



Open Research Online

The Open University's repository of research publications
and other research outputs

Grounded Language Interpretation of Robotic Commands through Structured Learning

Journal Item

How to cite:

Vanzo, Andrea; Croce, Danilo; Bastianelli, Emanuele; Basili, Roberto and Nardi, Daniele (2020). Grounded Language Interpretation of Robotic Commands through Structured Learning. *Artificial Intelligence*, 278, article no. 103181.

For guidance on citations see [FAQs](#).

© 2019 Elsevier Ltd.

Version: Accepted Manuscript

Link(s) to article on publisher's website:
<http://dx.doi.org/doi:10.1016/j.artint.2019.103181>

Copyright and Moral Rights for the articles on this site are retained by the individual authors and/or other copyright owners. For more information on Open Research Online's data [policy](#) on reuse of materials please consult the policies page.

oro.open.ac.uk

Journal Pre-proof

Grounded Language Interpretation of Robotic Commands through Structured Learning

Andrea Vanzo, Danilo Croce, Emanuele Bastianelli, Roberto Basili, Daniele Nardi

PII: S0004-3702(18)30293-5

DOI: <https://doi.org/10.1016/j.artint.2019.103181>

Reference: ARTINT 103181

To appear in: *Artificial Intelligence*

Received date: 6 June 2018

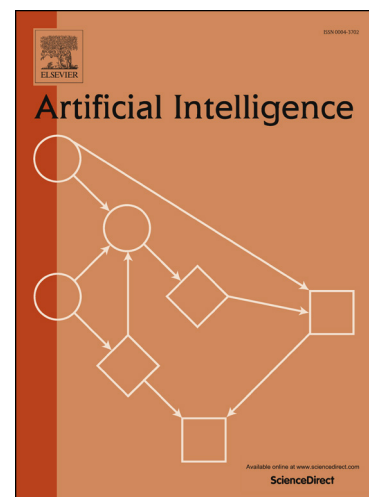
Revised date: 17 September 2019

Accepted date: 1 October 2019

Please cite this article as: A. Vanzo et al., Grounded Language Interpretation of Robotic Commands through Structured Learning, *Artif. Intell.* (2019), 103181, doi: <https://doi.org/10.1016/j.artint.2019.103181>.

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2019 Published by Elsevier.



Grounded Language Interpretation of Robotic Commands through Structured Learning

Andrea Vanzo^{a,*}, Danilo Croce^b, Emanuele Bastianelli^c, Roberto Basili^b,
Daniele Nardi^a

^a*Dept. of Computer, Control and Management Engineering “Antonio Ruberti”, Sapienza University of Rome, Italy*

^b*Dept. of Enterprise Engineering, University of Rome “Tor Vergata”, Italy*

^c*Knowledge Media Institute, The Open University, Milton Keynes, UK*

Abstract

The presence of robots in everyday life is increasing day by day at a growing pace. Industrial and working environments, health-care assistance in public or domestic areas can benefit from robots’ services to accomplish manifold tasks that are difficult and annoying for humans. In such scenarios, Natural Language interactions, enabling collaboration and robot control, are meant to be situated, in the sense that both the user and the robot access and make reference to the environment. Contextual knowledge may thus play a key role in solving inherent ambiguities of grounded language as, for example, the prepositional phrase attachment.

In this work, we present a linguistic pipeline for semantic processing of robotic commands, that combines discriminative structured learning, distributional semantics and contextual evidence extracted from the working environment. The final goal is to make the interpretation process of linguistic exchanges depending on physical, cognitive and language-dependent aspects. We present, formalize and discuss an adaptive Spoken Language Understanding chain for robotic commands, that explicitly depends on the

*Corresponding author now at Heriot-Watt University, Edinburgh, Scotland, UK.

**Work partially done while the author was at Dept. of Civil Engineering and Computer Science Engineering, University of Rome “Tor Vergata”, Italy

Email addresses: vanzo@diag.uniroma1.it (Andrea Vanzo),
croce@info.uniroma2.it (Danilo Croce), emanuele.bastianelli@open.ac.uk
(Emanuele Bastianelli), basili@info.uniroma2.it (Roberto Basili),
nardi@diag.uniroma1.it (Daniele Nardi)

operational context during both the learning and processing stages. The resulting framework allows to model heterogeneous information concerning the environment (e.g., positional information about the objects and their properties) and to inject it in the learning process. Empirical results demonstrate a significant contribution of such additional dimensions, achieving up to a 25% of relative error reduction with respect to a pipeline that only exploits linguistic evidence.

Keywords: Spoken Language Understanding, Automatic Interpretation of Robotic Commands, Grounded Language Learning, Human-Robot Interaction

1. Introduction

In the last decade, Human-Robot Interaction (HRI) is getting more and more attention within the AI and Robotics community. In fact, several different motivations are pushing forward the breakthroughs in the field. First, HRI embraces an incredibly wide range of research interests and topics. A domestic robot is expected of being able to: (i) navigate and self-localize within the environment, (ii) recognize people and objects (Vision capabilities), (iii) manipulate physical items (Grasping) and (iv) properly interact with human beings (Human-Robot Interaction). All these different challenges involve several capabilities (and so paradigms) that need to coherently interplay in order to design and build proper interactive robots. Second, domestic robots are going to be part of our everyday life in the very next future. Several robotic platforms have been already marketed and, at different level of specificity, they are able to support a variety of activities. The *iRobot Roomba* is probably the best among possible examples, due to its commercial success and the amount of innovation it contributed with. It is a vacuum cleaner capable of building a map of the environment, in order to autonomously plan and execute the cleaning of our homes.

However, though such a way of interacting with the robotic platform might be considered direct and accessible, human language is still one of the most natural ways of communication for its expressiveness and flexibility: the ability of a robot to correctly interpret users' commands is essential for proper HRI. For example, a spoken language interface would make the *Roomba* accessible to even more users.

An effective communication in natural language between humans and

robots is still challenging for the different cognitive abilities involved during the interaction. In fact, behind the simple command

“take the mug next to the keyboard” (1)

25 a number of implicit assumptions should be met in order to enable the robot
 26 to successfully execute the command. First, the user refers to entities that
 27 must exist into the environment, such as the *mug* and the *keyboard*. More-
 28 over, the robot needs a structured representation of the objects, as well as
 29 the ability to detect them. Finally, mechanisms to map lexical references to
 30 the objects must be available, in order to drive the interpretation process
 31 and the execution of a command.

We argue that the interpretation of a command must produce a logic form through the integrated use of sentence semantics, accounting for linguistic and contextual constraints. In fact, without any contextual information, the command 1 is ambiguous with respect to both syntax and semantics due to the Prepositional Phrase (PP) attachment ambiguity ([1, 2]). In the running example 1, the PP “next to the keyboard” can be attached either to the Noun Phrase (NP) or the Verb Phrase (VP), thus generating the following different syntactic structures

[VP take [NP the mug [PP next to the keyboard]]] (2)

[VP take [NP the mug] [PP next to the keyboard]] (3)

32 that evoke different meanings as well. In fact, due to the high ambiguity
 33 of the “take” word, i.e., it can be noun or verb with different meanings [3],
 34 whenever the syntactic structure of the running command is 2, “next to the
 35 keyboard” refers to “the mug”. Hence, the semantics of the command evokes
 36 a **Taking** action, in which the robot has to take the mug that is placed
 37 next to the keyboard. Conversely, if the syntactic structure is 3, “next to
 38 the keyboard” is attached to the verb phrase, indicating that the mug is
 39 located elsewhere far from the keyboard. In this case, the interpretation of
 40 the command refers to a **Bringing** action, in which robot has to bring the
 41 mug next to the keyboard, that is the goal of the action.

42 In fact, the structured representation of the environment is a discrimi-
 43 nating factor for resolving syntactic/semantic ambiguities of language, such
 44 as the attachment of the PP “next to the mug”, as well as for providing the

45 required knowledge in support of language grounding in a situated scenario.
46 While such ambiguities can be resolved through interactions, we believe that,
47 when useful resources are available, a knowledgeable system should exploit
48 them in order to minimise the user annoyance.

49 In conclusion, we foster an approach for the interpretation of robotic
50 spoken commands that is consistent with (i) the world (with all the enti-
51 ties therein), (ii) the robotic platform (with all its inner representations and
52 capabilities), and (iii) the linguistic information derived from the user's ut-
53 terance.

54 *1.1. Contributions and article outline*

55 The main contribution of this article consists of a framework for the
56 automatic understanding of robotic commands, aimed at producing inter-
57 pretations that coherently mediate among the world, the robotic platform
58 and the pure linguistic level triggered by a sentence. In fact, we support the
59 idea that the interpretation of a robotic command is not just an outcome of
60 a linguistic inference, but it is the result of a joint reasoning process involv-
61 ing both linguistic evidence and knowledge regarding the contextual physical
62 scenario. This work builds upon [4], that shows how the interpretation pro-
63 cess of a command can be made sensitive to the spatial position of perceived
64 entities within the environment. Here we make a step forward by proving
65 that the interpretation framework can be extended to richer feature spaces,
66 that allow for expressing domain properties of the involved entities, along
67 with spatial ones. To this end, this paper provides a robust formalization
68 of the Semantic Map, that collects all the semantic properties to be injected
69 in the language understanding process. Moreover, we prove the approach
70 to be language independent, with a more complete experimental session run
71 over a corpus in two different languages (i.e., English and Italian). Hence,
72 the proposed approach allows to (i) learn the interpretation function by rely-
73 ing on a corpus of annotated commands, (ii) inject grounded information
74 directly within the learning algorithm, thus integrating linguistic and con-
75 textual knowledge, and (iii) extend the features space as more specific and
76 rich information is made available. Experimental evaluations show that the
77 injection of these dimensions in the interpretation process is beneficial for
78 the correct interpretation of the real user intent, when perceptual knowledge
79 is paired with information coming from the operational domain.

80 We organize the manuscript in 7 sections. In the next section, the prob-
81 lem of natural language interpretation grounded in a robotic operating en-

82 vironment is discussed in the view of previous research and achievements in
83 literature. Section 3 provides a description of the knowledge resources, refer-
84 ring to both the linguistic assumptions and context modeling, while Section 4
85 describes the grounding process we designed, which allows to link linguistic
86 symbols to entities into the environment. In Section 5 a formal description of
87 the proposed system is provided, together with the adopted machine learning
88 techniques to feature modeling; results obtained through several experimen-
89 tal evaluations are reported in Section 6. Finally, in Section 7 we draw some
90 conclusions.

91 2. Related Work

92 The approach we propose makes use of grounded features extracted from
93 a Semantic Map [5] modeling the entities in the environment, as well as
94 *semantic* and *spatial* properties. Such features allow to drive the interpre-
95 tation process of the actions expressed by vocal commands. The realization
96 of robots that are able to intelligently interact with users within human-
97 populated environments requires techniques for linking language to actions
98 and entities into the real-world. Recently the research on this topic received
99 an incredible interest (see, for example, the workshops on Language Ground-
100 ing in Interactive Robotics [6, 7]).

101 Grounding language often requires the combination of the linguistic di-
102 mension and perception. For example, in [8], the authors make a joint use
103 of linguistic and perceptual information. Their approach leverages active
104 perception, so that linguistic symbols are directly grounded to elements ac-
105 tively perceived. Again, in [9], a Natural Language Understanding system
106 called Lucia is presented, based on Embodied Construction Grammar (ECG)
107 within the Soar architecture. Grounding is performed using knowledge from
108 the grammar itself, from the linguistic context, from the agent’s perception,
109 and from an ontology of long-term knowledge about object categories and
110 properties and actions the agent can perform. However, in these works per-
111 ceptual knowledge never modifies syntactic structures that can be generated
112 by the parser when they are incorrect. Conversely, our system is able to deal
113 with ambiguities at predicate level, allowing for selecting the interpretation
114 that is mostly coherent with the operational environment.

115 Similarly to our framework, the approaches in [10, 11] aim at ground-
116 ing language to perception through structured robot world knowledge. In
117 particular, in [11] the authors deal with the problem of using unknown out-

118 of-vocabulary words to refer to objects within the environment; the meaning
119 of such words is then acquired through dialog. Differently, we make use of
120 a mechanism based on Distributional Model of Lexical Semantics [12, 13]
121 together with phonetic similarity functions to achieve robustness (as in [14]),
122 while extracting grounded features through the lexical references contained
123 in the Semantic Map. Thanks to this mechanism, no further interactions
124 are required, and the acquisition of synonymic expressions is automatically
125 derived by reading large-scale document collections.

126 The problem of grounding semantic roles of a caption to specific areas
127 of the corresponding video is addressed in [15]. Grounding is performed on
128 both explicit and implicit roles. Semantic Role Labeling (SRL) follows a se-
129 quential tagging approach, implemented through Conditional Random Field
130 (CRF). The problem is further stressed in [16], where Gao and colleagues
131 studied a specific sub-category of the action verbs, namely the *result verbs*,
132 that are meant to cause a change of state in the *patient* referred by the verb
133 itself. In their framework, given a video and a caption, the aim is to ground
134 different semantic roles of the verb to objects in the video, relying on the
135 physical causality of verbs (i.e., physical changes that a verb may arouse
136 within the environment) as features in a CRF model. Similarly, in [17] the
137 problem of reasoning about an image and a verb is studied. In particular,
138 the authors aimed at picking the correct sense of the verb that describes the
139 action depicted into the image. In [18], the authors aim at resolving linguis-
140 tic ambiguities of a sentence paired with a video by leveraging sequential
141 labeling. The video paired with the sentence refers to one of the possible
142 interpretations of the sentence itself. Even though they make large use of
143 perceptual information to solve an SRL problem, their system requires an
144 active perception of the environment through RGB cameras. Hence, the
145 robot must have the capabilities for observing the environment at the time
146 the command is uttered. Again, in [19] the authors face the problem of
147 teaching a robot manipulator how to execute natural language commands by
148 demonstration, using video/caption pairs as valuable source of information.
149 Our system relies on a synthetic representation of the environment, acquired
150 through active interaction [20]. It allows the robot to make inferences on the
151 world it is working into, though it is not actively and directly observing the
152 surrounding environment. However, since the perception is injected in the
153 interpretation process as features for the learning machine, the framework
154 we propose can be scaled to active perception, whenever vision information
155 can be extracted and encoded into features in real-time.

