

Toward Human-Like Robot Learning

Sergei Nirenburg¹, Marjorie McShane¹, Stephen Beale¹, Peter Wood¹, Brian Scassellati²,
Olivier Magnin² and Alessandro Roncone²

¹ Language-Endowed Intelligent Agents Lab, Rensselaer Polytechnic Institute

² Social Robotics Lab, Yale University

Abstract. We present an implemented robotic system that learns elements of its semantic and episodic memory through language interaction with people. This human-like learning can happen because the robot can extract, represent and reason over the meaning of the user's natural language utterances. The application domain is collaborative assembly of flatpack furniture. This work facilitates a bi-directional grounding of implicit robotic skills in explicit ontological and episodic knowledge, and of ontological symbols in the real-world actions by the robot. In so doing, this work provides an example of successful integration of robotic and cognitive architectures.

Keywords: Learning based on language understanding, Integration of robotic and cognitive architectures, Memory management in artificial intelligent agents.

Setting the Stage

The ability of today's well-known artificial intelligent agents (Siri, Alexa, IBM's Watson, etc.) to react to vocal and written language communication relies ultimately on a variety of sophisticated methods of manipulating textual strings that are largely semantically uninterpreted (though a modicum of linguistic knowledge is used in their operation [39]). This methodology offers a solid breadth of coverage at a societally acceptable level of quality for select application areas. Still, as an approach to modeling human-level capabilities, this methodology has well-known limitations that are due to the dearth of knowledge resources and processing engines needed for extracting and representing the various kinds of meaning conveyed through natural language (NL) texts. To give just one example, knowledge-lean techniques find it difficult to resolve reference, as highlighted by the Winograd schema challenge competitions [11]. This paper presents a proof-of-concept system for teaching a robot in a human-like manner, relying on a combination of its NL understanding and reasoning skills. As a side effect, such a system demonstrates a methodology for overcoming the abovementioned knowledge acquisition bottleneck.

Within the field of cognitive systems, several recent projects have been devoted to aspects of language understanding [1, 6, 7, 9, 12, 22, 25, 35, 37, 38], a response to the fact that the knowledge-lean paradigm currently prevalent in NLP does not address the needs of sophisticated agent systems. Two characteristics set out system apart from the others. First, it simultaneously addresses the challenges of a) learning-oriented language-based human-robotic interaction, b) symbol grounding, c) linguistic meaning

extraction, and d) the enhancement and management of the episodic, semantic and procedural memory of a robot/agent. Second, the language processing component of the system and its associated knowledge resources address a broader set of meaning-related language phenomena, described at a finer grain-size of analysis.

To implement language-based learning in a social robotics environment, we must address the co-dependence among three capabilities: language understanding, learning, and task-oriented functioning. Language understanding requires knowledge, while the learning achieved through language understanding automatically adds to that knowledge. More knowledge means better language understanding, resulting in an increasingly effective, human-like process of lifelong learning.

To support a robot's interpretation of the results of its perceptual inputs as well as its decision-making and action, we must model the "minds" of the communicating agents: the types of objects and events in their world; the instances thereof mentioned in current and past communications; the agents' knowledge about their tasks and responsibilities; their beliefs about other agents; as well as their inventories of desires, goals, decision biases, etc. [12, 20, 22, 30]. Hence the need for modeling agent memory. Most cognitive architectures distinguish three kinds of memory: semantic – roughly, memory of known types of entities in the world or domain; episodic – remembered instances of entities; and procedural – roughly, uninterpreted skill-related routines.

Learning How to Build a Chair

The system we describe is a social robot collaborating with a human user. The experimental domain is furniture assembly (e.g., [10]), widely accepted as useful for demonstrating human-robot collaboration on a joint activity. Roncone et al. [34] report on a Baxter robot supplied with high-level specifications, represented in the HTN formalism [5], of basic actions implementing chair-building tasks. Using a keyboard command or pressing a button, the user could trigger the execution of basic actions by triggering the operation of low-level task planners that the robot could directly execute. The robot could not reason about its actions, which were stored in its procedural memory as uninterpreted skills. The system described here integrates the robotic architecture of [34] with a cognitive architecture [16]. The integrated system's language understanding and reasoning capabilities make it possible for the robot to a) learn the semantics of its hitherto uninterpreted basic actions; b) learn the semantics of operations performed by the robot's human collaborator and conveyed to the robot in natural language; c) learn, name and reason about meaningful groupings and sequences of the actions in a) and b) above and organize them hierarchically; and d) integrate the results of learning with knowledge stored in its semantic and episodic memory and establish connections between these memory modules and its procedural memory.

