with front or side grasps while plates always require a side grasp. After the *RobotStateChanged* event, the robot's gripper is at the cup's pose. The assertion in the world database that corresponds to the *RobotStateChanged* event positions the links of the robot's arm according to an inverse kinematics solution for placing the gripper at the cup's pose. The *ObjectAttached* event notifies the system that the object should be moved whenever the gripper link moves as long as the object is attached.

ObjectPerceived(?obj, ?s)	An object ?obj (described by an object designator) has been seen in a specific sensor ?s.
RobotStateChanged()	The robot has changed its state, i.e. it changed its position or the position of some links.
ObjectAttached(?obj, ?link)	An object ?obj has been attached to a specific link.
ObjectDetached(?obj, ?link)	An object ?obj has been detached from a specific link.
ObjectArticulationEvent(?obj, ?d)	An object changed its articulation status, e.g. a drawer or a cupboard has been opened by distance ?d.
ActionStarted(?m, ?d)	The process module ?m started to execute the action designator ?d
ActionFinished(?m, ?d)	The process module ?m finished executing the action designator ?d

TABLE I

OVERVIEW OF THE EVENTS USED FOR GENERATING A TIME LINE.

Table I gives an overview of the events that are used in the current system. The perception process module only generates *ObjectPerceived* events. The manipulation process module generates events of type *RobotStateChanged*, *ObjectAttached*, *ObjectDetached* and *ObjectArticulationEvent*. Navigation and PTU (Pan-Tilt-Unit) only generate *RobotStateChanged* events. When they start executing, all process modules generate the event *ActionStarted* and when they finish they generate the event *ActionFinished*.

Handling of time and incrementing time is especially critical in projection. While we want to produce predictions of the outcome of a specific plan quickly, we must not sacrifice the ability to project actions that happen concurrently. For instance, suppose the robot tries to see an object on the table while it is putting down an object next to it. If the destination location for the put-down action happens to be in front of the

object, the run time of both actions determines if perception succeeds or not and projection needs to account for that.

Each set of changes in the world database must cause an increment of the projection clock. This assumes that the time spent while executing code in a high-level plan takes almost no time. In fact, executing actual actions on the robot is the most time consuming part. For instance, a navigation action will take several seconds depending on the distance to drive and planning and executing an arm trajectory can take 30 seconds and more while triggering these actions in a high level plan can be done in fractions of seconds. By only allowing time increments in process modules, we keep plans general since they do not have to have explicit support for projection. The specific projection clock we use increments at a maximal rate, for instance 20 milliseconds. In other words, increments happen at most at 50 Hz. However, the amount of time by which the clock is incremented can vary depending on the action that is projected. That way, each action will take exactly 20 milliseconds real time while still adding events to the time line at approximations of the action's run time. This implementation is a good compromise to allow for concurrent execution and good performance. The projected run time of actions is calculated using heuristics and random numbers. For instance, navigation time is approximated by a linear function over the distance between the starting point of the action and the goal and noise generated by a random number generator. The same is done for manipulation while perception is modeled using a constant function. The heuristics are not deterministic to increase the variation in different projections of the same plan.

### C. Reasoning on Timelines

After a plan has been projected, the system needs to analyze the time line and search for flaws. We define two predicates, *Occurs* for reasoning about events and *Holds* for reasoning about world states over intervals of time. Figure 4 shows an example time line generated by a simple pick-up plan. The robot navigated to a location close to the table and picked up a cup.

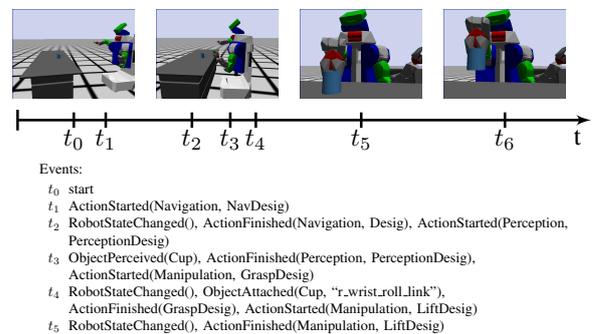


Fig. 4. Sample projection of a simple pick-up plan. The robot first navigates to the table, then moves the gripper to the location of the cup and finally lifts the cup.

The time line in our example contains four events. First the robot navigates to a location where it can reach the cup. Nav-

igation generates a *RobotStateChanged* event since the robot changed its location in the map. Perceiving the object generates an *ObjectPerceived* event. Grasping generates another *RobotStateChanged* event since the arm moved. In addition, the object is attached to the right gripper. Finally, lifting moves the robot arm again, so another *RobotStateChanged* event is generated. As already mentioned, the ordered sequence of events creates a time line. More specifically, the time line not only contains the events that were generated in the process modules and corresponding time stamps but also snapshots of the world database at the time when an event occurred. In contrast to classical purely symbolic projection, our approach does not require logical expressions that represent the world state to be added or removed when actions are performed. Instead, we calculate the truth values and variable bindings of predicates on demand, based on the stored geometric world databases. This approach avoids certain problems purely symbolic approaches suffer from. For instance, consider the predicate *Visible* as defined in the previous section before and after a put-down action. The put-down action might cause occlusions. In a purely symbolic representation of the world that only consists of logical formulas, it is hard to infer that an object is occluded by another object after a put down action, in particular because the put-down action only contains a reference to the manipulated object. In contrast, our system proves visibility using OpenGL and on demand and the problem of deciding which conditions do not hold anymore after executing an action is avoided.

To reason about events, we define the predicate *Occurs(?event, ?time)*. It unifies an event with a time stamp. For instance, to find all objects that have been grasped and the corresponding times we can state the following query:

*Occurs(ObjectAttached(?o, ?l), ?t)*

The resulting bindings of this query in the above example timeline are:

$?o = \text{Cup}, ?l = \text{"r\_wrist\_roll\_link"}, ?t = t_2$

The implementation of the *Occurs* predicate just iterates over the sequence of events on the time line and unifies the variables *?event* and *?time* with the corresponding event pattern and time stamp.