156 A different perspective has been addressed in [21], where the problem of
157 PP attachment ambiguity of images' caption is resolved by leveraging the
158 corresponding image. In particular, the authors propose a joint resolution of
159 both semantic segmentation of the image and prepositional phrase attach-
160 ment. In [22] the authors exploit an RGB-D image and its caption to im-
161 prove 3D semantic segmentation and co-reference resolution in the sentences.
162 However, while the above works leverage visual context for the semantic seg-
163 mentation of images or syntax disambiguation of captions, we use a synthetic
164 representation of the context to resolve semantic ambiguities of the human
165 language, with respect to a situated interactive scenario. Our approach is
166 thus able to cope with the correct semantics of a command that has been
167 uttered in a specific context.

168 It is worth noting that approaches making joint use of language and
169 perception have been proposed to model the language grounding problem
170 also when the focus is on grounded attributes, as in [23, 24, 25]. Although
171 the underlying idea of these works is similar to ours, our aim is to produce an
172 interpretation at the predicate level, that can in turn be grounded in a robotic
173 plan corresponding to the action expressed in an utterance. Therefore, the
174 findings of such works can be considered as complementary to our proposal,
175 as while they focus just on grounding linguistic symbols into entities and
176 attributes, we leverage such a process for linking the whole interpretation to
177 the current world.

178 To summarize, our work makes the following contributions with respect
179 to the presented literature.

- 180 • The perceptual information we leverage is extracted from a synthetic
181 representation of the environment. This allows the robot to include
182 information about entities that are not present in the same environment
183 the robot is operating into.
- 184 • The discriminative nature of the proposed learning process allows to
185 scale the feature space, and to include other dimensions without re-
186 structuring the overall system. Moreover, such property is useful to
187 evaluate the contributions provided by individual features.
- 188 • In our framework, perceptual knowledge is made essential to solve am-
189 biguities at predicate level, thus affecting the syntactic interpretation
190 of sentences according to dynamic properties of the operational envi-
191 ronment.

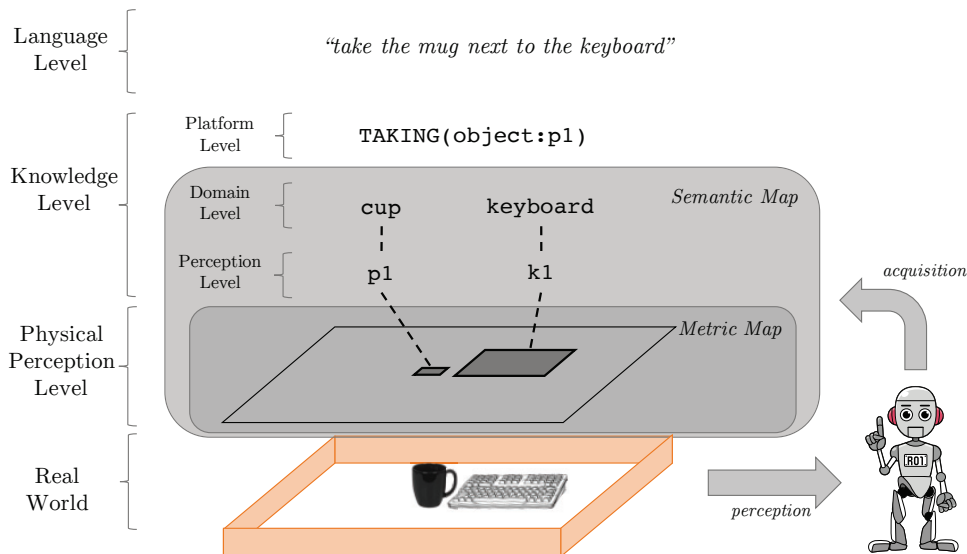


Figure 1: Layered representation of the knowledge involved in the interpretation of robotic commands

- 192
- 193
- 194
- 195
- 196
- 197
- 198
- 199
- 200
- The system is robust towards lexical variation and out-of-vocabulary words and no interaction is required to solve possible lexical ambiguities. This is achieved through Distributional Model of Lexical Semantics, used both as features for the tagging process and as principal component for grounding linguistic symbols to entities of the environment.
 - Since the grounding function is a pre-processing completely de-coupled step of the interpretation process, the mechanism is scalable to include further information that is not currently taken into account.

201

202

3. Knowledge, Language and Learning for Robotic Grounded Command Interpretation

203

204

205

206

207

208

While traditional language understanding systems mostly rely on linguistic information contained in texts (i.e., derived only from transcribed words), their application in HRI depends on a variety of other factors, including the perception of the environment. We categorize these factors into a layered representation as shown in Figure 1. First, we consider the *Language Level* as the governor of linguistic inferences: it includes observations (e.g., sequences

209 of transcribed words), as well as the linguistic assumptions of the speaker; the
 210 language level is modeled through frame-like predicates. Similarly, evidence
 211 involved by the robot’s perception of the world must be taken into account.
 212 The physical level, i.e., the *Real World*, is embodied into the *Physical Per-*
 213 *ception Level*: we assume that the robot has a synthetic image of its world,
 214 where existence and possibly other properties of entities are represented.
 215 Such representation is built by mapping the direct input of robot sensors
 216 into geometrical representations, e.g., *Metric Map*. These provide a struc-
 217 ture suitable for connecting to the *Knowledge Level*. Here *symbols*, encoded
 218 into the *Perception Level*, are used to refer to real-world entities and their
 219 properties inside the *Domain Level*. The latter comprises active concepts the
 220 robot sees, realized in a specific environment, plus general knowledge it has
 221 about the domain. All this information plays a crucial role during linguistic
 222 interactions. The integration of metric information with notions from the
 223 knowledge level provides an augmented representation of the environment,
 224 called *Semantic Map* [5]. In this map, the existence of real-world objects can
 225 be associated to *lexical* information, in the form of entity names given by a
 226 knowledge engineer or uttered by a user, as in Human-Augmented Mapping
 227 (HAM) [26, 20]. It is worth noting that the robot itself is a special entity
 228 described at this knowledge level: it does know its constituent parts as well
 229 as its capabilities that are the actions it is able to perform. To this end, we
 230 introduce an additional level (namely *Platform Level*), whose information is
 231 instantiated in a knowledge base called *Platform Model* (PM). The main aim
 232 of such a knowledge base is to enumerate all the actions the robot is able
 233 to execute. While SLU for HRI has been mostly carried out over evidence
 234 specific to the linguistic level, e.g., in [27, 28, 29, 30], this process should deal
 235 with all the aforementioned layers in a harmonized and coherent way. In fact,
 236 all linguistic primitives, including predicates and semantic arguments, corre-
 237 spond to perceptual counterparts, such as plans, robot’s actions, or entities
 238 involved in the underlying events.

239 In the following, we introduce the building blocks of our perceptually in-
 240 formed framework, defining the adopted interpretation formalism and shap-
 241 ing the perceptual information in a structured representation, i.e., the Se-
 242 mantic Map.

243 3.1. Frame-based Interpretation

244 A command interpretation system for a robotic platform must produce
 245 interpretations of user utterances. As in [31], the understanding process is

246 based on the Frame Semantics theory [32], which allows us to give a linguistic
 247 and cognitive basis to the interpretations. In particular, we consider the
 248 formalization promoted in the FrameNet [33] project, where actions expressed
 249 in user utterances are modeled as *semantic frames*. Each frame represents a
 250 micro-theory about a real-world situation, e.g., the actions of **Bringing** or
 251 **Motion**. Such micro-theories encode all the relevant information needed for
 252 their correct interpretation, represented in FrameNet via the so-called *frame*
 253 *elements*, whose role is to specify the participating entities in a frame, e.g.,
 254 the THEME frame element refers to the object that is taken in a **Bringing**
 255 action.

256 Let us consider the running example 1 “*take the mug next to the keyboard*”
 257 provided in Section 1. Depending on which syntactic structure is triggered
 258 by the contextual environment, this sentence can be intended as a command,
 259 whose effect is to instruct a robot that, in order to achieve the task, has to
 260 either

- 261 1. move towards a mug, and
- 262 2. pick it up,

263 OR

- 264 1. move towards a mug,
- 265 2. pick it up,
- 266 3. navigate to the keyboard, and
- 267 4. release the mug next to the keyboard.

To this end, a language understanding cascade should produce its FrameNet-annotated version, that can be

$$[take]_{Taking} [the\ mug\ next\ to\ the\ keyboard]_{THEME} \quad (4)$$

OR

$$[take]_{Bringing} [the\ mug]_{THEME} [next\ to\ the\ keyboard]_{GOAL} \quad (5)$$

268 depending on the configuration of the environment.

In the following, we introduce the notation used for defining an interpretation in terms of semantic frames and that will be useful to support the formal description of the proposed framework. In this respect, given a sentence s as a sequence of words w_i , i.e., $s = (w_1, \dots, w_{|s|})$, an interpretation

$\mathcal{I}(s)$ in terms of semantic frames determines a conjunction of predicates as follows:

$$\mathcal{I}(s) = \bigwedge_{i=1}^n p^i \quad (6)$$

where n is the number of predicates evoked by the sentence. Each predicate p^i is in turn represented by the pair

$$p^i = \langle f^i, Arg^i \rangle \quad (7)$$

269 where:

- 270 • $f^i \in F$ is the frame of the i^{th} predicate evoked by the sentence, where
271 F is the set of possible frames as defined in the Platform Model, e.g.,
272 **Taking, Bringing, ...**, and
- 273 • Arg^i is the set of arguments of the corresponding predicate p^i , e.g.,
274 $[the\ mug\ next\ to\ the\ keyboard]_{THEME}$ of the interpretation 4, while $[the\ mug]_{THEME}$
275 and $[next\ to\ the\ keyboard]_{GOAL}$ for the interpretation 5.

276 Every $arg_j^i \in Arg^i$ is identified by a triple $\langle a_j^i, r_j^i, h_j^i \rangle$ describing:

- 277 • the argument span a_j^i defined as subsequences of s : $a_j^i = (w_m, \dots, w_n)$
278 with $1 \leq m < n \leq |s|$, e.g., “the mug next to the keyboard” for 4 or
279 “the mug” and “next to the keyboard” for 5;
- 280 • the role label $r_j^i \in R^i$ (or frame element) associated to the current span
281 a_j^i and drawn from the vocabulary of frame elements R^i defined by
282 FrameNet for the current frame f^i , e.g., the semantic roles THEME or
283 THEME and GOAL associated to the interpretations 4 and 5, respec-
284 tively;
- 285 • the semantic head $h_j^i \in a_j^i$, as the meaning carrier word $w_k = h$ of the
286 frame argument, with $m \leq k \leq n$, e.g., “mug” for the single argu-
287 ment of interpretation 4 or “mug” and “keyboard” for the arguments
288 of interpretation 5.

289 Together with the arguments, Arg^i contains also the *lexical unit* LU that
290 anchors the predicate p_i to the text and is represented here through the same

291 structure of arguments, e.g., the verb *take*. The two different interpretations
 292 of the running example 1 will be represented through the following structures

$$\mathcal{I}(s) = \langle \mathbf{Taking}, \{ \langle (take), LU, take \rangle, \langle (the, mug, next, to, the, keyboard), THEME, mug \rangle \} \rangle$$

293 OR

$$\mathcal{I}(s) = \langle \mathbf{Bringing}, \{ \langle (take), LU, take \rangle, \langle (the, mug), THEME, mug \rangle, \langle (next, to, the, keyboard), GOAL, keyboard \rangle \} \rangle$$

294 depending on the configuration of the environment.

295 In conclusion, semantic frames can thus provide a cognitively sound bridge
 296 between the actions expressed in the language and the execution of such
 297 actions in the robot world, in terms of plans and behaviors.

298 3.2. Semantic Map

In this section we describe how to properly represent the environmental knowledge required for the interpretation process and provided by the robot. In line with [34] and according to the layered representation provided at the beginning of Section 3, we structure the Semantic Map (Figure 1) as the triple:

$$\mathcal{SM} = \langle \mathcal{R}, \mathcal{M}, \mathcal{P} \rangle \quad (8)$$

299 such as:

- 300 • \mathcal{R} is the global reference frame in which all the elements of the Semantic
 301 Map are expressed;
- 302 • \mathcal{M} is a set of geometrical elements obtained as raw sensor data ex-
 303 pressed in the reference frame \mathcal{R} and describing spatial information in
 304 a mathematical form;
- 305 • \mathcal{P} is the class hierarchy, a set of domain-dependent facts/predicates
 306 providing a semantically sound abstraction of the elements in \mathcal{M} .