The core prerequisite for human-like learning is the ability to automatically extract, represent and use the meaning of natural language texts – utterances, dialog turns, etc. This task is notoriously difficult: to approach human-level capabilities, intelligent agents must account for both propositional and discourse meaning; interpret both literal and non-literal (e.g., metaphorical or metonymical) meaning; offer a solution for reference resolution as well as implicature; and, particularly in informal genres, deal

with stops and starts, spurious repetitions, production errors, noisy communication channels and liberal (if unacknowledged) use of the least effort principle (e.g., [33]) by speakers and hearers. The language understanding module of our system, OntoSem [22], demonstrates progress on all of the above issues.

As its knowledge resources, OntoSem uses an ontology (world model) of some 9,000 concepts with on average of 16 properties each; a semantic lexicon for English covering about 25,000 lexical senses; and a frame-oriented formalism suitable for representing the semantics of robotic actions, natural language utterances and results of the robot's processing of other perceptual modalities (e.g., interoception, see [16]). The ontology constitutes the agent's semantic memory. It also maintains episodic memory, in the form of remembered instances of ontological concepts that it retained over its lifetime of processing perceptual input and carrying out reasoning.

In addition to the above static knowledge resources, OntoSem includes a variety of linguistic microtheories and associated processing algorithms covering phenomena such as lexical disambiguation and semantic dependency determination [13, 17]; multiword expressions [21]; sentence fragments [19]; reference and ellipsis [18, 40]; unexpected input [23]; mental model ascription [14]; speaker bias detection [20]; and processing non-literal language [31].

| | | |
|--------------------------|-------------------|-------------------------|
| SPEECH-ACT-1 | | |
| type | | command |
| scope | CHANGE-LOCATION-1 | |
| producer | | *speaker* |
| consumer | | ROBOT-0 |
| time | | time-0 ; time of speech |
| CHANGE-LOCATION-1 | | |
| agent | | ROBOT-0 |
| theme | SCREWDRIVER-1 | |
| effect | BESIDE | |
| | (AGENT.LOCATION | |
| | THEME.LOCATION) | |
| time | | > time-0 |
| token | | <i>fetch</i> |
| from-sense | | move-v2 |
| HUMAN-1 | | |
| agent-of | CHANGE-LOCATION-1 | |
| token | | <i>you</i> |
| from-sense | | you-n1 |
| SCREWDRIVER-1 | | |
| theme-of | CHANGE-LOCATION-1 | |
| token | | <i>screwdriver</i> |
| from-sense | | screwdriver-n1 |

Fig. 1. Meaning representation for the utterance *Now you will fetch a screwdriver* (simplified).

representations (MRs) of these utterances, generated by OntoSem. An example MR is presented in Fig. 1.

The Process.

At the beginning of the learning process, the robot can a) visually recognize parts of the future chair (the back, the seat, two types of dowels, and the tool (screwdriver) to be used in chair assembly and b) perform the following preprogrammed basic actions: GET(object) (e.g. a bracket or a screwdriver) from storage area to workspace; HOLD(object) to facilitate the human's actions and RELEASE(object) to the work surface.¹ While these actions are not primitive in the robot's procedural memory, they are conceptualized (and will be remembered in the semantic memory) as primitive event, with no constituent events listed in their description, as there is no need for the robot to reason about parts of these actions.

Input sequences for learning basic routines consist of the following types of elements: 1) commands to perform basic actions (to be executed by the robot's motor control module) and 2) natural language utterances by the user. The robot learns by reasoning on the basis of the meaning

¹ The robot also knows to CLEAR the workspace, but this is not used when assembling furniture.

In our task, the user teaches the robot three types of things: a) concept grounding: the connection between a basic action the robot knows how to perform and this action’s mental representation, constructed on the basis of relevant MRs and the robot’s stored knowledge; b) how sequences of basic actions can be combined to form representations of complex actions describing the robot’s work processes; and c) how these sequences (and their subsequences) can be encoded as concepts in the robot’s world model. Three different algorithms (modules) implement these three kinds of learning (see Fig. 2.). The algorithms are not applied in the order listed but rather are called whenever the input string licenses or requires them being triggered.

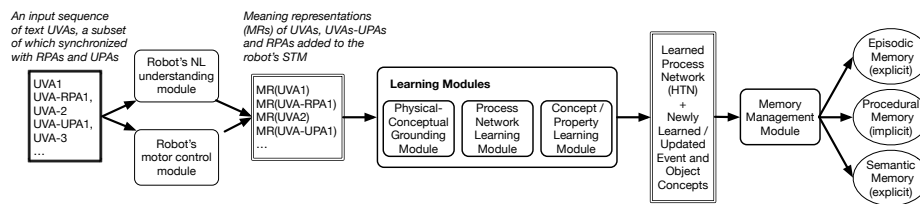


Fig. 2. The core learning process. Input is a sequence of user verbal actions (UVAs) which explain user physical actions (UPAs) and issue commands to the robot, thus verbalizing robot’s physical actions (RPAs), which facilitates grounding the former in the latter. UVAs are interpreted into uniform meaning representation and provide input to grounding, process network and concept/property learning modules (LMs). The memory management module (MMM) incorporates the results of learning into the episodic and semantic memories of the robot and mutually grounds RPAs in the robot’s procedural memory and corresponding concepts in its semantic memory.