Events are transitions in the world database, i.e. they indicate changes in the world. In other words, occasions (i.e. states) hold over time intervals and change only at events. For instance, if an object has been attached, the Prolog expression *ObjectInHand(?object)* holds until the object is detached again. To reason about these conditions and the corresponding time intervals they hold in, we define the predicate *Holds(?occasion, ?interval)*. The system tries to prove a given logical expression *?occasion* at the time points described by *?interval*. Please note that *?occasion* and *?interval* cannot be free variables since we cannot enumerate all possible logical expressions and time intervals. However, *?occasion* and *?interval* can be terms that contain free variables. *?interval* must be one of the following three

expressions:

- *during*( $t_0, t_1$ ) indicating that the occasion must hold at least once in the interval  $[t_0, t_1)$ .
- *at*( $t$ ) for specifying a single point in time, i.e. an interval with length 0.
- *throughout*( $t_0, t_1$ ) indicating that the occasion must hold throughout the complete time interval  $[t_0, t_1)$ .

The implementation is based on proving occasions in the different world databases stored together with the events on the time line, assuming that changes in the world database only happen with events. As an example, if we want to find objects that were standing on the table at time  $t$  with  $t_2 < t < t_3$ , the corresponding query can be formulated as follows:

*Holds(On(?object, Table), at(t))*

The *Holds* predicate then finds and loads the world database instance on the time line whose time stamp is directly before  $t$  and proves the term *On(?object, Table)*. The implementation of *during* iterates over all valid world databases in the interval and proves the occasion in each of them, yielding all possible solutions. The implementation of *throughout* only generates solutions if the occasion holds in all world databases in the interval with the same variable bindings.

## VI. FINDING BEHAVIOR FLAWS

In order to improve plan parameters based on projection, we need to evaluate the time line. The analysis of the time lines finds flaws of various severity. Some flaws are critical for the overall outcome of a plan, for instance occluded objects. Critical flaws indicate that the plan will probably fail. Other flaws such as the distance to drive between different actions, the overall execution time of a plan and blocking objects are only a measure of the quality of solutions and indicate potential problems. In this paper, we will show the definitions of the flaws “*object occluded*” and “*object blocking*”.

The object occluded flaw basically means that an object that is to be perceived later is occluded because the put-down location of another object is in front of it. In the simplest case, an error indicating that an object could not be perceived is thrown during projection. The flaw definition does not require the perception plan to fail though. We first need to find out which object was to be perceived and if it was occluded by an object we put down previously. We define the flaw as follows:

*Task(?task<sub>1</sub>) ∧ Task(?task<sub>2</sub>)*  
*∧ TaskGoal(?task<sub>1</sub>, Perceive(?o<sub>1</sub>))*  
*∧ TaskGoal(?task<sub>2</sub>, Achieve(ObjectPlacedAt(?o<sub>2</sub>, ?l)))*  
*∧ TaskStatus(?task<sub>2</sub>, succeeded) ∧ TaskEnd(?task<sub>2</sub>, ?t)*  
*∧ Holds(Occluding(?w, ?o<sub>2</sub>, ?o<sub>1</sub>), at(?t))*

Please note that the variable *?w* is unbound in *Holds*. That means that the default database is used which is implicitly bound by *Holds* to match the data base as it was at time  $t$ . First, we query the task tree that was generated by plan execution for a failed perception task and a put-down task. Then we assert that the put-down was successful and we

bind the time when it finished to the variable  $?t$ . Finally, we check if the object that was put down is occluding the object that was to be perceived.

We define blocking objects as objects that the robot is in collision with while grasping an object (excluding the grasped object itself). While this flaw not necessarily leads to execution errors on the actual robot where grasp planning and motion planning is used, it is still an indicator for potential problems and performance penalties since a motion planner will need more time to find a valid plan when trajectories need to be complex. We define the flaw to find all objects that are in collision with the robot during the execution of a pick-up plan:

$$\begin{aligned} & \text{Task}(?tsk) \wedge \text{TaskGoal}(?tsk, \text{Achieve}(\text{ObjectInHand}(?o))) \\ & \wedge \text{TaskStart}(?tsk, t_s) \wedge \text{TaskEnd}(?tsk, t_e) \\ & \wedge \text{Holds}(\text{Blocking}(?w, PR2, ?o, ?b), \text{during}(t_s, t_e)) \end{aligned}$$

First we find all pick-up plans and get their start and end times. Then we use the *Holds* predicate to find blocking objects during the pick-up action.

## VII. CONCLUSIONS AND FUTURE WORK

In this paper, we presented a system for projecting high-level robot plans for pick-and-place actions in a human household. The system is based on the CRAM Plan Language execution environment and a Prolog reasoning engine that integrates a geometrically accurate world database and implements predicates for reasoning about stability, visibility and reachability using the Bullet physics engine, OpenGL rendering and inverse kinematics calculation. Projection is implemented by assertions and retractions in the world database and the generation of a time line containing events. The system can predict potential problems in plan execution such as unwanted occlusions or objects hindering future actions and is used to generate plan parameters such as the locations for placing objects or for placing the robot when performing an action.

The use of projected time lines is not restricted to finding plan parameters such as locations at run time. We will implement a transformational planner that will be able to predict and fix behavior flaws by applying structural transformations to a plan.

**Acknowledgements:** This work is supported in part by the EU FP7 Projects *RoboHow* (grant number 288533) and *SAPHARI* (grant number 287513) and within the DFG excellence initiative research cluster *Cognition for Technical Systems – CoTeSys*, see also [www.cotesys.org](http://www.cotesys.org). We thank Joris DeSchutter's group from KU Leuven (the authors of iTaSC) for insightful discussions.