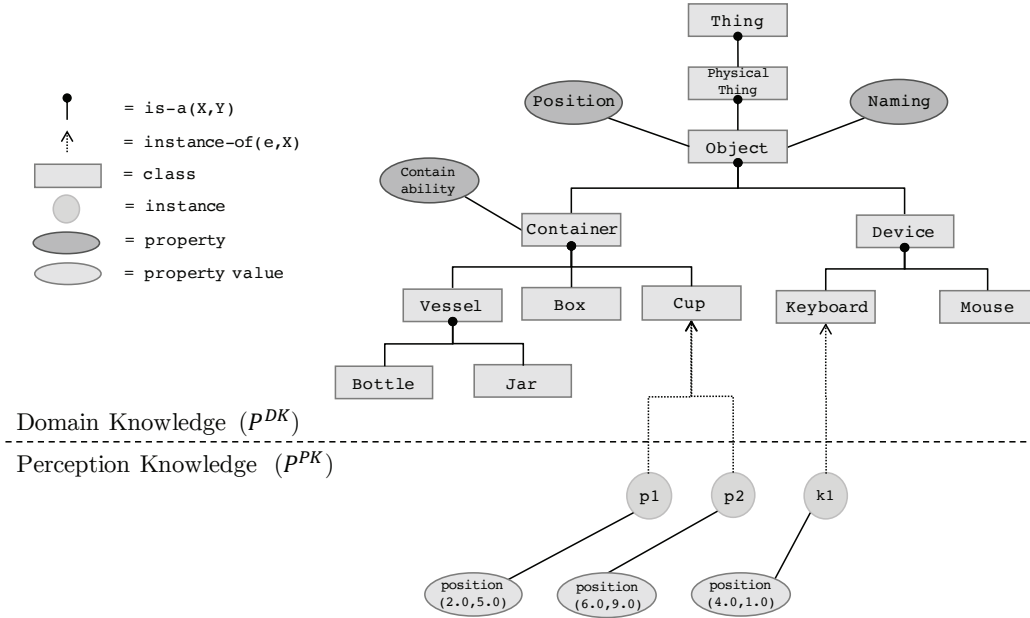


Figure 2: Sketch of the knowledge contained into a Semantic Map

\mathcal{P} is modeled as a (*Monotonic*) *Inheritance Network*. It is worth emphasizing that we do not require that the knowledge acquired through perception is fully consistent with the taxonomy of classes, as the Semantic Map is only used to support the linguistic processes addressed in this article. Hence, we decompose \mathcal{P} into two layers:

$$\mathcal{P} = \langle \mathcal{P}^{DK}, \mathcal{P}^{PK} \rangle \quad (9)$$

307 where:

- 308 • \mathcal{P}^{DK} is the *Domain Knowledge*, a conceptual knowledge base representing a hierarchy of classes, including their properties and relations, 309 *a priori* asserted to be representative of any environment; it might be 310 considered an intentional description of the robot's operation domain; 311
- 312 • \mathcal{P}^{PK} is the *Perception Knowledge*, collecting entities and properties 313 specific of the targeted environment and representing the extensional 314 knowledge, acquired by the robot.

315 The resulting structure of \mathcal{P} is shown in Figure 2, highlighting both the 316 Domain Knowledge \mathcal{P}^{DK} and the Perception Knowledge \mathcal{P}^{PK} .

317 The Semantic Map generation can follow different approaches: by rely-
 318 ing on hand-crafted ontologies and using traditional AI reasoning techniques
 319 [35, 36], by exploiting the purely automatic interpretation of perceptual out-
 320 comes [37, 38, 39], or by relying on interactions in a human-robot collabora-
 321 tion setting [40, 41]. However, the creation of the Semantic Map is out of the
 322 scope of this paper and we assume it as an available resource of the robotic
 323 system providing gold information. In fact, it is worth noting that the Se-
 324 mantic Map is an essential component of any real robot. Active perception
 325 mechanisms such as Computer Vision systems based on Deep Learning still
 326 lack in providing robust understanding of the surrounding world to support
 327 reasoning and planning mechanisms.

328 *Domain Knowledge.* The *Domain Knowledge* provides the terminology of
 329 the Semantic Map. It allows to define and structure the knowledge shared
 330 by different environments in the same domain. Such a resource can be either
 331 automatically generated consulting existing resources (e.g. WordNet [42] or
 332 ConceptNet [43]), extracted from unstructured documents (e.g. from texts
 333 present on the Web [44]), or manually created by a knowledge engineer.

334 In particular, the Domain Knowledge proposed here (Figure 2, upper
 335 part) is built upon the WordNet taxonomy and aims at modeling the hierar-
 336 chy of classes related to a domestic environment, and the domain-dependent
 337 semantic attributes.¹

338 To model the Domain Knowledge $\mathcal{P}^{\mathcal{DK}}$, we use **is-a** to define the hierarchy
 339 of classes, e.g., **is-a**(Cup, Container), and three specific properties: *Contain-*
 340 *ability*, *Naming* and *Position*. *Contain-ability* defines that all the elements
 341 of a given class might potentially contain something. *Naming* provides a set
 342 of words used to refer to a class. Conversely, *Position* is a property that is
 343 instantiated only whenever there exists an entity of the targeted class. In
 344 fact, it determines the position of the entity within the grid map of the envi-
 345 ronment. The following predicates are included into the Domain Knowledge:

- 346 • **is-contain-able**(C, t) denotes that the Contain-ability property holds
 347 for all the objects of the class C, e.g., **is-contain-able**(Cup, t);
- 348 • **naming**(C, N) defining N as the naming set, i.e., words that can be used
 349 to refer to the class C, e.g., **naming**(Table, {table, desk}).

¹We assume the attributes to be part of the Domain Knowledge, as active perception of those features is out of the scope of the article.

350 For the *Contain-able* property, the *Closed World Assumption* is applied, so
 351 that whenever the property is not defined for a class, it is assumed to be
 352 false, e.g., `is-contain-able(Keyboard, f)`.

353 It is worth noting that, for each class C , its naming can be defined by
 354 different modalities: it can be acquired through dialogic interaction, by rely-
 355 ing on the user’s preferred naming convention, extracted automatically from
 356 lexical resources or defined a priori by a knowledge engineer. In our setting,
 357 alternative naming has been provided by the combined analysis of Distri-
 358 butional Models and Lexical Databases (e.g., WordNet), and validated by a
 359 knowledge engineer.

360 *Perception Knowledge*. The Perception Knowledge \mathcal{P}^{PK} (Figure 2, lower
 361 part) is the *ABox* of the Semantic Map. It represents the actual config-
 362 uration of the current world. Hence, it is composed of elements that are
 363 actually present into the environment and perceived by the robot through
 364 its sensors.

365 \mathcal{P}^{PK} is defined through `instance-of(e, C)`, meaning that entity e is an
 366 entity of class C and inherits all the properties associated to C . Moreover,
 367 whenever a new entity is included into the Semantic Map, its corresponding
 368 *Position* must be instantiated. To this end, `position(e, x, y)` represents the
 369 value of the *Position* property for a given entity e within the grid map, in
 370 terms of (x, y) coordinates. Moreover, on top of the Semantic Map, the func-
 371 tion $d(e1, e2)$ allows to return the Euclidean distance among the entities $e1$
 372 and $e2$. This value is essential to determine whether two entities are far or
 373 near into the environment and possibly change the assumptions made during
 374 the interpretation of sentences making reference to these entities. For ex-
 375 ample, given two entities `entity-of(p1, Cup)` and `entity-of(k1, Keyboard)`
 376 whose positions are `position(p1, 2.0, 5.0)` and `position(k1, 4.0, 1.0)` respec-
 377 tively, their Euclidean distance will be $d(p1, k1) = 4.47$.

378 4. Grounding: a Side Effect of Linguistic Interpretation and Per- 379 ception

When interacting with a robot, users make references to the environment. In order for the robot to execute the requested command s , the corresponding interpretation $\mathcal{I}(s)$ must be grounded: semantic frames provided by $\mathcal{I}(s)$ are supposed to trigger grounded command instances that can be executed by the robot. Two steps are required for grounding an instantiated frame in

$\mathcal{I}(s)$. First, the frame f^i corresponding to predicate $p^i = \langle f^i, Arg^i \rangle \in \mathcal{I}(s)$ must be mapped into a behavior. Then, all the frame arguments $arg_j^i \in Arg^i$ must be explicitly associated to their corresponding actors in the plan. In fact, role labels r_j^i are paired just with the argument spans a_j^i and semantic heads h_j^i corresponding to frame elements. However, a_j^i and h_j^i play the role of anchors for the grounding onto the map: each lexical item can be used to retrieve a corresponding entity in the environment. In this respect, let $\mathcal{E}_{\mathcal{P}^{\mathcal{PK}}}$ be the set of entities populating $\mathcal{P}^{\mathcal{PK}}$, collected as:

$$\mathcal{E}_{\mathcal{P}^{\mathcal{PK}}} = \{\mathbf{e} \mid \text{instance-of}(\mathbf{e}, \cdot)\} \quad (10)$$

Then, for each entity \mathbf{e} , its corresponding naming can be gathered from the Domain Knowledge as follows:

$$\mathcal{N}(\mathbf{e}) = \{w_{\mathbf{e}} \mid \text{instance-of}(\mathbf{e}, \mathbf{C}) \wedge \text{naming}(\mathbf{C}, \mathbf{N}) \wedge w_{\mathbf{e}} \in \mathbf{N}\} \quad (11)$$

380 that is: given the entity \mathbf{e} and type \mathbf{c} , $\mathcal{N}(\mathbf{e})$ includes all the words in the
 381 naming set \mathbf{N} associated to \mathbf{c} that is defined into the Domain Knowledge
 382 $\mathcal{P}^{\mathcal{DK}}$.

The proposed linguistic grounding function $\Gamma : arg_j^i \times \mathcal{P}^{\mathcal{PK}} \rightarrow \mathcal{G}_{arg_j^i}$ is carried out by estimating to what extent the argument arg_j^i matches the naming provided for the entities in $\mathcal{P}^{\mathcal{PK}}$. Hence, $\Gamma(arg_j^i, \mathcal{P}^{\mathcal{PK}})$ produces a set of entities $\mathcal{G}_{arg_j^i}$ maximizing the lexical distance between arg_j^i and $w_{\mathbf{e}} \in \mathcal{N}(\mathbf{e})$, ordered depending on the real-valued lexical distance. Such lexical distance $g : h_j^i \times w_{\mathbf{e}} \rightarrow \mathbb{R}$ is indeed estimated as the cosine similarity between word embeddings vectors of the semantic head h_j^i (associated to arg_j^i) and the words $w_{\mathbf{e}}$ [14]. Hence, the set of grounded entities $\mathcal{G}_{arg_j^i}$ can be defined as:

$$\Gamma(arg_j^i, \mathcal{P}^{\mathcal{PK}}) \rightarrow \mathcal{G}_{arg_j^i} = \{\mathbf{e} \in \mathcal{E}_{\mathcal{P}^{\mathcal{PK}}} \mid \exists w_{\mathbf{e}} \in \mathcal{N}(\mathbf{e}) \wedge g(h, w_{\mathbf{e}}) > \tau\} \quad (12)$$

383 where τ is an empirically estimated threshold obeying to application-specific
 384 criteria.

385 The lexical semantic vectors are acquired through corpus analysis, as in
 386 Distributional Lexical Semantic paradigms. They allow to control references
 387 to elements modeling synonymy or co-hyponymy, when arguments spans,
 388 such as *cup*, are used to refer to entities with different names, e.g., a *mug*.
 389 However, depending on how the function g is modeled, it is possible to inject
 390 non-linguistic features that might be meaningful for the grounding itself. In
 391 fact, at the moment only semantic head h_j^i and naming $w_{\mathbf{e}}$ are taken into

392 account; hence, g neglects the contribution that, for example, adjectival mod-
 393 ifiers may carry, e.g., the color of an entity can be helpful in disambiguating
 394 the grounded entity, whenever two entities of the same class are present
 395 into the environment and they have different colors. The maximization of
 396 the similarity g between semantic head and entity naming corresponds to
 397 the minimization of the distance between the corresponding lexical semantic
 398 vectors and it can be extensively applied to grounding. Hence, g measures
 399 the confidence associated with individual groundings over the relevant lexical
 400 vectors.

401 It is worth noting that the grounding mechanism is here used to support
 402 the disambiguation of ambiguous commands, and it does not constitute the
 403 main contribution of the paper. Moreover, being such a process completely
 404 decoupled from the semantic parsing, different approaches for g (and there-
 405 fore of Γ) can be designed by relying on just linguistic evidence [14] or visual
 406 features [45]. However, the proposed mechanism is extensively used in this
 407 article to locate candidate grounded entities in the Semantic Map and to
 408 code them into perceptual features in the understanding process, described
 409 below.

410 5. Perceptually Informed Interpretation: the Language Understand- 411 ing Cascade

412 The interpretation framework we propose is based on a cascade of statisti-
 413 cal classification processes, modeled as sequence labeling tasks (Figure 3).
 414 The classification is applied to the entire sentence and is modeled as the
 415 Markovian formulation of a structured SVM (i.e., SVM^{hmm} proposed in [46]).
 416 In general, this learning algorithm combines a local discriminative model,
 417 which estimates the individual observation probabilities of a sequence, with
 418 a global generative approach to retrieve the most likely sequence, i.e., tags
 419 that better explain the whole sequence.

420 In other words, given an input sequence $\mathbf{x} = (\vec{x}_1 \dots \vec{x}_l) \in \mathcal{X}$, where \mathbf{x} is a
 421 sentence and $\vec{x}_i \in \mathbb{R}^n$ is a feature vector representing a word, the model pre-
 422 dicta a tag sequence $\mathbf{y} = (y_1 \dots y_l) \in \mathcal{Y}^+$ after learning a linear discriminant
 423 function. Note that labels y_i are specifically designed for the interpretation
 424 $\mathcal{I}(s)$. In fact, this process is obtained through the cascade of the Frame De-
 425 tection and Argument Labeling steps, where the latter is further decomposed
 426 in the Argument Identification and Argument Classification sub-steps. Each
 427 of these is mapped into a different SVM^{hmm} sequence labeling task.