At the beginning of the learning process, our robot channels Monsieur Jourdain from Molière’s *Le Bourgeois Gentilhomme* who was amazed to learn that all his life he was able to talk prose but didn’t know that he did so. Our robot can perform a number of basic actions but at the outset of the learning session does not know that it is performing these actions. As an illustration of concept grounding, consider an input subsequence consisting of the utterance of Fig.1 followed by a command for the robot to execute the basic action GET(screwdriver). OntoSem will generate the MR in Fig.1. When the robot executes GET(screwdriver), the physical-conceptual grounding learning module (LM) will link this procedure call with the representation of CHANGE-LOCATION-1 in the MR, thus linking the robotic and the cognitive architectures. This grounding is implemented as part of the concept/property LM of the robot (see Fig. 2) by adding the property PM-LINK (‘pm’ stands for “procedural memory”) to CHANGE-LOCATION-1, with the filler GET(screwdriver). The immediate purpose of this linking of the robotic and the cognitive architectures is to make the robot capable in its subsequent functioning to a) trigger basic actions autonomously on the basis of language input alone, without having the user issue the corresponding command and b) learn complex event sequences by just being told, without actually performing the actions comprising the complex event.

In parallel with grounding basic actions, the robot learns the legal sequences of actions that successfully complete the task at hand. This kind of learning is facilitated by the robot’s language processing capability in that it can learn by understanding the user’s utterances in their context. The robot organizes action sequences hierarchically

and makes sure that any non-terminal nodes in the resulting process network represent meaningful complex actions. If the robot does not have specifications for these complex actions in its stored knowledge, it learns new concepts for them, using the MRs obtained by processing the relevant user utterances.

Team tasks such as our chair assembly task typically involve joint actions by team members as well as individual actions by each of them. We treat joint tasks as complex tasks and require the system to decompose them into subtasks carried out by each of the team members. The individual tasks in our system include the robot's basic actions (labeled Robot Physical actions, RPAs) as well as actions performed by its user (User Physical Actions, UPAs), provided the latter are described by the user through an accompanying User Verbal Action, UVA. The RPAs and UPAs appear as terminal nodes in the process network being learned.² The robot's activity that includes all the kinds of learning it does as well as updating its memory structures comes under the rubric of Robot's Mental Action, RMS. Due to space constraints we cannot illustrate a complete process of assembling a chair (even the shortest version of the process numbers over 150 steps). So, we will present a small subset of this process – assembling the

| | |
|--|--|
| UVA1 | We will now build the right back leg |
| UVA2 | Get another foot bracket |
| RPA1 | GET(bracket-foot) |
| RPA2 | RELEASE(bracket-foot) |
| UVA3 | Get the right back bracket |
| RPA3 | GET(bracket-back-right) |
| RPA4 | RELEASE(bracket-back-right) |
| UVA4 | Get and hold another dowel |
| RPA5 | GET(dowel) |
| RPA6 | HOLD(dowel) |
| UVA5 | I am mounting the third set of brackets on a dowel |
| UPA1 | The user affixes the foot and the right back brackets to the dowel |
| UVA6 | Finished. Release the dowel |
| RPA7 | RELEASE(dowel) |
| UVA7 | OK, done assembling right back leg |
| RMA1 | Learns action subsequence ASSEMBLE-RIGHT-BACK-LEG |
| RMA2 | learns the object RIGHT-BACK-LEG with BRACKET-FOOT, BRACKET-BACK-RIGHT and DOWEL as fillers of HAS-OBJECT-AS-PART slot of RIGHT-BACK-LEG |
| RMA3 | Adds RIGHT-BACK-LEG as a filler of HAS-OBJECT-AS-PART of CHAIR |
| Fig. 2. Assembling the right back leg of the chair. | |

third of the four legs of the chair – accompanied by associated robotic learning, as illustrated in Fig. 3. All UVAs are first analyzed and their meanings are represented as MRs. UVA1 signals the beginning of a subsequence and, together with UVA7, marks the boundaries of the complex action. All the RPAs and the UPA occurring within this span, in the order of their occurrence, will form the set of the terminal nodes in the subset of the overall process network, becoming children of the non-terminal designating the complex action of building the right back leg. Once this (sub)hierarchy is constructed, the non-terminal node at its root must be named. As the robot assembles the back leg of this type of chair for the first time, its process network