## REFERENCES

- [1] D. McDermott, "Robot Planning," *AI Magazine*, vol. 13, no. 2, pp. 55–79, 1992.
- [2] M. Ghallab, A. Howe, C. Knoblock, D. McDermott, A. Ram, M. Veloso, D. Weld, and D. Wilkins, "PDDL—the planning domain definition language," *AIPS-98 planning committee*, 1998.
- [3] M. Fox and D. Long, "PDDL2.1: An extension of PDDL for expressing temporal planning domains," *Journal of Artificial Intelligence Research*, vol. 20, pp. 61–124, 2003.
- [4] D. Smith, Ed., *Special Issue on the 3rd International Planning Competition*, ser. Journal of Artificial Intelligence Research, vol. 20, 2003.
- [5] J.-C. Latombe, *Robot Motion Planning*. Boston, MA: Kluwer Academic Publishers, 1991.
- [6] S. M. LaValle, *Planning Algorithms*. Cambridge University Press, 2006.
- [7] M. Stilman and J. Kuffner, "Navigation among movable obstacles: Real-time reasoning in complex environments," in *Proceedings of the 2004 IEEE International Conference on Humanoid Robotics (Humanoids)*, vol. 1, December 2004, pp. 322–341.
- [8] S. Cambon, F. Gravot, and R. Alami, "asymov: Towards more realistic robot plans," in *International Conference on Automated Planning and Scheduling (ICAPS 2004)*, 2004.
- [9] F. Stulp, A. Fedrizzi, and M. Beetz, "Action-related place-based mobile manipulation," in *Proceedings of the International Conference on Intelligent Robots and Systems (IROS)*, 2009, pp. 3115–3120.
- [10] L. Mösenlechner and M. Beetz, "Using physics- and sensor-based simulation for high-fidelity temporal projection of realistic robot behavior," in *19th International Conference on Automated Planning and Scheduling (ICAPS'09)*, 2009.
- [11] L. Kunze, M. E. Dolha, E. Guzman, and M. Beetz, "Simulation-based temporal projection of everyday robot object manipulation," in *Proc. of the 10th Int. Conf. on Autonomous Agents and Multiagent Systems (AAMAS 2011)*, Yolum, Tumer, Stone, and Sonenberg, Eds. Taipei, Taiwan: IFAAMAS, May, 2–6 2011.
- [12] M. Beetz, L. Mösenlechner, and M. Tenorth, "CRAM – A Cognitive Robot Abstract Machine for Everyday Manipulation in Human Environments," in *IEEE/RSJ International Conference on Intelligent Robots and Systems*, Taipei, Taiwan, October 18–22 2010, pp. 1012–1017.
- [13] L. Mösenlechner and M. Beetz, "Parameterizing Actions to have the Appropriate Effects," in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, San Francisco, CA, USA, September 25–30 2011.
- [14] S. Zickler and M. Veloso, "Efficient physics-based planning: sampling search via non-deterministic tactics and skills," in *AAMAS '09: Proceedings of The 8th International Conference on Autonomous Agents and Multiagent Systems*. Richland, SC: IFAAMAS, 2009, pp. 27–33.
- [15] C. Dornhege, M. Gissler, M. Teschner, and B. Nebel, "Integrating symbolic and geometric planning for mobile manipulation," in *IEEE International Workshop on Safety, Security and Rescue Robotics (SSRR)*, 2009. [Online]. Available: <http://www.informatik.uni-freiburg.de/~ki/papers/dornhege-et-al-ssr09.pdf>
- [16] L. P. Kaelbling and T. Lozano-Perez, "Unifying perception, estimation and action for mobile manipulation via belief space planning," in *IEEE Conference on Robotics and Automation (ICRA)*, 2012. [Online]. Available: <http://lis.csail.mit.edu/pubs/tlp/ICRA12.1803.FI.pdf>
- [17] P. Michel, C. Scheurer, J. Kuffner, N. Vahrenkamp, and R. Dillmann, "Planning for robust execution of humanoid motions using future perceptive capability," in *Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS'07)*, October 2007, pp. 3223–3228.
- [18] L. Marin, E. A. Sisbot, and R. Alami, "Geometric tools for perspective taking for human-robot interaction," in *Mexican International Conference on Artificial Intelligence (MICAI 2008)*, 2008.
- [19] F. Zacharias, C. Borst, and G. Hirzinger, "Capturing robot workspace structure: representing robot capabilities," in *Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2007, pp. 3229–3236. [Online]. Available: <http://infoscience.epfl.ch/record/109336>
- [20] S. Hanks, "Practical temporal projection," in *Proc. of AAAI-90*, 1990, pp. 158–163.
- [21] D. McDermott, "An algorithm for probabilistic, totally-ordered temporal projection," in *Spatial and Temporal Reasoning*, O. Stock, Ed. Dordrecht: Kluwer Academic Publishers, 1997.
- [22] D. Pangercic, M. Tenorth, B. Pitzer, and M. Beetz, "Semantic object maps for robotic housework - representation, acquisition and use," in *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, Vilamoura, Portugal, October, 7–12 2012.
- [23] L. Mösenlechner, N. Demmel, and M. Beetz, "Becoming Action-aware through Reasoning about Logged Plan Execution Traces," in *IEEE/RSJ International Conference on Intelligent Robots and Systems*, Taipei, Taiwan, October 18–22 2010, pp. 2231–2236.

# Beyond Pick-and-Place — Naive Physics Reasoning for Handling Fluids in Everyday Manipulation

Lars Kunze, Reinhard Klapfer  
Intelligent Autonomous Systems  
Technische Universität München  
kunzel@cs.tum.edu

Michael Beetz  
Institute for Artificial Intelligence & TZI\*  
University of Bremen  
beetz@cs.uni-bremen.de

**Abstract**—Personal robot assistants that are to accomplish an open-ended set of everyday manipulation tasks like making pancakes are required to understand the physical effects of their own actions. In particular challenging are tasks that involve the handling of fluids.

In this paper, we investigate how robots can infer the consequences of their parameterized manipulation actions (here particularly pouring and mixing) in order to make competent and failure-aware decisions during their course of action. The proposed system allows robots to determine action parameters that lead to the desired effects by asking queries using a first-order language. The queries are answered based on interval-based first-order representations, called timelines, and learned decision trees which are grounded in detailed physics-based simulations of parameterized robot control programs.

## I. INTRODUCTION

Robotic agents acting in human environments have to handle an open environment with many different kinds of objects and stuff. If we do not want to hardcode how each kind of object should be handled we have to equip the robots with knowledge about these objects.

Consider, for example, a robotic assistant that is to help in a household. Such a robot has to handle fluids in various ways. It has to carry cups and bottles containing fluids, pour fluids from one container into another one, mix fluids, wipe them off the table, and so on. To perform these actions competently the robot has to know how to hold containers depending on their shapes, depending on whether they are open or closed, depending on the viscosity of the fluids and many other factors.

An important research question is in which form this knowledge should be encoded in robot control programs.

In Artificial Intelligence researchers proposed to state such knowledge explicitly and assert it as facts and axioms to a knowledge base [7], [8].

Generally, knowledge about the behavior of fluids is continuous in nature. But discrete knowledge about objects such as containers, their types and their openings is important, too. The latter is even important for seemingly simple pick-and-place tasks, where the size of a container opening can make a big difference for the execution.