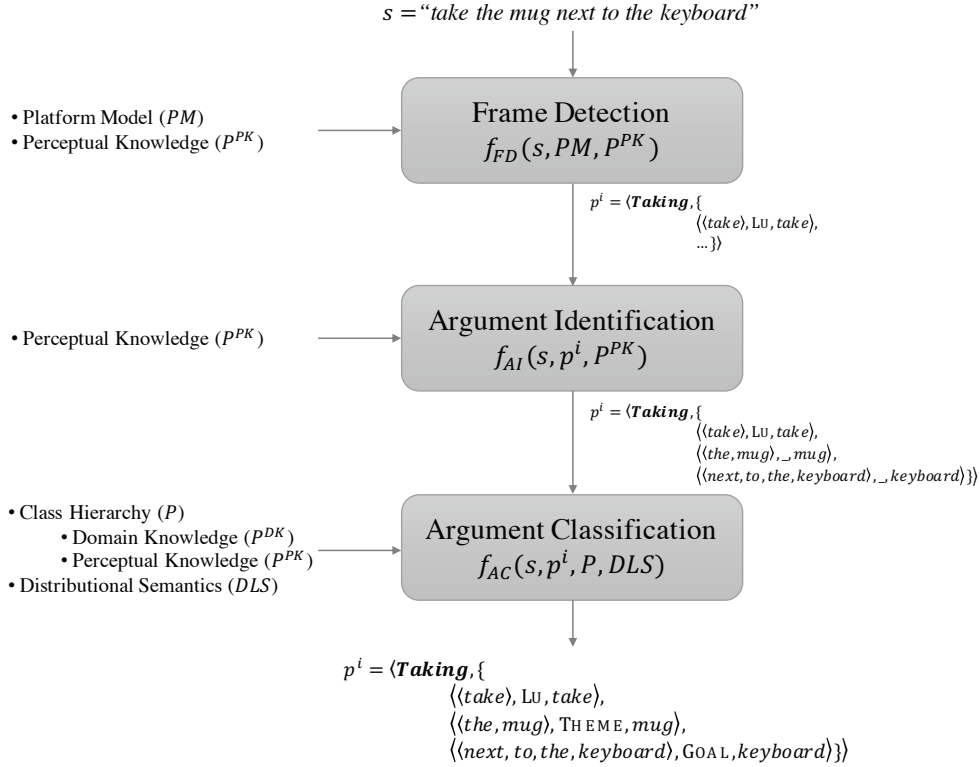


Figure 3: Processing cascade modeling the interpretation task

428 In the following, we first introduce the ML approach and then address its
 429 application to each step of the cascade.

430 5.1. The Learning Machinery

The aim of a Markovian formulation of SVM is to make the classification of a word x_i dependent on the label assigned to the previous elements in a history of length k , i.e., x_{i-k}, \dots, x_{i-1} . Given this history, a sequence of k step-specific labels can be retrieved, in the form y_{i-k}, \dots, y_{i-1} . In order to make the classification of x_i dependent also from the history, we augment the feature vector of x_i introducing a vector of transitions $\psi_{tr}(y_{i-k}, \dots, y_{i-1}) \in \mathbb{R}^l$: ψ_{tr} is a boolean vector where the dimensions corresponding to the k labels preceding the target element x_i are set to 1. A projection function $\phi(x_i)$ is

defined to consider both the observations, i.e., ψ_{obs} and the transitions ψ_{tr} in a history of size k by concatenating the two representation as follows:

$$x_i^k = \phi(x_i; y_{i-k}, \dots, y_{i-1}) = \psi_{obs}(x_i) \parallel \psi_{tr}(y_{i-k}, \dots, y_{i-1}) \quad (13)$$

431 with $x_i^k \in \mathbb{R}^{n+l}$ and $\psi_{obs}(x_i)$ does not interfere with the original feature space.
 432 Notice that the vector concatenation is here denoted by the symbol \parallel , and
 433 that linear kernel functions are applied to different types of features, ranging
 434 from linguistic to world-specific features.

The feature space operated by ψ_{obs} is defined by linear combinations of kernels to integrate independent properties. In fact, through the application of linear kernels, the space defined by the linear combination is equivalent to the space obtained by juxtaposing the vectors on which each kernel operates. More formally, assuming that K is a linear kernel, e.g., the inner product, and being x_i, x_j two instances, each composed by two vector representations a and b (i.e., $x_{i_a}, x_{i_b}, x_{j_a}, x_{j_b}$), then the resulting Kernel $K(x_i, x_j)$ will be the combination of the contributions given by Kernels working on the two representations (i.e., $K_a(x_{i_a}, x_{j_a})$ and $K_b(x_{i_b}, x_{j_b})$, respectively), that can be approximated through the concatenation of vectors $x_{i_a} \parallel x_{i_b}$ and $x_{j_a} \parallel x_{j_b}$:

$$K(x_i, x_j) = K_a(x_{i_a}, x_{j_a}) + K_b(x_{i_b}, x_{j_b}) = \langle x_{i_a} \parallel x_{i_b}, x_{j_a} \parallel x_{j_b} \rangle \quad (14)$$

435 Conversely, $\psi_{obs}(x_i) = x_{i_a} \parallel x_{i_b}$.²

436 At training time, we use the SVM learning algorithm LibLinear, pro-
 437 posed in [47] and implemented in KeLP [48] in a One-Vs-All schema over
 438 the feature space derived by ϕ , so that for each y_j a linear classifier $f_j(x_i^k) =$
 439 $w_j \phi(x_i; y_{i-k}, \dots, y_{i-1}) + b_j$ is learned. The ϕ function is computed for each
 440 element x_i by exploiting the gold label sequences. At classification time, all
 441 possible sequences $\mathbf{y} \in \mathcal{Y}^+$ should be considered in order to determine the
 442 best labeling $\hat{\mathbf{y}} = F(\mathbf{x}, k)$, where k is the size of the history used to enrich
 443 x_i , that is:

$$\begin{aligned} \hat{\mathbf{y}} = F(\mathbf{x}, k) &= \arg \max_{\mathbf{y} \in \mathcal{Y}^+} \left\{ \sum_{i=1 \dots m} f_j(x_i^k) \right\} \\ &= \arg \max_{\mathbf{y} \in \mathcal{Y}^+} \left\{ \sum_{i=1 \dots m} w_j \phi(x_i; y_{i-k}, \dots, y_{i-1}) + b_j \right\} \end{aligned}$$

444

²Before concatenating, each vector composing the observation of an instance, i.e., $\psi_{obs}(x_i)$, is normalized to have unitary norm, so that each representation equally contributes to the overall kernel estimation.

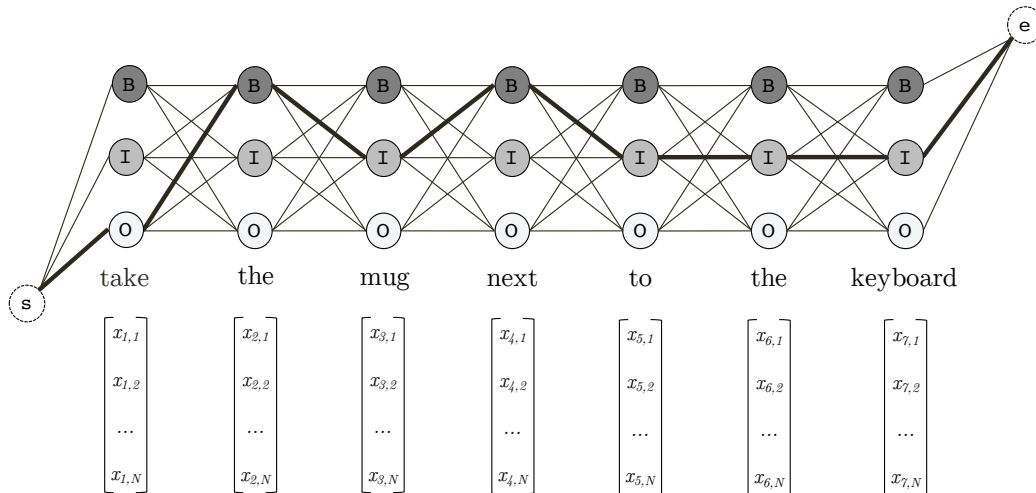


Figure 4: Viterbi decoding trellis of the Argument Identification step (Section 5.3), for the running command “take the mug next to the keyboard”, when the interpretation 5 is evoked. The label set refers to the IOB2 scheme, so that $y_i \in \{B, I, O\}$. Feature vectors x_i are obtained through the ϕ function. The best labeling $\mathbf{y} = (O, B, I, B, I, I, I) \in \mathcal{Y}^+$ is determined as the sequence maximizing the cumulative probability of individual predictions.

445 In order to reduce the computational cost, a *Viterbi-like decoding algorithm*
 446 (Figure 4) is adopted³ to derive the sequence, and thus build the
 447 augmented feature vectors through the ϕ function.

448 In the following, the different steps of the processing cascade are addressed
 449 individually.

450 5.2. Frame Detection

451 Our processing cascade starts with the **Frame Detection** (FD) step,
 452 whose aim is to find all the frames evoked by the sentence s . It corresponds
 453 to the process of filling the elements p^i in $\mathcal{I}(s)$, and can be represented as
 454 a function $f_{FD}(s, PM, \mathcal{P}^{PK})$, where s is the sentence, PM is the Platform
 455 Model and \mathcal{P}^{PK} is the Perception Knowledge. Assuming $s = \text{“take the mug$

³When applying $f_j(x_i^k)$ the classification scores are normalized through a softmax function and probability scores are derived.

456 *next to the keyboard*”, then

$$f_{FD}(s, PM, \mathcal{P}^{\mathcal{PK}}) = p^1 = \langle \mathbf{Taking}, \{ \langle \langle take \rangle, LU, take \rangle, \dots \} \rangle$$

457 for interpretation 4, while

$$f_{FD}(s, PM, \mathcal{P}^{\mathcal{PK}}) = p^1 = \langle \mathbf{Bringing}, \{ \langle \langle take \rangle, LU, take \rangle, \dots \} \rangle$$

458 for interpretation 5.

459 As already explained, the labeling process depends on linguistic informa-
 460 tion, as well as on the information derived from the Platform Model (i.e.,
 461 actions the robot is able to execute) and perceptual features extracted from
 462 the $\mathcal{P}^{\mathcal{PK}}$. In our Markovian framework states reflect frame labels, and the
 463 decoding proceeds by detecting lexical units w_k to which the proper frame
 464 f^i is assigned. This association is represented as a pair $\langle w_k, f^i \rangle$, e.g., *take-*
 465 **Taking**, *take-Bringing*. A special null label “_” is used to express the status
 466 of all other words, e.g., *the-* or *mug-*.

467 In the FD phase, each word is represented as a feature vector systemat-
 468 ically defined to be a composition between linguistic, robot-dependent and
 469 environmental observations, as hereafter detailed.

470 5.2.1. Linguistic features

471 Linguistic features here include lexical features (such as the surface or
 472 lemma of the current word and its left and right lexical contexts) and syntac-
 473 tic features (e.g., the POS-tag of the current word or the contextual POS-tag
 474 n -grams).

475 5.2.2. Robot-dependent features

476 Information about the robot coming from the PM are used to represent
 477 executable actions: these are mapped into frames through their correspond-
 478 ing LUs. The PM thus defines a set of pairing between LUs and frames,
 479 according to which boolean features are used to suggest possibly activated
 480 frames for each word in a sentence. In particular, if w_k is a verb, and $F^k \subseteq F$
 481 is the subset of frames that can be evoked by a word w_k (according to what
 482 stated in the PM), then, for every frame $f^i \in F^k$, the corresponding i -th
 483 feature of the w_k is set to **true**.

484 *5.2.3. Perceptual features*

485 In addition, features derived from the perceptual knowledge are used
 486 in the FD step as they are extracted from the $\mathcal{P}^{\mathcal{PK}}$. These “perception-
 487 based” features combine the information derived by the lexical grounding
 488 function with the syntactic dependency tree associated with s . In particu-
 489 lar, let v_h be a verb. Let $n(v_h)$ be the set of nouns governed by the verb
 490 v_h , $n(v_h) = \{w_k \mid POS(v_h) == \text{VB} \wedge POS(w_k) == \text{NN} \wedge w_k \text{ is rooted in } v_h$
 491 $\text{in the dependency (sub)tree}\}$. Let $t(v_h)$ be the set of tokens governed by
 492 the verb v_h , $t(v_h) = \{t_k \mid POS(v_h) == \text{VB} \wedge t_k \text{ is rooted in } v_h \text{ in the depen-}$
 493 $\text{dependency (sub)tree}\}$. Then the following perceptual features are extracted and
 494 associated to each token of the sentence.

495 *Grounded entities.* The number $|n(v_h)|$ of nouns governed by v_h is added as a
 496 feature to the representation of all the tokens $t_k \in t(v_h)$. Even though this is
 497 not a piece of perceptual evidence, its contribution must be considered when
 498 paired with another feature, whose aim is to explicit the number of entities
 499 that have been grounded by the tokens $w_k \in n(v_h)$. This feature is again
 500 added to the representation of all the tokens $t_k \in t(v_h)$. Formally, its value is
 501 defined as the cardinality of the grounded sets union $|\bigcup_{\forall w_k \in \text{arg}_j^i \wedge w_k \in n(v_h)} \mathcal{G}_{\text{arg}_j^i}|$.