LM learns the composition of this complex action (RMA1) and labels the parent node of this small subhierarchy with the name of the concept, ASSEMBLE-RIGHT-BACK-LEG, newly learned by the robot's concept/property LM. The latter module also learns

² This effectively establishes a particular grain size of description. Should the application require it, the actions that we consider to be basic at this time can be further decomposed.

the new object-type concept RIGHT-BACK-LEG, whose existence is the effect of the above action (RMA2). It also updates the concept of CHAIR by adding RIGHT-BACK-LEG as a filler of that concept's HAS-OBJECT-AS-PART property (RMA3). For a detailed description of the algorithms for process network construction and naming the newly learned concepts, see [32]. The newly learned concepts are illustrated in Fig. 4. Note that the results of the operation of the process network LM are recorded in the HAS-EVENT-AS-PART property of a result of the operation of the concept LM. At this stage in the process, the fillers of some of the properties in the concepts are tentative and are expected to be modified/tightened at the memory management stage.

| | |
|-------------------------|---|
| ASSEMBLE-RIGHT-BACK-LEG | |
| IS-A | PHYSICAL-EVENT |
| AGENT | HUMAN, ROBOT |
| THEME | RIGHT-BACK-LEG |
| INSTRUMENT | SCREWDRIVER |
| HAS-EVENT-AS-PART | GET(ROBOT, BRACKET-FOOT) RELEASE(ROBOT, BRACKET-FOOT) GET(ROBOT, BRACKET-BACK-RIGHT) RELEASE(ROBOT, BRACKET-BACK-RIGHT) GET(ROBOT, DOWEL) HOLD(ROBOT, DOWEL) MOUNT(USER, {BRACKET-FOOT, BRACKET-BACK-RIGHT}, DOWEL) |
| PART-OF-EVENT | ASSEMBLE-CHAIR |
| EFFECT | RIGHT-BACK-LEG ;default effects are events; if filler of effect is an ;object, this means the effect is its existence |
| RIGHT-BACK-LEG | |
| IS-A | CHAIR-PART |
| HAS-OBJECT-AS-PART | BRACKET-FOOT, BRACKET-BACK-RIGHT, DOWEL |
| PART-OF-OBJECT | CHAIR |

Fig. 4. Concepts learned as a result of processing the sequence of UVAs in Fig. 3.

Memory management. Knowledge learned by the robot during each session with a human trainer (such as the sequence in Fig. 3) must be remembered so they can be used in subsequent functioning. Mutual grounding of basic actions and corresponding ontological events is recorded both in the robot's procedural memory (by augmenting the procedures implementing the robot's basic motor actions with links to their corresponding concepts in semantic memory) and in its semantic memory (by adding pm-links, see above). Newly learned process sequences and objects (such as ASSEMBLE-RIGHT-BACK-LEG and RIGHT-BACK-LEG of Fig. 4) must be incorporated in the robot's long-term semantic and episodic memories. Due to space constraints, in what follows we give an informal description of the process. A more comprehensive description describing the relevant algorithms in full detail will be published separately.

For each newly learned concept, the memory management module (MMM) first determines whether this concept should be a) added to the robot's semantic memory or b) merged with an existing concept (instead of being added as a separate one). To make this choice, the MMM uses an extension of the concept matching algorithm of [4, 29]. This algorithm is based on unification, with the added facility for naming concepts and determining their best position in the hierarchy of the ontological world model in the robot's long-term semantic memory.

- physical-event
- absorb
- apply-force
- apply
- bury
- change-location
- chemical-reaction
- cover
- demonstrate
- device-event
- disaster-event
- display
- dispose-of
- energy-event
- fasten
- filter
- gather
- laminate
- living-event
- maintain
- motion-event
- natural-event
- overcome
- produce
- ring-event
- search
- simulate
- supernatural-event
- technology-event
- tie
- untie
- wave-energy-event
- wrap

Fig. 5.

- all
- event
 - physical-event
 - produce
 - create-artifact
 - assemble
 - social-event
 - work-activity
 - manufacturing-activity
 - assemble
 - produce
 - create-artifact
 - assemble

Fig. 6.