Representing knowledge about fluids and their behavior dependent on the type of containers, their shape, etc. explicitly within a knowledge base seems to be an infeasible approach. Hence, in this work, we rather suggest an implicit

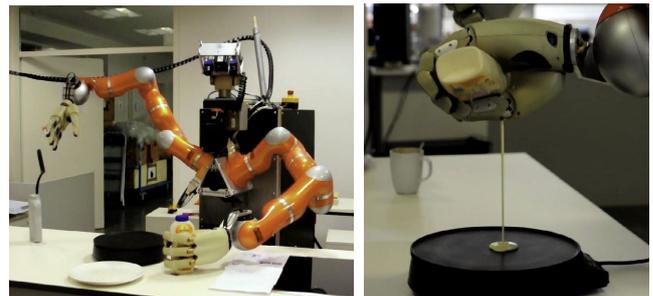


Fig. 1. The robot Rosie grasps a bottle of pancake mix and pours it onto the pancake maker.

representation of this body of knowledge. We only represent the core knowledge about fluids explicitly, namely the knowledge about its physical behavior. We then employ detailed physical simulations to derive the knowledge dependent on a particular context. Basically, the robot mentally envisions the outcome of its manipulation actions whereby the effects of its actions simply emerge from the laws of physics.

Humans are able to reason about these physical processes and adapt their behavior intuitively based on experiences, common sense and mental simulations, according to the simulation theory of cognition [9]. Understanding everyday physical phenomena, that is representing and reasoning about them, is an endeavor in the field of Artificial Intelligence which dates at least back to the work of Hayes [7]. More recently, there has been work on physical reasoning problems like “Cracking an egg” ([15]) which is listed on the common sense problem page<sup>1</sup>. In [2], Davis presents a formal solution to the problem of pouring liquids and in his work on the representation of matter [3], he investigated the advantages and disadvantages of various representations including those for liquids. In [4], he claims that it is tempting to use simulations for spatial and physical reasoning problems. But he also argues that simulations are not suitable for the interpretation of natural language texts because many entities in texts are highly underspecified. However, in the context of robotics entities in the environment can often be sufficiently recognized by sensors and represented using internal models. Therefore we believe that if we equip robots with Naive Physics, that is knowledge about the characteristics of

\*The Centre for Computing Technologies (TZI).

<sup>1</sup>[http://www.common-sense-reasoning.org/problem\\_page.html](http://www.common-sense-reasoning.org/problem_page.html)

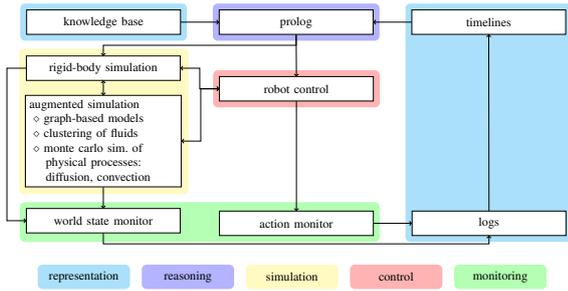


Fig. 2. The simulation-based temporal projection framework.

physical processes, objects, and substances, we will enable them to make informed decisions in context of planning, diagnosis and learning.

In this work, we build on the concept of simulation-based temporal projections as proposed in [13], [12]. Everyday robot manipulation tasks are simulated with varying parameterizations, world states of the simulation and states of the robot control programs are monitored and logged. The resulting logs are translated into a first-order representations, called timelines. These timelines are then used to answer logical queries on the resulting data structures in order to understand the physical effects of the robot’s actions. The main contribution of this work is the design and implementation of data structures and algorithms for representing and reasoning about fluids within this framework.

This paper is structured as follows: In Section II we give an overview of the simulation-based temporal projection framework and briefly explain its components. Section III describes how fluids are represented and simulated. The experimental results are presented in Section IV. Finally, we discuss the present work and conclude in Section V.

## II. SIMULATION-BASED TEMPORAL PROJECTION

This section gives a brief overview of the simulation-based temporal projection framework as introduced in [13], [12].

The framework, its components as well as the data flow are depicted in Figure 2. The framework is based on state-of-the-art technologies such as ROS<sup>2</sup>, the Gazebo simulator<sup>3</sup> and the point cloud library PCL<sup>4</sup>. A manipulation scenario can be specified in a knowledge base using predicates of a first-order language. Based on this description, Prolog instantiates an adequate physical simulation. Within the simulation a robot can freely navigate and interact with objects. The behavior of the robot is specified by a robot control program. In the simulator we represent, e.g., a pancake mix as particles using the data structures of Gazebo. Given that we are particularly interested in analyzing the behavior of liquids we group the simulated particles by an Euclidean clustering technique. Having obtained information of clusters makes it possible to reason about the fusion or division of volumes or chunks

of liquids. The clustering is realized as node located at an augmented simulation layer. As Gazebo uses ODE<sup>5</sup> which is only capable of dealing with rigid bodies the simulation of liquids is only an approximation. Therefore we use the information about the clusters to initialize a more accurate simulation of liquids by considering physical aspects such as molecular motion due to diffusion and convection. The robot’s actions, its interactions with the objects, the state of the liquid, the clusters and the state of the environment (world) are crucial information for the reasoning framework. For example, the world state comprises information about the position, orientation, linear and angular velocities, and the bounding box of an object at a point in time and is denoted as follows:

*World state* :  $\langle time, obj, pos, orient, lin\_vel, ang\_vel, bbox \rangle$ .

Additionally, we also monitor contact events:

*Contact state* :  $\langle time, o1, o2, num, force, torque, normal \rangle$ ,

whereby we observe the number of contact points as well as the forces, torques and normals between them. We have implemented monitoring routines as Gazebo controllers that keep track of the dynamics of objects and write this information to log files.

These logs are then translated into interval-based first-order representations. We access and evaluate the data structures from Prolog using the following predicate

$$SimulatorValue(\overbrace{position(o, pos)}^{Function}, \overbrace{t}^{Time\ point}, \overbrace{tl}^{Timeline}),$$

whereby different functions are available for accessing the time-stamped information of the world and contact states.

In order to define more high-level predicates we use concepts similar to those in the Event Calculus [11]. The notation is based on two concepts, namely fluents and events. Fluents are conditions that change over time, e.g., a mug contains a pancake mix:  $contains(mug, mix)$ . Events (or actions) are temporal entities that have effects and occur at specific points in time, e.g., consider the action of pouring the mix from the mug onto the pancake maker:  $occurs(pour(mug, pan))$ . Logical statements about both fluents and events are expressed by using the predicate  $Holds(f, t, tl)$  where  $f$  denotes a fluent or event,  $t$  simply denotes a point in time, and  $tl$  a timeline. Using the similar predicate  $Holds_{tt}$  we can query for a time interval throughout the fluent holds. Logical queries to the framework are basically answered through Prolog’s backtracking mechanism over a set of timelines.

Having given an overview of the overall system, the next section will lead us to the representation and simulation of fluids within this framework.