502 *Spatial features.* This is probably the key contributing feature among the
 503 perceptual ones. In fact, it tries to capture the spatial configuration of the
 504 involved entities populating the environment, by allowing an active control of
 505 the predicate prediction, whenever the distance between objects is the only
 506 discriminating factor. Operationally, $\forall w_k \in \text{arg}_j^i \wedge w_k \in n(v_h)$, their corre-
 507 sponding grounding sets $\mathcal{G}_{\text{arg}_j^i}$ are extracted. Then, from each $\mathcal{G}_{\text{arg}_j^i}$, the most
 508 promising candidate entities (i.e., the one maximizing g) are considered and
 509 the average Euclidean spatial distance between them is computed, by relying
 510 on the predicate $\text{distance}(\mathbf{e1}, \mathbf{e2}, \mathbf{d})$. The resulting feature is a discretized
 511 version of the averaged distance (i.e., **near/far**). Such a discrete value is
 512 obtained by comparing the Euclidean distance \mathbf{d} against an empirically eval-
 513 uated threshold ϵ .

514 *5.3. Argument Identification*

515 For each identified predicate $p^i \in \mathcal{I}(s)$, the **Argument Identification**
 516 (AI) step predicts all its arguments arg_j^i , by detecting the corresponding ar-
 517 gument span a_j^i and semantic head h_j^i . This process starts filling the missing

518 elements of each j -th argument $arg_j^i \in Arg^i$. More formally, for a given sen-
 519 tence s , the i^{th} identified predicate p^i , the AI process can be summarized as
 520 the function $f_{AI}(s, p^i, \mathcal{P}^{\mathcal{PK}})$ updating the structure of $\mathcal{I}(s)$ as follows:

$$f_{AI}(s, p^i, \mathcal{P}^{\mathcal{PK}}) = p^1 = \langle \mathbf{Taking}, \{ \\ \langle \langle take \rangle, LU, take \rangle, \\ \langle \langle the, mug, next, to, the, keyboard \rangle, -, mug \rangle \} \rangle$$

521 for interpretation 4, or

$$f_{AI}(s, p^i, \mathcal{P}^{\mathcal{PK}}) = p^1 = \langle \mathbf{Bringing}, \{ \\ \langle \langle take \rangle, LU, take \rangle, \\ \langle \langle the, mug \rangle, -, mug \rangle, \\ \langle \langle next, to, the, keyboard \rangle, -, keyboard \rangle \} \rangle$$

522 for interpretation 5.

In the proposed Markovian framework, states now reflect argument boundaries between individual $arg_j^i \in Arg^i$. Following the IOB2 notation, the Begin (B), Internal (I) or Outer (O) tags are assigned to each token. For example, the result of the AI over the sentence “take the mug next to the keyboard” would be

O-take B-the I-mug I-next I-to I-the I-keyboard (Interpr. 4)

or

O-take B-the I-mug B-next I-to I-the I-keyboard (Interpr. 5)

523 5.3.1. Linguistic features

524 In this step, the same morpho-syntactic features adopted for the FD are
 525 used together with the frame f^i detected in the previous step. For each token,
 526 its lemma, right and left contexts are considered as purely lexical features.
 527 Conversely, the syntactic features used are POS-tag of the current token and
 528 left and right contextual POS-tags n -grams (see Section 5.2.1).

529 5.3.2. Perceptual features

530 Similarly to the FD step, the following dedicated features derived from
 531 the perceptual knowledge are introduced.

532 *Grounded entities.* For each noun $w_k \in arg_j^i$ such that $\mathcal{G}_{arg_j^i} \neq \emptyset$, a boolean
 533 feature is set to **true**. It is worth reminding that $\mathcal{G}_{arg_j^i}$ contains candidate
 534 entities referred by arg_j^i . Moreover, for each preposition $prep_k$, given their
 535 syntactic dependent $w_k^{dep} \in arg_j^i$, a boolean feature is set to **true** if and
 536 only if $\mathcal{G}_{arg_j^i} \neq \emptyset$. Again, for each preposition $prep_k$, the number of nouns
 537 $w_k \in arg_j^i$ on the left and on the right of $prep_k$, whose $\mathcal{G}_{arg_j^i} \neq \emptyset$, are also
 538 used as features in its corresponding feature vector.

539 *Spatial features.* For each preposition $prep_k$, we also retrieve its syntactic
 540 governor in the tree $w_f^{gov} \in arg_j^i$ and measure the average Euclidean distance
 541 in $\mathcal{P}^{\mathcal{PK}}$ between entities in $\mathcal{G}^{dep} \cup \mathcal{G}^{gov}$. As well as for the FD feature, if this
 542 score is under a given threshold ϵ , the spatial feature is set to **near**, replacing
 543 the default value of **far**.

544 5.4. Argument Classification

545 In the **Argument Classification** (AC) step, for each the frame $p^i =$
 546 $\langle f^i, Arg^i \rangle \in \mathcal{I}(s)$, all the $arg_j^i \in Arg^i$ are labeled according to their semantic
 547 role $r_j^i \in arg_j^i$, e.g., **THEME** to the argument *the mug next to the keyboard*,
 548 or **THEME** and **GOAL** to arguments *the mug* and *next to the keyboard*, re-
 549 spectively. In fact, in this step states correspond to role labels. The main
 550 novelty of this work with respect to [4] is that classification here exploits both
 551 linguistic features and semantic information about the application domain
 552 extracted from the $\mathcal{P}^{\mathcal{DK}}$. This is possible thanks to the proposed framework,
 553 which allows to inject new features that might possibly contribute to the
 554 task achievement. Consequently, AC predictions will reflect also information
 555 extracted from the Domain Knowledge.

556 Given a predicate $p^i = \langle f^i, Arg^i \rangle$, the class hierarchy \mathcal{P} , and the Distri-
 557 butional Lexical Semantics (*DLS*), the AC function can thus be written as
 558 $f_{AC}(s, p^i, \mathcal{P}, DLS)$ and produces the following complete structure

$$f_{AC}(s, p^i, \mathcal{P}, DLS) = p^1 = \langle \mathbf{Taking}, \{ \\ \langle (take), LU, take \rangle, \\ \langle (the, mug, next, to, the, keyboard), \mathbf{THEME}, mug \rangle \} \rangle$$

559 for interpretation 4, or

$$f_{AC}(s, p^i, \mathcal{P}, DLS) = p^1 = \langle \mathbf{Bringing}, \{ \\ \langle (take), LU, take \rangle, \\ \langle (the, mug), \mathbf{THEME}, mug \rangle, \\ \langle (next, to, the, keyboard), \mathbf{GOAL}, keyboard \rangle \} \rangle$$

560 for interpretation 5.

561 5.4.1. Linguistic features

562 Again, the same morpho-syntactic features adopted in both FD and AI
 563 are obtained from s , together with the frame p^i and the IOB2 tags coming
 564 from the previous stages. For each token, its lemma, right and left contexts
 565 are considered as purely lexical features. The POS-tag of the current token
 566 and left and right contextual POS-tag n -grams are used as the syntactic
 567 features (see Section 5.2.1).

568 In addition, Distributional Lexical Semantics (*DLS*) is applied to gener-
 569 alize the argument semantic head h_j^i of each argument arg_j^i : the distribu-
 570 tional (vector) representation for h_j^i is thus introduced to extend the feature
 571 vector corresponding to each $w_k \in a_j^i$, where a_j^i is a member of the triple
 572 $\langle a_j^i, r_j^i, h_j^i \rangle = arg_j^i \in Arg^i$, representing the argument span.

573 5.4.2. Domain-dependent features

574 Semantic features have been extracted from $\mathcal{P}^{\mathcal{DK}}$ to link the interpreta-
 575 tion $\mathcal{I}(s)$ to the Domain Knowledge. However, grounded entities must be
 576 provided in order to extract such attributes from the Domain knowledge.
 577 Consequently, there is an implicit dependence of the AC on the $\mathcal{P}^{\mathcal{PK}}$. In
 578 particular, the following features have been designed to further generalize
 579 the model proposed in [4].

580 *Entity-type attribute.* The *Entity-type attribute* helps in generalizing the se-
 581 mantic head of an argument through the class the corresponding grounded
 582 entity belongs to. Again, for each p^i and for each $arg_j^i \in Arg^i$, the semantic
 583 head h_j^i is grounded into a set of possible candidate entities through $\mathcal{G}_{arg_j^i}$.
 584 The most promising candidate \mathbf{e} , i.e., $\max_e g(h_j^i, w_e)$, is extracted and its
 585 class \mathbf{C} , obtained through the predicate $\text{is-a}(\mathbf{e}, \mathbf{C})$, is applied to the semantic
 586 head feature vector.

587 *Contain-ability attribute.* The *Contain-ability attribute* is a domain-dependent
 588 semantic attribute, meaning that all the elements of \mathbf{C} can contain something.
 589 To this end, for each p^i and for each $arg_j^i \in Arg^i$, the semantic head h_j^i is
 590 grounded into a set of possible candidate entities through $\mathcal{G}_{arg_j^i}$. The most
 591 promising candidate \mathbf{e} , i.e., $\max_e g(h_j^i, w_e)$, is then extracted and a boolean
 592 feature is applied to the semantic head feature vector, reflecting the value of
 593 $\text{is-contain-able}(\mathbf{C}, \mathbf{t})$, where \mathbf{C} is the class the entity \mathbf{e} belongs to.

FEATURE	FD	AI	AC
<i>Linguistic features</i>	✓	✓	✓
<i>Platform Model (PM)</i>	✓	✗	✗
<i>Domain Knowledge ($\mathcal{P}^{\mathcal{DK}}$)</i>	✗	✗	✓
<i>Perception Knowledge ($\mathcal{P}^{\mathcal{PK}}$)</i>	✓	✓	✗
<i>Distributional Lexical Semantics (DLS)</i>	✗	✗	✓

Table 1: Feature modeling of the three steps (i.e., FD, AI and AC)

594 A reader-friendly sum up is provided in Table 1 where, for each step of
 595 the processing cascade, features and resources used are shown. In particular,
 596 while AI uses only *Linguistic features* and *Perception Knowledge* $\mathcal{P}^{\mathcal{PK}}$, in FD
 597 even the Platform Model PM is exploited. Conversely, due to the nature of
 598 the task the AC step mostly relies on Domain Knowledge $\mathcal{P}^{\mathcal{DK}}$ and Distri-
 599 butional Lexical Semantics *DLS*, in order to provide effective generalization
 600 capability while choosing the correct semantic role.

601 6. Experimental Evaluation

602 The scalability of the proposed framework towards the systematic in-
 603 troduction of perceptual information has been evaluated in the semantic
 604 interpretation of utterances in a house Service Robotics scenario. The eval-
 605 uation is carried out using the Human-Robot Interaction Corpus (HuRIC),
 606 presented in Appendix A.

607 The *DLS* vectors used in the grounding mechanism $g(\cdot, \cdot)$ have been ac-
 608 quired through a Skip-gram model [13], through the `word2vec` tool. By
 609 applying the settings $min-count=50$, $window=5$, $iter=10$ and $negative=10$
 610 onto the UkWaC corpus [49], we derived 250 dimensional word vectors for
 611 more than 110,000 words. The SVM^{hmm} algorithm has been implemented
 612 within the KeLP framework [48].

613 Measures have been carried out on four tasks, according to a 10-fold
 614 evaluation schema. The first three correspond to evaluating the individ-
 615 ual interpretation steps, namely the FD, AI and AC, (Sections 6.1, 6.2 and
 616 6.3). In these tests, we assume gold annotations as input information for
 617 the task, even if they depend on a previous processing step. The last test
 618 (Section 6.4) concerns the analysis of the end-to-end interpretation chain.
 619 It thus corresponds to the ability of interpreting a fully grounded and exe-
 620 cutable command and reflects the behavior of the system in a real scenario.

621 While Perception Knowledge $\mathcal{P}^{\mathcal{PK}}$ is involved in both the FD and AI
 622 tasks, AC relies just on the Domain Knowledge $\mathcal{P}^{\mathcal{DK}}$ and the Distributional
 623 Model DLS . Hence, in order to emphasize the contribution of such informa-
 624 tion, we considered two settings.

625 The first relies just on linguistic features and information from the Se-
 626 mantic Map is neglected. We call this setting *Pure Linguistic* ($pLing$), as the
 627 interpretation is driven just by lexical/syntactic observation of the sentence.
 628 It refers to a configuration in which only the features corresponding to the
 629 first two rows of Table 1 are considered.

630 The second is a *Grounded* (*Ground*) setting. It is built upon the features
 631 designed around the Semantic Map, that has been encoded into a set of
 632 predicates \mathcal{P} , and the Distributional Model DLS , represented by Word Em-
 633 beddings. In order to enable for the extraction of meaningful properties from
 634 \mathcal{P} , grounding is based on the set \mathcal{G} of entities populating the environment
 635 and is built using the grounding function $\Gamma(arg_j^i, \mathcal{P}^{\mathcal{PK}})$. $\mathcal{P}^{\mathcal{PK}}$ features are
 636 injected into the FD and AI steps, while $\mathcal{P}^{\mathcal{DK}}$ features together with Word
 637 Embeddings are used into the AC process. Hence, this setting applies all the
 638 features defined in Table 1.

639 Results obtained in every run are reported in terms of Precision, Recall
 640 and F-Measure (F1) as micro-statistics across the 10 folds. The contribu-
 641 tion of Semantic Map information is emphasized in terms of Relative Error
 642 Reduction (RER) over F-measure with respect to the $pLing$ setting, relying
 643 just on linguistic information.