As an illustration, suppose, the concept matching algorithm has determined that the newly learned concept ASSEMBLE-RIGHT-BACK-LEG must be added to the robot's semantic memory. In this case, the algorithm also suggests the most appropriate position for the concept in the ontological hierarchy. This is determined by comparing a) the inventory and b) sets of fillers of the properties defined for the new concepts and for the potential parents of the new concepts in the ontological hierarchy. In the example of Fig. 4, the LM used the safest, though at the same time least informative, filler (PHYSICAL-EVENT) for the IS-A property of ASSEMBLE-RIGHT-BACK-LEG). To determine the appropriate parent, the algorithm traverses the ontological hierarchy from PHYSICAL-EVENT down until it finds the closest match that does not violate recorded constraints. Fig. 5 shows the immediate children of PHYSICAL-EVENT (triangles next to concepts indicate that they are non-terminals) in our current ontology. The algorithm eventually reaches the area of the descendants of the concept EVENT presented in Fig. 6. The figure illustrates the fact that our ontology supports multiple inheritance – notably, the concept ASSEMBLE is a descendant of both CREATE-ARTIFACT and MANUFACTURING-ACTIVITY, while CREATE-ARTIFACT, by virtue of being a child of the concept PRODUCE, is a descendant of both PHYSICAL-EVENT and SOCIAL-EVENT. This mechanism facilitates an economical representation of different aspects of the meaning of ontological concepts and allows the representation of both differences and similarities among concepts.³

Returning to the example of incorporating ASSEMBLE-RIGHT-BACK-LEG into the robot's semantic memory, the MMM attempts to make the newly learned concept a child of the existing ontological hierarchy as possible, to be able to inherit the values of as many of its properties as possible without the extra work of determining them one by one). This attempt fails because the filler of the AGENT property in ASSEMBLE, ASSEMBLY-LINE-WORKER, is more constrained than the corresponding filler in the newly learned ASSEMBLE-RIGHT-BACK-LEG,⁴ which is HUMAN or ROBOT. The algorithm backtracks and succeeds in making ASSEMBLE-RIGHT-BACK-LEG a child of CREATE-ARTIFACT (and, thus, a sibling of ASSEMBLE).

Suppose now that a concept RIGHT-BACK-LEG already exists in the ontology. If this concept is as illustrated in Figure 7, then, after comparing this concept with the newly learned concept (see Fig. 4), the MMM will, instead of adding a new (possibly renamed) concept to the robot's semantic memory, just add an *optional* filler BRACKET-BACK-RIGHT to the HAS-OBJECT-AS-PART property of the standing concept of Fig. 7, thus

³ This mechanism requires additional representational means to block spurious inheritance whose description is outside the scope of this paper.

⁴ This example shows that in our formalism the names of concepts are just labels and do not carry meaning just by themselves.

| | |
|--------------------|------------------------|
| RIGHT-BACK-LEG | |
| IS-A | CHAIR-PART |
| HAS-OBJECT-AS-PART | BRACKET-FOOT, DOWEL |
| PART-OF-OBJECT | CHAIR |

Fig. 7. A possible RIGHT-BACK LEG.

| | |
|--------------------|------------------------|
| RIGHT-BACK-LEG | |
| IS-A | SOFA-PART |
| HAS-OBJECT-AS-PART | BRACKET-FOOT, DOWEL |
| PART-OF-OBJECT | SOFA |

Fig. 8. An alternative RIGHT-BACK-LEG.

| | |
|----------------------|---|
| RIGHT-BACK-LEG-CHAIR | |
| IS-A | CHAIR-PART |
| HAS-OBJECT-AS-PART | BRACKET-FOOT, DOWEL, BRACKET-BACK-RIGHT |
| PART-OF-OBJECT | CHAIR |

Fig. 9. The right-back leg of a chair.

| | |
|---------------------|------------------------|
| RIGHT-BACK-LEG-SOFA | |
| IS-A | SOFA-PART |
| HAS-OBJECT-AS-PART | BRACKET-FOOT, DOWEL |
| PART-OF-OBJECT | SOFA |

Fig.10. The right-back leg of a sofa.

merging the standing and the newly learned knowledge. If, however, the standing concept is as illustrated in Fig. 8, then, because of the mismatch of the fillers of part-of properties between the newly learned and the standing concept, the MMM will yield two new concepts (Figs. 9,10). Note the need of modifying the names of the concepts. An important case of merging several versions of a concept in one representation is the system's ability to represent the content of an action's HAS-EVENT-AS-PART property as an HTN, augmented with the means of expressing temporal ordering, optionality and valid alternative action sequences.

Semantic memory stores the robot's knowledge of concept types. So, for example, it will contain a description of all that the robot knows about chairs and chair legs. This knowledge will be used to feed the reasoning rules the robot will use during language processing, learning and decision-making. To make the robot more human-like, we also support case-based reasoning by analogy.

For this purpose, the MMM records sequences of RPAs, UPAs and UVAs that the robot represents and carries out during specific sessions of interacting with specific users in the robot's long-term episodic memory. The contents of the episodic memory will also support the robot's ability to "mindread" its various users [2, 14, 26, 27, 28] and, as a result, be able to anticipate their needs at various points during joint task execution as well as interpret their UVAs with higher confidence.