## III. REPRESENTATION AND SIMULATION OF FLUIDS

Simulating liquids is of great interest in physics and chemistry. As some processes occur very fast, events might

<sup>2</sup><http://www.ros.org>

<sup>3</sup><http://www.gazebosim.org>

<sup>4</sup><http://www.pointclouds.org>

<sup>5</sup><http://www.ode.org>

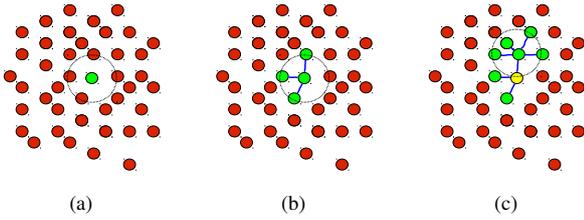


Fig. 3. Generating a deformable pancake model from liquid particles. Illustration of the algorithm’s procedure: (a) Radial search from the seed Point. (b) Creation of hinge joints to the neighbors. (c) Radial search and creation of joints in a recursive step.

not be observable in all its details in reality. The purpose of simulating liquids in our work is to observe the impact of the robot’s action with respect to the liquid’s behavior, which is of great importance when, e.g., pouring and mixing liquids. Different approaches have been incorporated to simulate liquids depending on the required level of accuracy needed. In this work we propose two complementary approaches for simulating liquids, (1) a graph-based model similar to [10] and (2) a Monte-Carlo simulation for modeling diffusion and convection [5]. Both do not simulate liquids in all their aspects but provide enough information for making logical inferences about qualitative phenomena.

#### A. Representing Fluids using Graph-based Models

The model for representing fluids was adapted from the work of Johnston [10]. Originally, it was designed to simulate a wide range of physical phenomena including diverse domains such as physical solids or liquids as hyper-graphs where each vertex and edge is annotated with a frame that is bound to a clock and linked to update rules that respond to discrete-time variants of Newton’s laws of mechanics.

Our pancake mix model can be in two states: first, the mix is liquid, and second, the mix becomes a deformable pancake after cooking. In the simulation we use a graph-based model for representing the mix and the pancake. The vertices of the graph are particles where each particle is defined by a round shape with an associated diameter, a mass and a visual appearance model. The benefit of this model is that it is realized as graph with no connection between the vertices whenever the state is liquid. This means that the individual particles could move freely to some extent. This was useful for performing the pouring task. Due to the fact of the particles not being connected with joints, the simulated liquid can be poured over the pancake maker where it dispenses due to its round shape. A controller was attached to the spheres that applies small forces to the particles in order to simulate the viscosity of the pancake mix. Currently, we do not consider heat as the trigger of transforming the liquid to a solid pancake but simply assume the event to occur after constant time. We identified all particles on the pancake maker and created the pancake based on a graph traversal algorithm starting at the cluster center (Figure 3).

#### B. Clustering of Fluid Particles

The basic idea of applying clustering methods is as follows. Let us, for example, assume that someone pours some pancake mix onto a pancake maker as illustrated in Figure 4. After the pouring action some particles reside in the container, some are spilled onto the table, and some others are on the pancake maker which will eventually form the pancake. If we want to address the particles in these three locations it perfectly makes sense to group them in chunks (clusters). This reflects also how humans address fluids like milk or sugar in natural language, e.g., there is some milk spilled onto the table. Therefore the behavior and the contact information of clusters of particles in everyday manipulation tasks are of particular interest. We decided to use a Euclidean Clustering strategy for computing the groups of particles as shown in Algorithm 1.

---

#### Algorithm 1 Euclidean clustering of particles.

---

- 1) Set up an empty list of clusters  $Clst$
  - 2) For every particle  $p_i \in P$  do
    - Add  $p_i$  to the current cluster  $C$
    - For every point  $p_j \in C$  do
      - Find particle  $p_k$  using a radial search around particle  $p_j$
      - For each particle  $p_k$  add it to  $C$  if not processed, yet
      - Terminate if all  $p_j \in C$  have been processed
    - Add  $C$  to the list of clusters  $Clst$  and reset  $C$  to an empty list
  - 3) The algorithm terminates if all particles have been processed and are part of a cluster  $c_i \in Clst$
- 

Instead of looking at the individual particles when interpreting the outcome of a manipulation scenario we look at clusters of particles. For every cluster we compute information such as mean, covariance, size (number of particles), and its bounding box. Since we have full knowledge about every particle and its belonging to a cluster, we can keep track of it, i.e., if its pose or extension change over time. However, whenever new particles become part of or are separated from a cluster we assign a new ID to it. That is, clusters of particles have only a limited time during which they exist. Hence, we can recognize which actions cause changes to clusters.

#### C. Monte Carlo Simulation of Fluids

Deformable bodies are seen as a big challenge in simulation and usually require a lot of computational power [1]. The physical simulation approach [5] uses a Monte-Carlo process to simulate diffusion of liquids. Molecular movement is either provoked from heat or from a difference in potential.

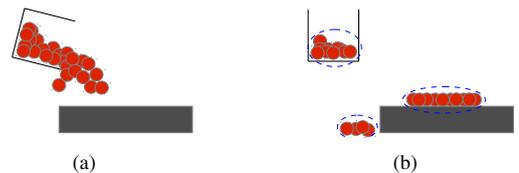


Fig. 4. Basic idea of the clustering approach: during simulation we identify clusters of particles. For example, after pouring, one cluster resides still in the mug, a second is on the pancake maker and a third is spilled onto the table. We are able to extract information including contacts, position, extension, and size of the individual clusters.

The rate of change depends on the diffusion coefficient and its respective change. This is a well known concept in physics described by equation 1 and denoted as the *macroscopic diffusion* equation or *Fick's second law* of diffusion. This differential equation takes into consideration a change of concentration over time.

$$\frac{\partial C}{\partial t} = D \cdot \frac{\partial^2 C}{\partial x^2} \quad (1)$$

It can be shown [5] that *Random Walk* gives one particular solution for the above partial differential equation. Motivated by this idea we applied Algorithm 2 proposed by Frenkel et al. to simulate this physical effect. The algorithm follows the Metropolis scheme and uses a probability function to decide if a particle is going to be displaced or not. The Leonard-

---

#### Algorithm 2 Metropolis scheme.

---

- 1) Select a particle  $r$  at random and calculate its energy potential  $U(r^N)$
- 2) Give the particle a random displacement,  $r' = r + \Delta$
- 3) Calculate the new energy potential  $U(r'^N)$
- 3) Accept the move from state  $r^N$  to  $r'^N$  with probability

$$\text{acc}(r^N \mapsto r'^N) = \min\left(1, \exp\left(-\beta \left[U(r'^N) - U(r^N)\right]\right)\right)$$