644 6.1. Frame Detection

645 In this experiment, we aim at evaluating the performance of the system
 646 in recognizing the actions evoked by the command. This step represents the
 647 entry point of the interpretation cascade: minimizing the error at this stage
 648 is essential to avoid error propagation throughout the whole pipeline.

		FD			
		Precision	Recall	F1	RER
EN	<i>pLing</i>	94.52% \pm 0.04	94.32% \pm 0.08	94.41% \pm 0.05	-
	<i>Ground</i>	95.59% \pm 0.02	96.31% \pm 0.05	95.94% \pm 0.03	27.42%
IT	<i>pLing</i>	94.84% \pm 0.22	95.58% \pm 0.19	95.19% \pm 0.19	-
	<i>Ground</i>	95.14% \pm 0.17	95.54% \pm 0.15	95.32% \pm 0.14	2.52%

Table 2: FD results: evaluating the whole span

649 Table 2 reports the results obtained for the two settings *pLing* and *Ground*,
 650 over the two datasets (i.e., English and Italian). In this case, we count a pre-
 651 diction as correct only whenever all the tokens belonging to the lexical unit
 652 LU have been correctly classified.

653 First, it is worth emphasizing that the *F1* is always higher than 94%.
 654 This means that the system will be (almost) always able to detect the correct
 655 action expressed by the command. In fact, linguistic features seem to already
 656 model the problem with a good coverage of the phenomena.

657 However, when perceptual features (extracted from the Perception Knowl-
 658 edge $\mathcal{P}^{\mathcal{PK}}$) are injected, the *F1* increases up to 95.94%, with a Relative Error
 659 Reduction of 27.42%. The contribution of such evidence is mainly due to
 660 one of the most frequent errors, concerning the ambiguity of the “take” verb.
 661 In fact, as explained in Section 1, due to the *PP attachment* ambiguity, the
 662 interpretation of such verb may differ (i.e., either *Bringing* or *Taking*)
 663 depending on the spatial configuration of the environment. As the *pLing*
 664 setting does not rely on any kind of perceptual knowledge, the system is not
 665 able to correctly discriminate among them. Hence, the resulting interpreta-
 666 tion is more likely to be wrong, as it does not reflect the semantics carried
 667 by the environment.

668 On the other hand, the Italian dataset does not seem to benefit from
 669 these features. In fact, the *RER* in such a configuration is 2.52% (i.e., from
 670 95.19% to 95.32%). This is probably due to the absence of the above linguistic
 671 phenomena in the Italian dataset.

672 6.2. Argument Identification

673 In this section, we evaluate the ability of the AI classifier in identifying
 674 the argument spans of the commands’ predicates. According to the results
 reported in Table 3, this task seems to be the most challenging one. In fact,

		AI			
		Precision	Recall	F1	RER
EN	<i>pLing</i>	89.62% ± 0.11	91.61% ± 0.03	90.59% ± 0.05	-
	<i>Ground</i>	90.04% ± 0.16	91.33% ± 0.10	90.67% ± 0.12	0.86%
IT	<i>pLing</i>	82.89% ± 0.84	85.51% ± 0.58	84.14% ± 0.68	-
	<i>Ground</i>	83.41% ± 0.84	86.30% ± 0.56	84.77% ± 0.66	4.02%

Table 3: AI results: evaluating the whole span

675

676 the F1 settles just under the 91% on the English dataset, with the *pLing* and
 677 *Ground* settings scoring 90.59% and 90.67% respectively. Moreover, in this
 678 case the Perception Knowledge does not seem to substantially contribute to
 679 the correct classification of the argument boundaries.

680 On the other hand, in the Italian setting the F1 does not exceed 85%
 681 (84.14% and 84.77% for the *pLing* and *Ground* settings). However, the per-
 682 ceptual information contributes to a slightly larger gain with respect to the
 683 one obtained on English. This is probably due to the presence of commands
 684 where the spatial configuration of the environment is essential to correctly
 685 chunk the argument spans. For example, for a command like “*porta il li-*
 686 *bro sul tavolo in cucina*” (“*bring the book on the table in the kitchen*”), the
 687 fragment *il libro sul tavolo* (*the book on the table*) may correspond to one
 688 single argument in which *sul tavolo* (*on the table*) is a spatial modifier of *il*
 689 *libro* (*the book*). In this case, *in cucina* (*in the kitchen*) composes another
 690 semantic argument. This interpretation is spatially correct whenever, within
 691 the corresponding Semantic Map, *the book* is *on the table* and the latter is
 692 outside the *kitchen*. Conversely, if *the book* is not *on the table* which is, in
 693 turn, into *the kitchen*, then *sul tavolo in cucina* (*on the table in the kitchen*)
 694 will constitute an entire argument span.

695 6.3. Argument Classification

696 For the scope of the article this experiment is the most interesting one, as
 697 here we inject the novel information extracted from the Domain Knowledge
 698 \mathcal{P}^{DK} , regarding the Contain-ability property and the class of the grounded
 699 entity.

		AC			
		Precision	Recall	F1	RER
EN	<i>pLing</i>	94.46% \pm 0.05	94.46% \pm 0.05	94.46% \pm 0.05	-
	<i>Ground</i>	95.49% \pm 0.05	95.49% \pm 0.05	95.49% \pm 0.05	18.65%
IT	<i>pLing</i>	91.52% \pm 0.23	91.52% \pm 0.23	91.52% \pm 0.23	-
	<i>Ground</i>	92.21% \pm 0.11	92.21% \pm 0.11	92.21% \pm 0.11	8.14%

Table 4: AC results: evaluating the whole span

700 As reported in Table 4, the system is able to recognize the involved entities
 701 with high accuracy, with a F1 higher than 91.50% in both the English and
 702 Italian datasets. This result is surprising when analyzing the complexity of

703 the task. In fact, the classifier is able to cope with a high level of uncertainty,
 704 as the amount of possible semantic roles is sizable, i.e., 34 for the English
 705 dataset, 27 for the Italian one.

706 Besides obtaining high accuracy in all the configurations, a twofold contri-
 707 bution is achieved when distributional information about words and domain-
 708 specific evidence is adopted. On the one hand, the *DLS* injects beneficial
 709 lexical generalization into training data: frame elements of arguments whose
 710 semantic heads are close in the vector space are seemingly tagged. For ex-
 711 ample, given the training sentence “*take the book*”, if *the book* is the THEME
 712 of a **Taking** frame, similar arguments for the same frame will receive the
 713 same role label as *volume* in “*grab the volume*”. Moreover, we provide further
 714 lexical generalization by including the class name of the grounded entity in
 715 the feature space, so that lexical references like *tv*, *tv set*, *television set*, and
 716 *television* refer to the same class *Television*.

717 On the other hand, information related to domain-dependent attributes
 718 of a given class might be helpful to solve specific errors of the AC process.
 719 For example, when including the *Contain-ability* property as a feature, we
 720 are implicitly suggesting to the learning function that an object can con-
 721 tain something. Consequently, this information allows to better discriminate
 722 whether an object must be labeled as “*Containing_object*” rather than “*Con-*
 723 *tainer_portal*”.

724 6.4. End-to-End Processing Cascade

725 In this section, we conclude our experimental evaluation by reporting the
 726 results obtained through the end-to-end processing cascade. In this case, each
 727 step is fed with the labels coming from the previous one: it thus represents
 728 a real scenario configuration, when the system is operating on a robot.

		Precision	Recall	F1	RER
		AC			
EN	<i>pLing</i>	86.12% ± 0.16	81.41% ± 0.29	83.67% ± 0.22	-
	<i>Ground</i>	89.25% ± 0.11	86.39% ± 0.22	87.77% ± 0.14	25.10%
IT	<i>pLing</i>	77.10% ± 0.81	76.08% ± 0.80	76.47% ± 0.72	-
	<i>Ground</i>	78.33% ± 0.85	77.23% ± 0.53	77.67% ± 0.60	5.09%

Table 5: Evaluating the end-to-end chain against the whole span

729 In this configuration, we chose to report only the results of the AC step
 730 (Table 5), as its output represents the end of the pipeline. Moreover, we

731 are implicitly estimating the error propagation, as each step is fed the in-
 732 formation output from the previous one. These results give thus an idea of
 733 the performance of the whole system. Note that the *DLS* and the domain-
 734 dependent features (*Ground* setting) boost the performance for both lan-
 735 guages. More specifically, the *Ground* configuration consistently outperforms
 736 the *pLing* one for English, suggesting the benefits given by the promoted fea-
 737 ture space. This behavior is less evident over the Italian dataset, even tough
 738 results confirm the general trend.

		Precision	Recall	F1	RER
		AC			
EN	<i>pLing</i>	91.04% \pm 0.07	91.54% \pm 0.07	91.28% \pm 0.06	-
	<i>Ground</i>	92.90% \pm 0.04	93.34% \pm 0.04	93.11% \pm 0.02	20.89%
IT	<i>pLing</i>	83.07% \pm 0.41	87.30% \pm 0.30	85.07% \pm 0.31	-
	<i>Ground</i>	84.15% \pm 0.33	88.83% \pm 0.27	86.35% \pm 0.24	8.58%

Table 6: Evaluating the end-to-end chain against the semantic head

In order to provide an even more realistic evaluation of the system, we measured the performance of the system by considering only the prediction over the semantic heads (Table 6). This evaluation wants to reproduce the usage of the framework, where just the semantic head is adopted to instantiate and execute a plan. For example, given the command “take the mug next to the keyboard”, together with one of its interpretations

[take]_{Taking} [the mug next to the keyboard]_{THEME},

739 only two information are required in order for the robot to execute the re-
 740 quested action, namely the type of the action **Taking** and the object to be
 741 taken, *mug*, which is pointed by the semantic head of the THEME argument.

742 The results reported in Table 6 are extremely encouraging for the applica-
 743 tion of the proposed framework in realistic scenarios. In fact, over the English
 744 dataset the F1 is always higher than 91% in the recognition of the correct
 745 label of the semantic head, along with semantic predicates and boundaries
 746 used to express intended actions. Moreover, the recognition of the full com-
 747 mand benefits from Semantic Map features, with a F1 score increasing to
 748 93.11%. In addition, the low variance suggests a good stability of the system
 749 against random selection of the training/tuning/testing sets.

750 Though with lower results, such a trend is confirmed over the Italian
 751 dataset. In fact, the difference between the two dataset is due to two reasons:

		Precision	Recall	F1	RER
		AC			
	<i>pLing</i>	91.04% \pm 0.07	91.54% \pm 0.07	91.28% \pm 0.06	-
	<i>spFeat</i>	92.63% \pm 0.05	93.07% \pm 0.05	92.83% \pm 0.03	17.75%
EN	<i>Contain</i>	92.72% \pm 0.03	93.17% \pm 0.05	92.93% \pm 0.02	18.83%
	<i>Entity</i>	92.81% \pm 0.03	93.26% \pm 0.05	93.02% \pm 0.02	19.87%
	<i>Ground</i>	92.90% \pm 0.04	93.34% \pm 0.04	93.11% \pm 0.02	20.89%
	<i>pLing</i>	83.07% \pm 0.41	87.30% \pm 0.30	85.07% \pm 0.31	-
	<i>spFeat</i>	83.21% \pm 0.44	87.83% \pm 0.36	85.39% \pm 0.34	2.10%
IT	<i>Contain</i>	83.42% \pm 0.42	88.05% \pm 0.35	85.60% \pm 0.32	3.54%
	<i>Entity</i>	83.90% \pm 0.33	88.58% \pm 0.30	86.10% \pm 0.25	6.91%
	<i>Ground</i>	84.15% \pm 0.33	88.83% \pm 0.27	86.35% \pm 0.24	8.58%

Table 7: Ablation study of the end-to-end chain against the semantic head

752 first, the different linguistic phenomena and ambiguities present in the two
753 languages do not allow to directly compare the two empirical evaluations;
754 second, the small number of examples used to train/test the models biases
755 the final results, being the Italian dataset composed of only 241 commands.
756 However, the system seems to be deployable on a real robot, with the best
757 configuration obtaining an F1 of 86.36%.

758 6.5. Ablation study

759 In order to assess the contribution of the different properties extracted
760 from the Semantic Map, we performed an ablation study of the end-to-end
761 cascade. The performance are measured by considering only the prediction
762 over the semantic head. We tested different configurations of the learning
763 function by incrementally adding the proposed features, finally reaching the
764 complete *Ground* model. The *spFeat* setting refers to a learning function,
765 where spatial features and the Distributional Model *DLS* are used along
766 with the standard linguistic features; this configuration is then extended
767 with either the Contain-ability property (*Contain*) or the *Entity type* of the
768 grounded entities (*Entity*), as discussed in Section 5. Finally, the *Ground*
769 setting that integrates *all* features has been tested.

Results are shown in Table 7. Over the English dataset we observed that the injection of spatial features reduces the relative error by 17.75% (92.83% F1). This set of features allows to solve most of the PP attachment ambiguities, like the ones mentioned before. Further improvements are obtained with

the *Contain* configuration (18.83% RER - 92.93% F1). This feature has been proven to be useful in the *Closure* frame prediction. In fact, sentences like “close the jar” and “close the door” generate two different interpretations in terms of frame elements:

[close]_{Closure} [the jar]_{CONTAINING_OBJECT}

and

[close]_{Closure} [the door]_{CONTAINER_PORTAL}

770 Marking the semantic head with the *Contain-ability* property of the grounded
 771 object allows to drive the final interpretation towards the correct one. When
 772 the *Entity type* of the grounded object is injected as a feature, we get an error
 773 reduction of 19.87% (93.02% F1). In this feature space, entities are clustered
 774 in categories, explicitly providing further generalization in the learning func-
 775 tion.