Discussion and Future Work

Summary. The system presented concentrates on robotic learning through language understanding. This learning results in extensions to and modifications of the three kinds of robotic memory – the explicit semantic and episodic memory and the implicit (skill-oriented) procedural memory. We believe that deciding to separate the learning and the operation modes of robot's functioning (cf. [1]) unduly complicates the learning process by requiring the robot to do extra work to distinguish between material that must be learned and known material and also because utterances that humans use while teaching are different pragmatically from those used in communication during regular functioning. The expected practical impact of the ability to learn and reason will include the robot's ability to a) perform complex actions without the user having to spell out a complete sequence of basic and complex actions; b) reason about task allocation between itself and the human user; and c) test and verify its knowledge through dialog

with the user, avoiding the need for large numbers of training examples required by learning by demonstration only. The inability of the state-of-the-art deep learning-based systems to provide human-level explanations is a well-known constraint on the utility of such systems. The cognitive robots we develop will still be capable of sophisticated reasoning by analogy but will be also capable of explaining their decisions and actions. Finally, our approach to learning does not depend on the availability of “big data” training materials. Instead, we model the way people learn since early childhood and throughout their lives – by being taught using natural language.

Evaluation-Related Issues. The standard evaluation approaches that rely on comparisons with human performance on selected tasks have attracted criticism in the field of cognitive systems. Thus, Jones et al. [8] argue for practical evaluation of *integrated* cognitive systems that “involves not only measuring a system’s task competence but also the properties of *adaptivity, directability, understandability, and trustworthiness.*” Cassimatis et al. [3] make a convincing argument for going beyond what they call “model fit” evaluations of cognitive systems in terms of how well their output fits human behavior in an experimental setting. Instead, they propose that a model of higher-order cognition should be evaluated on the basis of “(a) its ability to reason, solve problems, converse and learn as well as people do; (b) the breadth of situations in which it can do so; and (c) the parsimony of the mechanisms it posits.”

The critique of the model fit approach is motivated by the realization that good-quality comparison measures are difficult to build and that they are typically not explanatory – they do not provide reasons for performance discrepancies. At the same time, alternative methods of quantifying the factors suggested in [3, 8] must be developed before these recommendations can be followed in practice.

The demonstration of our system’s performance which we intend to present at the conference will include not only a view of the joint task performance accompanied by learning but will also overtly detail and illustrate the processes of language understanding, reasoning, learning and memory management by our cognitive robot. Specifically, we will demonstrate: a) the robot’s understanding of UVAs; b) mutual grounding of basic RPAs and appropriate ontological concepts, mediated by MRs generated from relevant UVAs; c) learning task-oriented sequences of RPAs; d) learning new ontological concepts, both new object types and new (complex) event types (the abovementioned sequences of RPAs will be learned as fillers of the HAS-EVENT-AS-PART properties of the newly learned complex event types); and e) incorporating the newly grounded and learned concepts into the robot’s semantic and episodic memories.

Language Understanding. Language understanding in the area of cognitive systems is not to be equated with current mainstream natural language processing (NLP), a thriving R&D area that methodologically is almost entirely “knowledge-lean” [15]. Some of the more application-oriented projects in cognitive systems support their language processing needs with such knowledge-lean methods, thus agreeing to a lower level of quality in exchange for broader coverage and faster system development. Longer-term, more theoretically motivated projects seek to develop explanatory models of human-level language processing that require knowledge (see [16] for a discussion]). The knowledge in such models supports not only language understanding but also reasoning and decision-making [30]. Indeed, deep language analysis requires knowledge

that is not overtly mentioned in the text or dialog. To be able to interpret language input, a language understanding agent must, thus, be able to reason about the world, the speech situation and other agents present in it. It must also routinely make decisions both about the interpretation of components of the input, what is implied by the input and what is omitted from the input and about whether and if so to what depth to analyze the input. (Human-level agents must be able to disregard parts of language inputs. This ability is essential for humans – it explains, for example, why we habitually interrupt one another in conversations.)

Memory Modeling. Another theoretical contribution of our work is overt modeling of the robot’s memory components. These components include an implicit memory of skills and explicit memories of concepts (objects, events and their properties) and of instances of sequences of events (episodes, represented in our system as HTNs). The link established between the implicit and explicit layers of memory allows the robot to reason about its own actions. Scheutz et al. [36] discuss methodological options for integrating robotic and cognitive architectures and propose three “generic high-level interfaces” between them – the perceptual interface, the goal interface and the action interface. In our work, the basic interaction between the implicit robotic operation and explicit cognitive operation is supported by interactions among the three components of the memory system of the robot.