---

Jones Potential Function (Equation 2) was used to model the interaction among the particles in the liquid, i.e., to model the particles' behavior according to the concentration of particles in their neighborhood. The parameters  $\sigma$  and  $\epsilon$  are used to shape the function and  $r$  is the distance to neighboring particles.

$$U(r) = 4\epsilon \left[ \left(\frac{\sigma}{r}\right)^{12} - \left(\frac{\sigma}{r}\right)^6 \right] \quad (2)$$

Stirring a material is another type of mass transfer called convection. Convection is the movement of mass due to forced fluid movement. Convective mass transfer is a faster mass transfer than diffusion and happens when stirring is involved. The faster the fluid moves, the more mass transfer and therefore the less time it takes to mix the ingredients together [6]. We simulated this physical property by simply introducing an impulse in stirring direction to the particles in the point cloud that are in reach of the cooking spoon. In this way, we could achieve with this simple model the behavior of molecular motion due to forced fluid movement.

#### D. Measuring the Homogeneity of Mixed Fluids

Particular interest is the homogeneity of the liquid when stirred was involved in the conducted experiments. It was decided to use the local density of the particles represented as point cloud as a measure of divergence, while using the assumption that the inverse of this is a measure of homogeneity. This distance measure [14] is known as the Jensen-Shannon divergence and used widely in information theory. The Jensen-Shannon divergence is defined as:

$$JS(P, Q) = \frac{1}{2} S\left(P, \frac{P+Q}{2}\right) + \frac{1}{2} S\left(Q, \frac{P+Q}{2}\right) \quad (3)$$

where  $S(P, Q)$  is the Kullback divergence shown in equation 4, and  $P$  and  $Q$  two probability distributions defined over a discrete random variable  $x$ .

$$S(P, Q) = \sum_x P(x) \log\left(\frac{P(x)}{Q(x)}\right) \quad (4)$$

We propose the division of the point cloud in a three-dimensional grid (Figure 5). Each cell of the grid represents a discrete probability distribution  $x$  defined on the mixed probabilities of the two classes  $P$  and  $Q$ , that could be computed as the relative frequency.

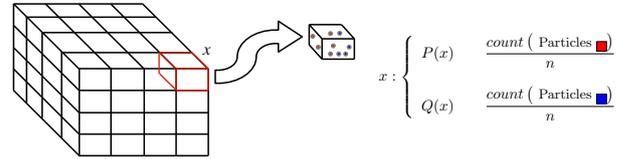
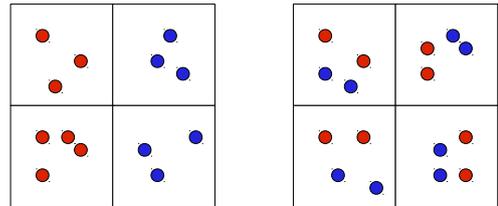


Fig. 5. Density grid used for discretization and local density estimation.

The following example emphasizes the usage of this distance function related to the homogeneity of a liquid which consists of two classes of particles. If we assume a perfect separation of the two classes as shown in Figure 6, we would expect a high divergence and a low homogeneity as we define the homogeneity as its inverse.



(a) Maximum divergence, minimal homogeneity. (b) Minimal divergence, maximal homogeneity.

Fig. 6. A simple 2D density grid for a two class problem.

## IV. EXPERIMENTAL RESULTS

In this section we are going to highlight the results of both experiments, namely mixing ingredients in a bowl while measuring the level of homogeneity<sup>6</sup>, and second, pouring the mix onto a pancake maker and reasoning about the resulting size and whether some mix was spilled<sup>7</sup>.

#### A. Mixing Liquids — Analysis of Homogeneity

We used the Monte Carlo method previously described in Section III-C to simulate the physical effects when mixing liquids with different trajectories.

We selected the coefficients to represent two viscous liquids. Figure 7 and Figure 8 show the course of homogeneity when the robot stirs the liquids using (1) an elliptic trajectory, (2) a spiral trajectory, (3) a lineal trajectory, and (4) no trajectory (without stirring). As we expect, the ingredients do not mix very well when the robot does not stir the liquids. Hence, the result of the experiment confirms our hypothesis: Stirring increases the homogeneity of mixed liquids.

Furthermore, the result shows that with an elliptic trajectory the best result could be achieved. Although Figure 7

<sup>6</sup>Video (mixing): <http://www.youtube.com/watch?v=ccHXmkKT8CE#>

<sup>7</sup>Video (pouring): <http://www.youtube.com/watch?v=tzQk7SSPRaY>

visualizes continuous data of the homogeneity, we are mainly interested in qualitative effect models. These qualitative models of *homogeneous*, *semi-homogeneous* and *inhomogeneous* regions are simply defined by thresholding the quantitative data. Given the knowledge of homogeneous, semi-homogeneous and inhomogeneous regions a robot could adapt the trajectory dynamically by applying techniques known from Reinforcement Learning.

### B. Pouring Fluids — Reasoning about Clusters

In this experiment we address the scenario of pouring some pancake mix located in a container onto a pancake maker: the robot grasps a mug containing pancake mix from the table, lifts it and pours the content onto a pancake maker (Figure 9). In this experiment we used the resulting timelines to analyze the qualitative outcome of the executed action. The parameterization of the task included the gripper position, the pouring angle and the pouring time. We also looked at different container types and fill levels. The task was considered to be successful if no pancake mix has been spilled, i.e. the liquid resides on the pancake maker or in the container and not on other objects such the kitchen table after the pouring action ends. We used the resulting clusters and their corresponding contact and spatial information to examine the outcome. Figure 10 shows exemplarily how clusters of pancake mix are spatially related to other objects before, during and after the pouring action.

The following Prolog expression shows how information about clusters can be retrieved from timelines (*TL*):

```
?- holds_tt(occurs(pour(Params)),I,TL), [_ ,End] = I,
    partOf(X,pancake_mix), holds(on(X,pmaker),Time,TL),
    after(Time,End),
    simulator_value(size(X,Size),Time,TL).
```

where  $X$  denotes a cluster of pancake mix in contact with a pancake maker after a pouring action has been carried out.