776 Conversely, over the Italian dataset we see that spatial properties do
 777 not improve consistently the performance, reducing the F1 error of 2.10%
 778 (85.39%). This result is probably biased by the language itself, with a small
 779 amount of PP attachment ambiguities in the dataset. Instead, a larger contri-
 780 bution is provided by the two domain-dependent features. For example, the
 781 *Contain* setting gets an error reduction of 5.47% (85.89% F1), by handling
 782 the same ambiguities found in the English dataset. As in the experiment over
 783 the English section, the *Entity type* provides a further improvement (6.34%
 784 RER - 86.02% F1), due to the generalization of the semantic head. Again,
 785 such a discrepancy in the results is mainly due to the different linguistic
 786 phenomena therein.

787 However, in both datasets the best performance are obtained when the
 788 full set of features is used, thus providing evidence on (i) the contribution of
 789 the different properties, and (ii) the compositionality of the feature spaces.

790 7. Conclusion

791 In this work, we presented a comprehensive framework for the definition
 792 of robust natural language interfaces for Human-Robot Interaction, specifi-
 793 cally designed for the automatic interpretation of spoken commands towards
 794 robots in domestic environments. The proposed solution allows to inject
 795 domain-dependent and environment-specific evidence into the interpretation

796 process. It relies on Frame Semantics and supports a structured learning ap-
797 proach to language processing, able to produce meaningful commands from
798 individual sentence transcriptions. A hybrid discriminative-generative learn-
799 ing method is proposed to map the interpretation process into a cascade of
800 sentence annotation tasks.

801 Starting from [4], we defined a systematic approach to enriching the ex-
802 ample representation with additional feature spaces not directly addressable
803 by the linguistic level. Our aim is to leverage the knowledge derived from
804 a semantically-enriched implementation of a robot map (i.e., its Semantic
805 Map), by expressing information about the existence and position of entities
806 surrounding the robot, along with their semantic properties. Observations
807 extracted from the Semantic Map to support the interpretation are then ex-
808 pressed through a feature modeling process. Thanks to the discriminative
809 nature of the adopted learning mechanism, such features have been injected
810 directly in the algorithm. As a result, command interpretation is made de-
811 pendent on the robot’s perception of the environment.

812 The proposed machine learning processes have been trained by using an
813 extended version of *HuRIC*, the Human Robot Interaction Corpus. The cor-
814 pus, originally composed of examples in English, now contains also a subset
815 of examples in Italian. Moreover, each example has been paired with the
816 corresponding Semantic Map, linking the command to the environment in
817 which it has been uttered and enabling the extraction of valuable contextual
818 features. This novel corpus promotes the development of the proposed in-
819 terpreting cascade in more languages, but, most importantly, it will support
820 the research in grounded natural language interfaces for robots.

821 The empirical results obtained over both languages are promising, espe-
822 cially when the system is evaluated in a real scenario (end-to-end cascade
823 evaluated against the semantic head); a closer analysis brings about sev-
824 eral observations. First, the results confirm the effectiveness of the proposed
825 processing chain, even when only linguistic information is exploited. Second,
826 they prove the effect of contextual features extracted from the Semantic
827 Map, which contributed, with different extent, to the improvement of each
828 sub-task. Finally, the results promote the application of the same approach
829 in different languages. In fact, the systematic extraction of both linguistic
830 and contextual features makes the system extendable to other languages.

831 Clearly, there is room to further develop and improve the proposed frame-
832 work, starting from an extension of *HuRIC* with additional sentences and
833 semantic features, in order to consider a wider range of robotic actions and

834 properties. Specifically, future research will focus on the extension of the pro-
835 posed methodology [4], e.g., by considering spatial relations between entities
836 in the environment or their physical characteristics, such as their color, in
837 the grounding function. In conclusion, we believe that the proposed solution
838 will support further and more challenging research topics in the context of
839 HRI, such as interactive question answering or dialogue with robots.

840 References

- 841 [1] T. Aikawa, C. Quirk, L. Schwartz, Learning prepositional attachment
842 from sentence aligned bilingual corpora, Association for Machine Trans-
843 lation in the Americas, 2003.
- 844 [2] K. Church, R. Patil, Coping with syntactic ambiguity or how to put the
845 block in the box on the table, *Computational Linguistics* 8 (3-4) (1982)
846 139–149.
- 847 [3] T. Wasow, A. Perfors, D. Beaver, The puzzle of ambiguity, *Morphology*
848 and the web of grammar: Essays in memory of Steven G. Lapointe
849 (2005) 265–282.
- 850 [4] E. Bastianelli, D. Croce, A. Vanzo, R. Basili, D. Nardi, A discrimina-
851 tive approach to grounded spoken language understanding in interactive
852 robotics, in: *Proceedings of the 25th International Joint Conference on*
853 *Artificial Intelligence (IJCAI)*, New York, New York, USA, 2016.
- 854 [5] A. Nüchter, J. Hertzberg, Towards semantic maps for mobile robots,
855 *Robot. Auton. Syst.* 56 (11) (2008) 915–926.
- 856 [6] M. Bansal, C. Matuszek, J. Andreas, Y. Artzi, Y. Bisk (Eds.), *Pro-*
857 *ceedings of the First Workshop on Language Grounding for Robotics*,
858 Association for Computational Linguistics, Vancouver, Canada, 2017.
859 URL <http://www.aclweb.org/anthology/W17-28>
- 860 [7] G. Salvi, S. Dupont (Eds.), *Proceedings GLU 2017 International Work-*
861 *shop on Grounding Language Understanding*, Stockholm, Sweden, 2017.
862 doi:10.21437/GLU.2017.
- 863 [8] S. Tellex, T. Kollar, S. Dickerson, M. Walter, A. Banerjee, S. Teller,
864 N. Roy, Approaching the symbol grounding problem with probabilistic
865 graphical models, *AI Magazine* 32 (4) (2011) 64–76.

- 866 [9] P. Lindes, A. Mininger, J. R. Kirk, J. E. Laird, Grounding language
867 for interactive task learning, in: Proceedings of the First Workshop on
868 Language Grounding for Robotics, 2017, pp. 1–9.
- 869 [10] F. Kaplan, Talking AIBO: First experimentation of verbal interactions
870 with an autonomous four-legged robot, in: Proceedings of the CELE-
871 Twente workshop on interacting agents, 2000.
- 872 [11] J. Thomason, S. Zhang, R. Mooney, P. Stone, Learning to interpret
873 natural language commands through human-robot dialog, in: Proceed-
874 ings of the 2015 International Joint Conference on Artificial Intelligence
875 (IJCAI), Buenos Aires, Argentina, 2015, pp. 1923–1929.
- 876 [12] M. Sahlgren, The word-space model, Ph.D. thesis, Stockholm University
877 (2006).
- 878 [13] T. Mikolov, K. Chen, G. Corrado, J. Dean, Efficient estimation of word
879 representations in vector space, CoRR abs/1301.3781.
- 880 [14] E. Bastianelli, D. Croce, R. Basili, D. Nardi, Using semantic models
881 for robust natural language human robot interaction, in: AI* IA 2015,
882 Advances in Artificial Intelligence, Springer International Publishing,
883 2015, pp. 343–356.
- 884 [15] S. Yang, Q. Gao, C. Liu, C. Xiong, S.-C. Zhu, J. Y. Chai, Grounded
885 semantic role labeling, in: Proceedings of the 2016 Conference of the
886 North American Chapter of the Association for Computational Lin-
887 guistics: Human Language Technologies, Association for Computational
888 Linguistics, 2016, pp. 149–159.
- 889 [16] Q. Gao, M. Doering, S. Yang, J. Y. Chai, Physical causality of action
890 verbs in grounded language understanding., in: ACL (1), The Associa-
891 tion for Computer Linguistics, 2016.
- 892 [17] S. Gella, M. Lapata, F. Keller, Unsupervised visual sense disambigua-
893 tion for verbs using multimodal embeddings, in: Proceedings of the
894 2016 Conference of the North American Chapter of the Association for
895 Computational Linguistics: Human Language Technologies, Association for
896 Computational Linguistics, 2016, pp. 182–192. doi:10.18653/v1/
897 N16-1022.

- 898 [18] Y. Berzak, A. Barbu, D. Harari, B. Katz, S. Ullman, Do you see what I
899 mean? visual resolution of linguistic ambiguities, CoRR abs/1603.08079.
- 900 [19] M. Alomari, P. Duckworth, M. Hawasly, D. Hogg, A. Cohn, Natural
901 language grounding and grammar induction for robotic manipulation
902 commands (August 2017).
- 903 [20] G. Gemignani, R. Capobianco, E. Bastianelli, D. Bloisi, L. Iocchi,
904 D. Nardi, Living with robots: Interactive environmental knowledge ac-
905 quisition, *Robotics and Autonomous Systems* 78 (2016) 1–16.
- 906 [21] G. Christie, A. Laddha, A. Agrawal, S. Antol, Y. Goyal, K. Kochers-
907 berger, D. Batra, Resolving vision and language ambiguities together:
908 Joint segmentation & prepositional attachment resolution in captioned
909 scenes, *Computer Vision and Image Understanding* 163 (2017) 101 – 112,
910 language in Vision. doi:[https://doi.org/10.1016/j.cviu.2017.09.](https://doi.org/10.1016/j.cviu.2017.09.001)
911 001.
- 912 [22] C. Kong, D. Lin, M. Bansal, R. Urtasun, S. Fidler, What are you talking
913 about? text-to-image coreference, in: *Proceedings of the 2014 IEEE*
914 *Conference on Computer Vision and Pattern Recognition, CVPR '14,*
915 *IEEE Computer Society, Washington, DC, USA, 2014,* pp. 3558–3565.
916 doi:10.1109/CVPR.2014.455.
- 917 [23] C. Matuszek, N. FitzGerald, L. S. Zettlemoyer, L. Bo, D. Fox, A joint
918 model of language and perception for grounded attribute learning., in:
919 *ICML, icml.cc / Omnipress, 2012.*
- 920 [24] J. Krishnamurthy, T. Kollar, Jointly learning to parse and perceive:
921 Connecting natural language to the physical world., *TACL* 1 (2013)
922 193–206.
- 923 [25] Y. Yu, A. Eshghi, O. Lemon, Learning how to learn: An adaptive di-
924 alogue agent for incrementally learning visually grounded word mean-
925 ings, in: *Proceedings of the First Workshop on Language Grounding*
926 *for Robotics, Association for Computational Linguistics, Vancouver,*
927 *Canada, 2017,* pp. 10–19.
- 928 [26] A. Diosi, G. R. Taylor, L. Kleeman, Interactive SLAM using laser and
929 advanced sonar, in: *Proceedings of the 2005 IEEE International Con-*

- 930 ference on Robotics and Automation, ICRA 2005, April 18-22, 2005,
931 Barcelona, Spain, 2005, pp. 1103–1108.
- 932 [27] D. L. Chen, R. J. Mooney, Learning to interpret natural language navi-
933 gation instructions from observations, in: Proceedings of the 25th AAAI
934 Conference on AI, 2011, pp. 859–865.
- 935 [28] C. Matuszek, E. Herbst, L. S. Zettlemoyer, D. Fox, Learning to parse
936 natural language commands to a robot control system, in: J. P. Desai,
937 G. Dudek, O. Khatib, V. Kumar (Eds.), ISER, Vol. 88 of Springer Tracts
938 in Advanced Robotics, Springer, 2012, pp. 403–415.
- 939 [29] E. Bastianelli, G. Castellucci, D. Croce, R. Basili, D. Nardi, Effective
940 and robust natural language understanding for human-robot interaction,
941 in: Proceedings of ECAI 2014, IOS Press, 2014.
- 942 [30] E. Bastianelli, G. Castellucci, D. Croce, R. Basili, D. Nardi, Structured
943 learning for spoken language understanding in human-robot interaction,
944 The International Journal of Robotics Research 36 (5-7) (2017) 660–683.
- 945 [31] B. J. Thomas, O. C. Jenkins, Roboframenet: Verb-centric semantics for
946 actions in robot middleware, in: 2012 IEEE International Conference on
947 Robotics and Automation, 2012, pp. 4750–4755. doi:10.1109/ICRA.
948 2012.6225172.
- 949 [32] C. J. Fillmore, Frames and the semantics of understanding, *Quaderni di*
950 *Semantica* 6 (2) (1985) 222–254.
- 951 [33] C. F. Baker, C. J. Fillmore, J. B. Lowe, The berkeley framenet project,
952 in: Proceedings of ACL and COLING, 1998, pp. 86–90.
- 953 [34] R. Capobianco, J. Serafin, J. Dichtl, G. Grisetti, L. Iocchi, D. Nardi,
954 A proposal for semantic map representation and evaluation, in: *Mobile*
955 *Robots (ECMR)*, 2015 European Conference on, IEEE, 2015, pp. 1–6.
- 956 [35] D. Pangercic, M. Tenorth, B. Pitzer, M. Beetz, Semantic object maps
957 for robotic housework - representation, acquisition and use, in: 2012
958 IEEE/RSJ International Conference on Intelligent Robots and Systems
959 (IROS), Vilamoura, Portugal, 2012.