Exploring the implicit/explicit boundary. The learning system we present has the potential to support further investigations of the interactions between explicit and implicit memories and their use. Mercier and Sperber [24] argue that human reasoning is typically triggered only when people must explain and justify the decisions that they (or others) made, while the decisions themselves more often than not are made on the basis of implicit skills. The environment in which our system operates offers a potential testbed of this hypothesis. We plan to expand our system’s capabilities to include first clarification dialogs and, next, explanation and justification dialogs.

Next Steps. The first enhancement of the current learning system will consist in demonstrating how, after RPAs are mutually grounded in ontological concepts, the robot will be able to carry out commands or learn new action sequences by acting on UVAs, without any need for direct triggering through software function calls or hardware operations. Next, we intend to add text generation capabilities, both to allow the robot a more active role in the learning process (by asking questions) and to enrich interaction during joint task performance with a human user. Another novel direction of work will involve adapting to particular users – modeling robots’ individuality and related phenomenological (“first-person” view) aspects of its internal organization and memory, developing and making use of mindreading capabilities [2, 14] that will in turn facilitate experimentation in collaboration among agents with different “theories of minds of others.”

Acknowledgements. This work was supported in part by Grant N00014-17-1-221 from the U.S. Office of Naval Research. Any opinions or findings expressed in this material are those of the authors and do not necessarily reflect the views of the Office of Naval Research.

References

1. Allen, J., Chambers, N., Ferguson, G., Galescu, L., Jung, H., Swift, M., and Taysom, W. (2007). PLOW: A Collaborative Task Learning Agent. Proceedings of the Twenty-Second Conference on Artificial Intelligence (AAAI-07). Vancouver, Canada.
2. Bello, P. 2011. Cognitive Foundations for a Computational Theory of Mindreading. *Advances in Cognitive Systems* 1:1-6.
3. Cassimatis, N., P. Bello and P. Langley. 2008. Ability, Breadth and Parsimony in Computational Models of Higher-Order Cognition. *Cognitive Science* 32: 1304-1322.
4. English, J. and S. Nirenburg. 2007. Ontology Learning from Text Using Automatic Ontological-Semantic Text Annotation and the Web as the Corpus. Proceedings of the AAAI-07 Spring Symposium on Machine Reading.
5. K. Erol, J. Hendler, and D. S. Nau 1994. HTN Planning: Complexity and Expressivity. Proceedings of AAAI-94.
6. Forbus, K., Riesbeck, C., Birnbaum, L., Livingston, K., Sharma, A., and Ureel, L. 2007. Integrating natural language, knowledge representation and reasoning, and analogical processing to learn by reading. Proceedings of AAAI-07.
7. Forbus, K., Lockwood, K. and Sharma, A. 2009. Steps towards a 2nd generation learning by reading system. AAAI Spring Symposium on Learning by Reading.
8. Jones, R.M., R.E. Wray and M. van Lent. 2012. Practical Evaluation of Integrated Cognitive Systems. *Advances in Cognitive Systems* 1:83-92.
9. Jung, H., J. Allen, L. Galescu, N. Chambers, M. Swift and W. Taysom. 2008. Utilizing Natural Language for One-Shot Task Learning. *Journal of Logic and Computation*, 18:3, 475-493.
10. Knepper, R.A., T. Layton, J. Romanishin, and D. Rus. 2013. IkeaBot: An autonomous multi-robot coordinated furniture assembly system. *IEEE International Conference on Robotics and Automation*.
11. Levesque, H., E. Davis, and L. Morgenstern. 2011. The Winograd schema challenge. Proceedings of AAAI Spring Symposium on Logical Formalizations of Commonsense Reasoning.
12. Lindes, P. and J. Laird. 2016. Toward Integrating Cognitive Linguistics and Cognitive Language Processing. Proceedings of ICCM-2016.
13. Mahesh, K., S. Nirenburg, S. Beale, E. Viegas, V. Raskin and B. Onyshkevych. 1997. Word Sense Disambiguation: Why Statistics When You Have These Numbers? Proceedings of TMI-97, Santa Fe.
14. McShane, M. 2014. Parameterizing mental model ascription across intelligent agents. *Interaction Studies*, 15(3): 404-425.
15. McShane, M. 2017. Natural Language Understanding (NLU, not NLP) in Cognitive Systems. AI Magazine Special Issue on Cognitive Systems.
16. McShane, M. and Nirenburg, S. 2012. A knowledge representation language for natural language processing, simulation and reasoning. *International Journal of Semantic Computing*, 6(1): 3-23.
17. McShane, M. and Nirenburg, S. 2015. Decision-making during language understanding by intelligent agents. *Artificial General Intelligence*, Volume 9205 of the series Lecture Notes in Computer Science, pp. 310-319.
18. McShane, M., and Babkin, P. 2015. Automatic ellipsis resolution: Recovering covert information from text. Proceedings of AAAI-15, pp. 572-578.
19. McShane, M., Nirenburg, S. and Beale, S. 2005. Semantics-based resolution of fragments and underspecified structures. *Traitement Automatique des Langues*, 46(1): 163-184.
20. McShane, M., Nirenburg, S., and Jarrell, B. 2013. Modeling decision-making biases. *Biologically-Inspired Cognitive Architectures (BICA) Journal*, 3: 39-50.
21. McShane, M., Nirenburg, S. and Beale, S. 2015. The Ontological Semantic treatment of multiword expressions. *Linguisticæ Investigationes*, 38(1): 73-110.