We used logical queries such as above to extract data for learning decision trees in order to classify pancake sizes and pouring angles. Although decision trees can also be learned on continuous data, we mapped the numerical data to qualitative concepts such as, for example, *Small*, *Medium* and *Large* pancakes. Thereby, the resulting qualitative models are intuitively interpretable by humans. For learning we used

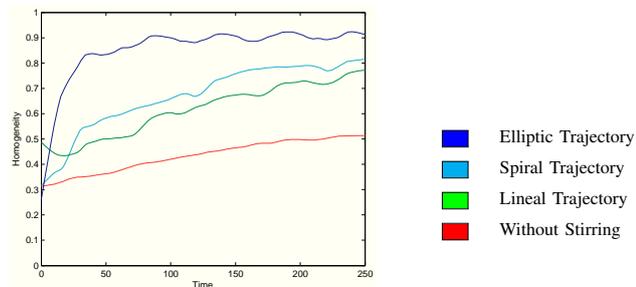


Fig. 7. Homogeneity over time of different stirring trajectories. The graph shows the change of homogeneity on the vertical axis for different trajectories in direct comparison with the result of scenario of not stirring over time.

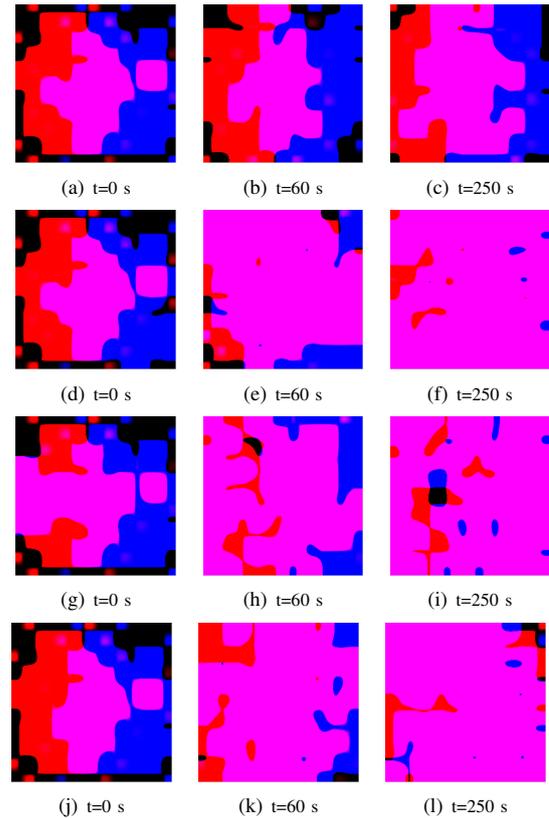


Fig. 8. The color coded images show the spatial distribution of homogeneity of two liquids as a 2D slice of the density grid depicted in Figure 5. Black stands for uncovered regions, red and blue for inhomogeneous liquids of corresponding classes and purple homogeneous regions. Stirring trajectories: (a-c) without stirring, (d-f) elliptic, (g-i) lineal, (j-l) spiral.

Weka’s *J48* algorithm in its default parameterization [16].

The resulting decision tree for pancake sizes is visualized in Figure 11. Overall, the learned model achieves an accuracy of 92.41%. The most decisive attribute is the fill level (particles). If there are only a few particles available, the robot can only make small pancakes. In case of many particles, the size of the pancake depends first on the type of container, and second on the tilting angle. Only if the container is a bottle with a small opening and the angle is high, the robot can make pancakes of different sizes by varying the time.

In a second experiment we have learned an action model for predicting the pouring angle. The learned decision tree, depicted in Figure 12, achieves an accuracy of 64.35%. Given a desired size of a pancake and a context determined by the container’s type and its fill level, the robot can infer an appropriate angle.

## V. CONCLUSIONS

The present work can be considered as interdisciplinary research of two fields: Robotics and Artificial Intelligence.

With our approach we enable robots to reason about the consequences of their own actions. We equip them with the capability of making appropriate decisions about their parameterizations throughout their activity using well-

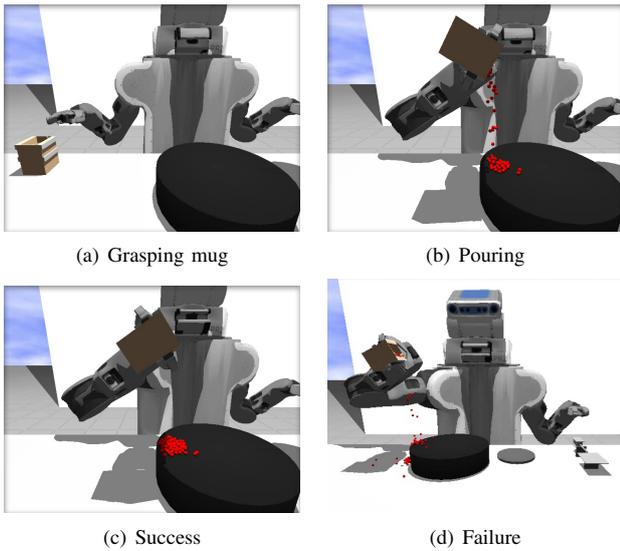


Fig. 9. PR2 pours mix onto pancake maker.

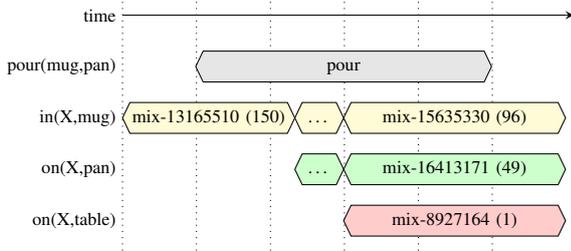


Fig. 10. Visualization of the main clusters of particles before, during and after the pouring action. The pancake mix was represented by 150 particles. The number of particles of a cluster is shown in parenthesis. During the simulation there were more than 20 clusters generated on this timeline.

established methods of AI and detailed physical simulations.

To this end, we have developed a system that simulates robot manipulation tasks, monitors relevant states and actions, and translates this information into first-order representations, called timelines. Then, we use the logic programming language Prolog to answer queries based on the data structures of the temporal projections.

The main contribution of this work is the extension of the framework with respect to data structures and algorithms for representing and simulating fluids. We conducted two experiments within the framework: mixing and pouring liquids. The resulting timelines of the experiments were evaluated with different performance criteria, e.g., the homogeneity of the mix, and used for learning compact decision trees for predicting the size of a pancake and a pouring angle.