- 960 [36] C. Galindo, A. Saffiotti, S. Coradeschi, P. Buschka, J.-A. Fernandez-
961 Madrigal, J. González, Multi-hierarchical semantic maps for mobile
962 robotics, in: Intelligent Robots and Systems, 2005.(IROS 2005). 2005
963 IEEE/RSJ International Conference on, IEEE, 2005, pp. 2278–2283.
- 964 [37] P. Buschka, A. Saffiotti, A virtual sensor for room detection, in:
965 IEEE/RSJ International Conference on Intelligent Robots and Systems,
966 Vol. 1, 2002, pp. 637–642.
- 967 [38] J. Wu, H. I. Christensen, J. M. Rehg, Visual place categorization: Prob-
968 lem, dataset, and algorithm, in: 2009 IEEE/RSJ International Confer-
969 ence on Intelligent Robots and Systems, 2009, pp. 4763–4770.
- 970 [39] O. M. Mozos, H. Mizutani, R. Kurazume, T. Hasegawa, Categorization
971 of indoor places using the kinect sensor, *Sensors* 12 (5) (2012) 6695–6711.
- 972 [40] G. Gemignani, R. Capobianco, E. Bastianelli, D. D. Bloisi, L. Iocchi,
973 D. Nardi, Living with robots: Interactive environmental knowledge ac-
974 quisition, *Robotics and Autonomous Systems* 78 (Supplement C) (2016)
975 1 – 16.
- 976 [41] D. Skočaj, M. Kristan, A. Vrečko, M. Mahnič, M. Janíček, G.-J. M.
977 Kruijff, M. Hanheide, N. Hawes, T. Keller, M. Zillich, et al., A system for
978 interactive learning in dialogue with a tutor, in: Intelligent Robots and
979 Systems (IROS), 2011 IEEE/RSJ International Conference on, IEEE,
980 2011, pp. 3387–3394.
- 981 [42] G. A. Miller, Wordnet: A lexical database for english, *Commun. ACM*
982 38 (11) (1995) 39–41.
- 983 [43] R. Speer, J. Chin, C. Havasi, Conceptnet 5.5: An open multilingual
984 graph of general knowledge, in: Proceedings of the Thirty-First AAAI
985 Conference on Artificial Intelligence, February 4-9, 2017, San Francisco,
986 California, USA., 2017, pp. 4444–4451.
987 URL [http://aaai.org/ocs/index.php/AAAI/AAAI17/paper/view/
988 14972](http://aaai.org/ocs/index.php/AAAI/AAAI17/paper/view/14972)
- 989 [44] A. Shah, V. Basile, E. Cabrio, S. K. S., Frame instance extraction and
990 clustering for default knowledge building, in: Proceedings of the 1st
991 International Workshop on Application of Semantic Web technologies

- 992 in Robotics co-located with 14th Extended Semantic Web Conference
 993 (ESWC 2017), Portoroz, Slovenia, May 29th, 2017., 2017, pp. 1–10.
- 994 [45] D. S. Chaplot, K. M. Sathyendra, R. K. Pasumarthi, D. Rajagopal,
 995 R. Salakhutdinov, Gated-attention architectures for task-oriented lan-
 996 guage grounding, in: Thirty-Second AAAI Conference on Artificial In-
 997 telligence, 2018.
- 998 [46] Y. Altun, I. Tsochantaridis, T. Hofmann, Hidden Markov support vector
 999 machines, in: Proc. of ICML, 2003.
- 1000 [47] R.-E. Fan, K.-W. Chang, C.-J. Hsieh, X.-R. Wang, C.-J. Lin, Liblinear:
 1001 A library for large linear classification, *Journal of Machine Learning*
 1002 *Research* 9 (2008) 1871–1874.
- 1003 [48] S. Filice, G. Castellucci, D. Croce, R. Basili, Kelp: a kernel-based
 1004 learning platform for natural language processing, in: *Proceedings of*
 1005 *ACL2015: System Demonstrations*, Beijing, China, 2015.
- 1006 [49] A. Ferraresi, E. Zanchetta, M. Baroni, S. Bernardini, Introducing and
 1007 evaluating ukwac, a very large web-derived corpus of english, in: *Pro-*
 1008 *ceedings of the 4th Web as Corpus Workshop (WAC-4) Can we beat*
 1009 *Google*, 2008, pp. 47–54.
- 1010 [50] E. Bastianelli, G. Castellucci, D. Croce, R. Basili, D. Nardi, Huric: a
 1011 human robot interaction corpus, in: *Proceedings of LREC 2014, Reyk-*
 1012 *javik, Iceland*, 2014.

1013 **Appendix A. HuRIC: a Corpus of Robotic Commands**

1014 The proposed computational paradigms are based on machine learning
 1015 techniques and strictly depend on the availability of training data. Hence, in
 1016 order to properly train and test our framework, we developed a collection of
 1017 datasets that together form the Human-Robot Interaction Corpus (HuRIC)⁴,
 1018 formerly presented in [50].

⁴Available at <http://sag.art.uniroma2.it/huric>. The download page also contains a detailed description of the release format.

1019 HuRIC is based on Frame Semantics and captures cognitive information
 1020 about situations and events expressed in sentences. The most interesting
 1021 feature is that HuRIC is not system or robot dependent both with respect to
 1022 the surface of sentences and with respect to the adopted formalism for both
 1023 representing and extracting the interpretation of the command. In fact, it
 1024 contains information strictly related to Natural Language Semantics and it
 1025 thus results decoupled from the specific system.

	English	Italian
<i>Number of examples</i>	656	241
<i>Number of frames</i>	18	14
<i>Number of predicates</i>	762	272
<i>Number of roles</i>	34	28
<i>Predicates per sentence</i>	1.16	1.13
<i>Sentences per frame</i>	36.44	17.21
<i>Roles per sentence</i>	2.02	1.90
<i>Entities per sentence</i>	6.59	6.97

Table A.8: HuRIC: some statistics

1026 The corpus exploits different situations representing possible commands
 1027 given to a robot in a house environment. HuRIC is composed of different
 1028 subsets, characterized by different order of complexity, designed to differ-
 1029 ently stress a labeling architecture. Each dataset includes a set of audio files
 1030 representing robot commands, paired with the correct transcription. Each
 1031 sentence is then annotated with: lemmas, POS tags, dependency trees and
 1032 Frame Semantics. Semantic frames and frame elements are used to represent
 1033 the meaning of commands, as, in our view, they reflect the actions a robot
 1034 can accomplish in a home environment. In this way, HuRIC can potentially
 1035 be used to train all the modules of the processing chain presented in Section
 1036 5.

1037 HuRIC provides commands in two different languages: English and Ital-
 1038 ian. While the English subset contains 656 sentences, 241 commands are
 1039 available in Italian. Almost all Italian sentences are translations of the orig-
 1040 inal commands in English and the corpus keeps also the alignment between
 1041 those sentences. We believe these alignments will support further researches
 1042 in further areas, such as in the context of Machine Translation. The number
 1043 of annotated sentences, number of frames and further statistics are reported
 1044 in Table A.8. Detailed statistics about the number of sentences for each frame

Frame	Ex	Frame	Ex	Frame	Ex
<i>Motion</i>	143	<i>Bringing</i>	153	<i>Cotheme</i>	39
GOAL	129	THEME	153	COTHEME	39
THEME	23	GOAL	95	MANNER	9
DIRECTION	9	BENEFICIARY	56	GOAL	8
PATH	9	AGENT	39	THEME	4
MANNER	4	SOURCE	18	SPEED	1
AREA	2	MANNER	1	PATH	1
DISTANCE	1	AREA	1	AREA	1
SOURCE	1				
<i>Locating</i>	90	<i>Inspecting</i>	29	<i>Taking</i>	80
PHENOMENON	89	GROUND	28	THEME	80
GROUND	34	DESIRED_STATE	9	SOURCE	16
COGNIZER	10	INSPECTOR	5	AGENT	8
PURPOSE	5	UNWANTED_ENTITY	2	PURPOSE	2
MANNER	2				
<i>Change_direction</i>	11	<i>Arriving</i>	12	<i>Giving</i>	10
DIRECTION	11	GOAL	11	RECIPIENT	10
ANGLE	3	PATH	5	THEME	10
THEME	1	MANNER	1	DONOR	4
SPEED	1	THEME	1	REASON	1
<i>Placing</i>	52	<i>Closure</i>	19	<i>Change_operational_state</i>	49
THEME	52	CONTAINING_OBJECT	11	DEVICE	49
GOAL	51	CONTAINER_PORTAL	8	OPERATIONAL_STATE	43
AGENT	7	AGENT	7	AGENT	17
AREA	1	DEGREE	2		
<i>Being_located</i>	38	<i>Attaching</i>	11	<i>Releasing</i>	9
THEME	38	GOAL	11	THEME	9
LOCATION	34	ITEM	6	GOAL	5
PLACE	1	ITEMS	1		
<i>Perception_active</i>	6	<i>Being_in_category</i>	11	<i>Manipulation</i>	5
PHENOMENON	6	ITEM	11	ENTITY	5
MANNER	1	CATEGORY	11		

Table A.9: Distribution of frames and frame elements in the English dataset

Frame	Ex	Frame	Ex	Frame	Ex
<i>Motion</i>	51	<i>Locating</i>	27	<i>Inspecting</i>	4
GOAL	28	PHENOMENON	27	GROUND	2
DIRECTION	20	GROUND	6	UNWANTED_ENTITY	2
DISTANCE	13	MANNER	2	DESIRED_STATE	2
SPEED	8	PURPOSE	1	INSTRUMENT	1
THEME	3				
PATH	2				
MANNER	1				
SOURCE	1				
<i>Bringing</i>	59	<i>Cotheme</i>	13	<i>Placing</i>	18
THEME	60	COTHEME	13	THEME	18
BENEFICIARY	31	MANNER	6	GOAL	17
GOAL	26	GOAL	5	AREA	1
SOURCE	8				
<i>Closure</i>	10	<i>Giving</i>	7	<i>Change_direction</i>	21
CONTAINER_PORTAL	6	THEME	7	DIRECTION	21
CONTAINING_OBJECT	5	RECIPIENT	6	ANGLE	9
DEGREE	1	DONOR	1	SPEED	9
<i>Taking</i>	22	<i>Being_located</i>	14	<i>Being_in_category</i>	4
THEME	22	LOCATION	14	ITEM	4
SOURCE	8	THEME	12	CATEGORY	4
<i>Releasing</i>	8	<i>Change_operational_state</i>	14		
THEME	8	DEVICE	14		
PLACE	3				

Table A.10: Distribution of frames and frame elements in the Italian dataset

1045 and frame elements are reported in Tables A.9 and A.10 for the English and
 1046 Italian subsets, respectively.

1047 The current release of HuRIC is made available through a novel XML-
 1048 based format, whose extension is `hrc`. For each command we are able to
 1049 store: (i) the whole sentence, (ii) the list of the tokens composing it, along
 1050 with the corresponding lemma and POS tag, (iii) the dependency relations
 1051 among tokens, (iv) the semantics, expressed in terms of Frames and Frame
 1052 elements, and (v) the configuration of the environment, in terms of entities
 1053 populating the Semantic Map (SM). In fact, since in the initial HuRIC ver-
 1054 sion linguistic information were provided without an explicit representation
 1055 of the environment, we extended the corpus by pairing each utterance with
 1056 a possible reference environment. Hence, each command is paired with a
 1057 automatically generated SM, reflecting the disposition of entities matching
 1058 the interpretation, so that perceptual features can be consistently derived
 1059 for each command. Extended examples are of the form $\langle s, SM \rangle$. The map
 1060 generation process has been designed to reflect real application conditions.
 1061 First, we built a reference Knowledge Base (KB) acting as domain model
 1062 and containing classes that describe the entities of a generic home environ-

1063 ment. Then, for each sentence s , the corresponding SM is populated with
1064 the set of referred entities, plus a control set of 20 randomly-generated ad-
1065 ditional objects, all taken from the KB. The naming function \mathcal{LR} has been
1066 defined simulating the lexical references introduced by a process of Human-
1067 Augmented Mapping. The set of possible lexical alternatives (from which
1068 such \mathcal{LR} draws) has been designed to simulate free lexicalization of entities
1069 in the SM. For every class name in the KB, a range of possible polysemic
1070 variations has been defined, by automatically exploiting lexical resources,
1071 such as WordNet [42], or by corpus-analysis. The final set has been then
1072 validated by human annotators.

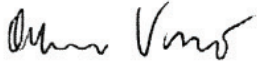

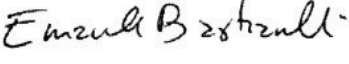
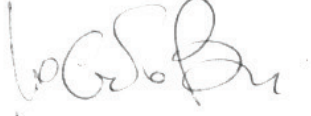

CONFLICT OF INTEREST

We wish to confirm that there are no known conflicts of interest associated with this publication and there has been no significant financial support for this work that could have influenced its outcome.

We confirm that the manuscript has been read and approved by all named authors and that there are no other persons who satisfied the criteria for authorship but are not listed. We further confirm that the order of authors listed in the manuscript has been approved by all of us. We confirm that we have given due consideration to the protection of intellectual property associated with this work and that there are no impediments to publication, including the timing of publication, with respect to intellectual property. In so doing we confirm that we have followed the regulations of our institutions concerning intellectual property.

We understand that the Corresponding Author is the sole contact for the Editorial process (including Editorial Manager and direct communications with the office). He/she is responsible for communicating with the other authors about progress, submissions of revisions and final approval of proofs. We confirm that we have provided a current, correct email address which is accessible by the Corresponding Author and which has been configured to accept email from vanzo@diag.uniroma1.it.

Signed by all authors as follows:

- Andrea Vanzo  06/05/2019
- Danilo Croce  07/05/2019
- Emanuele Bastianelli  06/05/2019
- Roberto Basili  7/5/2019
- Daniele Nardi  06/05/2019