22. McShane, M., Nirenburg, S. and Beale, S. 2016. Language understanding with Ontological Semantics. *Advances in Cognitive Systems* 4:35-55
23. McShane, M., K. Blissett, and I. Nirenburg. 2017 Treating Unexpected Input in Incremental Semantic Analysis. Proceedings of The Fifth Annual Conference on Advances in Cognitive Systems
24. Mercier, H. and D. Sperber. 2017. *The Enigma of Reason*. Cambridge, MA: Harvard University Press.
25. Mohan, S., A. Mininger, and J. Laird. 2013. Towards an indexical model of situated language comprehension for real-world cognitive agents. *Advances in Cognitive Systems* 3: 163-182.
26. Nirenburg, S. 2011. Toward a Testbed for Modeling the Knowledge, Goals and Mental States of Others. Proceedings of the Minds and Machines Symposium at the 2011 Meeting of the International Association for Computers and Philosophy.
27. Nirenburg, Sergei and Marjorie McShane. 2012. Agents modeling agents: Incorporating ethics-related reasoning. Proceedings of the Symposium on Moral Cognition and Theory of Mind, Birmingham, UK, July.
28. Nirenburg, S., McShane, M., Beale, S., English, J. and Catizone, R. 2010. Four kinds of learning in one agent-oriented environment. Proceedings of the First Annual Meeting of the BICA Society.
29. Nirenburg, S., T. Oates and J. English. 2007. Learning by Reading by Learning to Read. Proceedings of the International Conference on Semantic Computing.
30. Nirenburg, S. and McShane, M. 2015. The interplay of language processing, reasoning and decision-making in cognitive computing. Proceedings of the 20th International Conference on Applications of Natural Language to Information Systems (NLDB-2015).
31. Nirenburg, S. and McShane, M. 2016. Slashing metaphor with Occam's Razor. Proceedings of the Fourth Annual Conference on Advances in Cognitive Systems.
32. Nirenburg, S. and P. Wood. 2017. Toward Human-Style Learning in Robots. Proceedings of the AAAI Fall Symposium on Natural Communication with Robots.
33. Piantadosi, S. T., Tily, H. & Gibson, E. (2012). The communicative function of ambiguity in language. *Cognition* 122, 280–291.
34. Roncone, A., O. Mangin and B. Scassellati. 2017. Transparent Role Assignment and Task Allocation in Human Robot Collaboration. Proceedings of International Conference on Robotics and Automation. Singapore.
35. Sharma, A., N. H. Vo, S. Gaur and C. Baral. 2015. An approach to solve Winograd schema challenge using automatically extracted commonsense knowledge. Proceedings of AAAI Spring Symposium on Logical Formalizations of Commonsense Reasoning.
36. Scheutz, M. J. Harris and P. Schmermerhorn. 2013. Systematic Integration of Cognitive and Robotic Architectures. *Advances in Cognitive Systems* 2:277-296.
37. Scheutz, M., E. Krause, B. Oosterveld, T. Frasca, and R. Platt. 2017. Spoken Instruction-Based One-Shot Object and Action Learning in a Cognitive Robotic Architecture. Proceedings of AAMAS '17.
38. Schuller. P. 2014 Tackling Winograd schemas by formalizing relevance theory in knowledge graphs. Proceedings of the International Conference on Principles of Knowledge Representation and Reasoning.
39. Vanderwende, L., A. Menezes and C. Quirk. 2015. An AMR parser for English, French, German, Spanish and Japanese and a new AMR-annotated corpus. Proceedings of the 2015 NAACL Conference. Denver, CO.
40. Viegas, E. and S. Nirenburg. 1995. The Semantic Recovery of Event Ellipsis: Its Computational Treatment. Proceedings of the 14th International Joint Conference on Artificial Intelligence. Montreal, Quebec.