#### Impact on Robotics

The concept of simulation-based temporal-projections can be used to equip robots with the ability to understand physical phenomena. In particular, robots can employ the qualitative information about physical aspects of their manipulation tasks to answer questions in the following contexts:

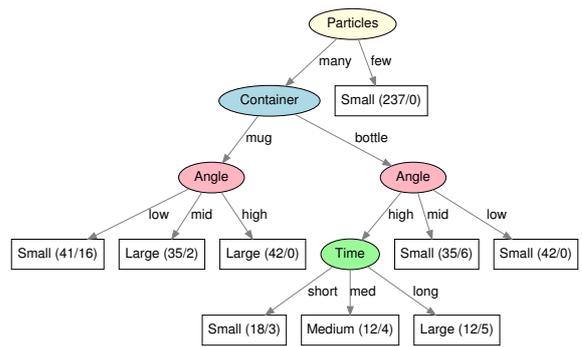


Fig. 11. Decision tree for predicting the size of a pancake. The size is discretized in three classes, namely *Small*, *Medium*, and *Large*. The tree is learned from 474 instances and classifies 438 instances correctly (92.41%) and 36 incorrectly (7.59%).

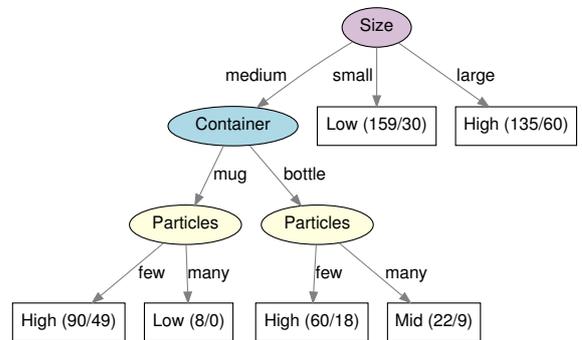


Fig. 12. Decision tree for selecting an angle. The angle is discretized in three classes, namely *Low*, *Mid*, and *High*. The tree is learned from 474 instances and classifies 305 instances correctly (64.35%) and 169 incorrectly (35.65%).

**Monitoring** What is the expected outcome of an action?

**Planning** Which action will lead to the intended goal?

**Diagnosis** What has caused something to happen?

**Question Answering** Why has an action being performed?

**Reinforcement Learning** How to explore the parameter space of an action effectively?

This incomplete list of contexts and queries illustrates where the developed framework and learned models can be employed in order to adjust the behavior and to improve the overall performance of robots. Hence, we believe that the underlying idea and the developed methods of this work can have a broad impact in field of robotics.

#### Acknowledgments

This work has been supported by the EU FP7 Projects *RoboHow* (grant number 288533) and *SAPHARI* (grant number 287513) and the DFG excellence initiative research cluster *Cognition for Technical Systems – CoTeSys*.

#### REFERENCES

- [1] J. Brown, S. Sorkin, C. Bruyns, J.-C. Latombe, K. Montgomery, and M. Stephanides. Real-time simulation of deformable objects: Tools and application. In *IN COMP. ANIMATION*, 2001.
- [2] E. Davis. Pouring liquids: A study in commonsense physical reasoning. *Artif. Intell.*, 172(12-13):1540–1578, Aug. 2008.

- [3] E. Davis. Ontologies and representations of matter. In M. Fox and D. Poole, editors, *Proceedings of the Twenty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2010, Atlanta, Georgia, USA, July 11-15, 2010*. AAAI Press, 2010.
- [4] E. Davis. Qualitative spatial reasoning in interpreting text and narrative. *Spatial Cognition and Computation*, 2012. Forthcoming.
- [5] D. Frenkel and B. Smit. *Understanding Molecular Simulation, Second Edition: From Algorithms to Applications (Computational Science)*. Academic Press, 2 edition, Nov. 2001.
- [6] H. Gould, J. Tobochnik, and C. Wolfgang. *An Introduction to Computer Simulation Methods: Applications to Physical Systems (3rd Edition)*. Addison-Wesley Longman Publishing Co., Inc., Boston, MA, USA, 2005.
- [7] P. Hayes. The naive physics manifesto. In D. Michie, editor, *Expert Systems in the Micro Electronic Age*, pages 242–270. Edinburgh University Press, 1979.
- [8] P. Hayes. Naive physics i: Ontology for liquids. In J. R. Hobbs and R. C. Moore, editors, *Formal Theories of the Commonsense World*, pages 71–107. Ablex, Norwood, NJ, 1985.
- [9] G. Hesslow. The current status of the simulation theory of cognition. *Brain Research Reviews*, 1428:71–79, 2012.
- [10] B. Johnston and M.-A. Williams. Comirit: Commonsense reasoning by integrating simulation and logic. In *Proceedings of the 2008 conference on Artificial General Intelligence 2008: Proceedings of the First AGI Conference*, pages 200–211, Amsterdam, The Netherlands, The Netherlands, 2008. IOS Press.
- [11] R. Kowalski and M. Sergot. A logic-based calculus of events. *New generation computing*, 4(1):67–95, 1986.
- [12] L. Kunze, M. E. Dolha, and M. Beetz. Logic Programming with Simulation-based Temporal Projection for Everyday Robot Object Manipulation. In *2011 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, San Francisco, CA, USA, September, 25–30 2011. Best Student Paper Finalist.
- [13] L. Kunze, M. E. Dolha, E. Guzman, and M. Beetz. Simulation-based temporal projection of everyday robot object manipulation. In Yolum, Tumer, Stone, and Sonenberg, editors, *Proc. of the 10th Int. Conf. on Autonomous Agents and Multiagent Systems (AAMAS 2011)*, Taipei, Taiwan, May, 2–6 2011. IFAAMAS.
- [14] A. P. Majtey, P. W. Lamberti, and D. P. Prato. Jensen-shannon divergence as a measure of distinguishability between mixed quantum states. *Phys. Rev. A*, 72:052310, Nov 2005.
- [15] L. Morgenstern. Mid-Sized Axiomatizations of Commonsense Problems: A Case Study in Egg Cracking. *Studia Logica*, 67(3):333–384, 2001.
- [16] I. H. Witten and E. Frank. *Data Mining: Practical Machine Learning Tools and Techniques*. Morgan Kaufmann, San Francisco, second edition, 2005.

# Bibliography

[Klapfer et al., 2012] Klapfer, R., Kunze, L., and Beetz, M. (2012). Pouring and mixing liquids—understanding the physical effects of everyday robot manipulation actions. *Human Reasoning and Automated Deduction*.

[Mösenlechner and Beetz, 2013] Mösenlechner, L. and Beetz, M. (2013). Fast temporal projection using accurate physics-based geometric reasoning. In *IEEE International Conference on Robotics and Automation (ICRA)*, Karlsruhe, Germany.