# The IMAGACT Cross-linguistic Ontology of Action.
# A New Infrastructure for Natural Language Disambiguation

**Massimo Moneglia[1], Monica Monachini[2], Omar Calabrese[3], Alessandro Panunzi[1], Francesca Frontini[2], Gloria Gagliardi[1], Irene Russo[2]**

[1] University of Florence (Italy), [2] ILC CNR Pisa (Italy), [3] University of Siena (Italy)

moneglia@unifi.it, monica.monachini@ilc.cnr.it, omar.calabrese@unisi.it, alessandro.panunzi@unifi.it, francesca.frontini@ilc.cnr.it, gloria.gagliardi@unifi.it, irene.russo@ilc.cnr.it

## Abstract

Action verbs, which are highly frequent in speech, cause disambiguation problems that are relevant to Language Technologies. This is a consequence of the peculiar way each natural language categorizes Action i.e. it is a consequence of semantic factors. Action verbs are frequently "general", since they extend productively to actions belonging to different ontological types. Moreover, each language categorizes action in its own way and therefore the cross-linguistic reference to everyday activities is puzzling. This paper briefly sketches the IMAGACT project, which aims at setting up a cross-linguistic Ontology of Action for grounding disambiguation tasks in this crucial area of the lexicon. The project derives information on the actual variation of action verbs in English and Italian from spontaneous speech corpora, where references to action are high in frequency. Crucially it makes use of the universal language of images to identify action types, avoiding the underdeterminacy of semantic definitions. Action concept entries are implemented as prototypic scenes; this will make it easier to extend the Ontology to other languages.

**Keywords:** Action verbs, Ontology, Imagery

## 1. Introduction

In all language modalities Action verbs bear the basic information that should be processed in order to make sense of a sentence. Especially in speech, they are the most frequent structuring elements of the discourse. But Action verbs are also the less predictable linguistic type for bilingual dictionaries and they cause major problems for MT technologies. This is not because of language specific phraseology, but is rather a consequence of the peculiar way in which each natural language categorizes events; i.e. it is a consequence of semantic factors (Majid et al., 2008). In ordinary languages the most frequent Action verbs are "general", since they extend to actions belonging to different ontological types. Each language categorizes action in its own way and therefore the cross-linguistic reference to everyday activities is puzzling (Moneglia & Panunzi, 2007). For instance, considering English and Italian, the high frequency verbs *to put* and *mettere* are both general, since they extend to many different types (1, 2, 3 of Table 1), but despite a rough translation relation, they are not coextensive, since *mettere* cannot be extended to 4, which is a type extended by *to put*.

This is one example of the crucial reasons for which natural language predications are a challenge for machine translation, since the ontological entities referred to by action verbs in simple sentences are not identified and there is no guarantee that two predicates in a bilingual dictionary pick up the same entity.
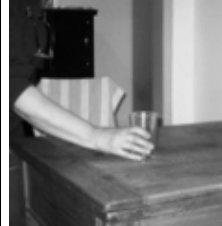
| ACTION TYPE | INSTANCES | EQUIVALENT VERBS |
|---|---|---|
|  | **Type 1** John puts the glass on the table  John mette il bicchiere sul tavolo | *to locate*  *collocare* |
|  | **Type 2** John puts the cap on the pen  John mette il tappo alla penna | *to fasten*  *inserire* |
|  | **Type 3** John puts water into the whisky  John mette l'acqua nel whisky | *to add*  *aggiungere* |
|  | **Type 4** *Mary mette su la mano  Mary puts her hand up | *to raise* |

Table 1. Action types of *to put* and *mettere*

Nevertheless, the application of general verbs to the action types in their extension is *productive*. For instance in events of type 1, *to put* will be translated into Italian with *mettere*, but no instance of type 4 using the English *put* can be translated into Italian with *mettere*:

(1) John puts a glass / a pot / a dress on the table / on the stove / on the arm-chair
(1') John mette un bicchiere / la pentola / sul tavolo / sul fornello / sulla poltrona

(2) John puts his hand / leg / shoulder up
(2') * John mette su la mano / gamba / spalla

The Italian usage of *mettere* with a body part always requires a point of reference e.g. John *puts* / *mette* his hand *in front of the picture*.
As far as the application of a verb to a type is productive, it should, in principle, be predictable. However, the range of productive variations of general verbs in various languages is unknown.
Existing repositories, and WordNet in particular (Fellbaum, 1998), may produce errors in disambiguation tasks for many reasons (Moneglia et al., 2012). The number of types recorded for each entry is high, but since the resource is not derived from corpora, peripheral meanings are not distinguished from those with high probabilities of occurrence. Moreover, descriptions given for each WordNet synset are too vague and are difficult to use for disambiguation tasks, even by expert annotators (Ng et al., 1999). Crucially, beyond these problems, the productivity of verb applications cannot be guaranteed by all *synsets* in the same manner. More specifically, verbs have various usages which depart from their actual meaning and in those meanings the translation relation cannot be predicted, since the usage is not productive. For instance, among the *synsets* of *to put* in WordNet the following is recorded:

S: (v) arrange, set up, put, order (arrange thoughts, ideas, temporal events)

We can see that translations do not run in parallel in all instances of the type; it works in 3, but for some idiosyncratic reason(s) it does not work in 4:

(3) I put my schedule in a certain way > Ho messo i miei impegni in un certo modo
(4) I put my life in a certain way > * Ho messo la mia vita in un certo modo

The distinction between productive and idiosyncratic types is crucial. Only primary usages like those in Table 1 are productive for sure, while phraseological or metaphorical usages are frequently not.
The IMAGACT project uses both corpus-based and competence-based methodologies for the simultaneous bootstrapping of a language independent action ontology from spontaneous speech resources of different languages.

The IMAGACT infrastructure will ground natural language disambiguation in all action types referred by action verbs for which a productive application can be foreseen. This paper sketches the key features of the IMAGACT project. Section 2 describes the corpus based strategy chosen for ontology building. Sections 3 and 4 briefly present the infrastructure for bootstrapping information from corpora, which comprises an innovative strategy to set up cross-linguistic correlations. In section 5, the strategy to extend the resource to an open set of languages is presented.

## 2. The Exploitation of spontaneous speech repositories

Actions specified by those verbs that are most frequently used in ordinary communication are also the actions which are more relevant for our everyday activities and constitute the universe of reference for the language. The actual use of Action oriented verbs in linguistic performance can therefore be appreciated by observing their occurrence in spontaneous speech, in which reference to action is primary. Spontaneous Speech Corpora published in the last two decades are exploited in IMAGACT to this end. IMAGACT focuses on high frequency verbs, which can provide sufficient variation in spoken corpora (500 highly ranked verbs referring to actions representing the basic verbal lexicon). IMAGACT identifies the variation of this set in the BNC-Spoken and, in parallel, in a collection of Italian Spoken corpora (C-ORAL-ROM; LABLITA; LIP; CLIPS) to get a higher probability of occurrence of relevant action types. Around 50,000 occurrences of this lexicon, derived from a 2Mw sampling of both corpora, are annotated.

## 3. The IMAGACT annotation infrastructure

### 3.1 Corpus Annotation Workflow
The corpus-based strategy relies on an induction process which separates the metaphorical and phraseological usages from proper occurrences and then classifies the proper occurrences into types, keeping granularity to its minimal level. This procedure foresees the annotation of verb occurrences in each language corpus and it is accomplished through a web based annotation interface. The procedure has been standardized in the specifications of the IMAGACT project (Moneglia & Panunzi, 2011).
Accordingly, the annotation consists of two shots leading from occurrences of each verb in a language corpus to the identification of the action types in which the verb occurs, and to the validation of the generated typology of actions productively referred to by the verb. The workflow accomplished through the annotation infrastructure follows:

*1. Standardization and gathering of occurrences into types*

1.1 - Generation of a simple sentence in the third person, representing the meaning of the instance in the corpus in a clear manner.

1.2 - Negative selection of occurrences which do not instantiate the verb in its own meaning (metaphorical or phraseological)

1.3 - Grouping of standardized proper occurrences into classes according to the number of equivalent synsets fitting with the group

1.4 - Selection of "best examples" representing the class in all possible argument structures

*2. Validation and Annotation of types*

2.1 - Comparison of the types to ensure that two claimed types do not refer to the same action (cutting granularity).

2.2 - Assignment of thematic roles and aspectual class to each best example

2.3 - Assessment that each instance of a type corresponds to the assigned best example (productivity of the type)

2.4 - Scripting of the type

In the following paragraphs the process of deriving action types from verb occurrences will be described, taking the English verb *to roll* as an example.

## 3.2. Standardization and gathering of occurrences into types

The first annotation task (1.1) is to derive from the oral context, which is frequently incomplete or fragmented, a simple sentence that properly represents the action to a possible reader and which allows a clear interpretation.

The annotator reads the context of the selected verbal occurrence (as in the dark part of Fig. 1) in order to grasp the meaning and mentally represent the referred action. The form of the standardization complies with some basic criteria:

- The sentence must be in its positive form, third person, present tense, active voice
- It can contain only essential arguments of the verb; possible specifiers of the verb arguments that are useful in grasping the meaning appear in brackets
- Generic expressions are not permitted in subject or argument position (e.g. "someone", "a man", "something" etc...);
- Basic level expression (Rosch 1978) should be preferred if available (hypernym or hyponym if necessary) or a proper name otherwise.
- Word order in sentences must be linear, with no embedding and/or distance relationships, for the purpose of being parsable

For instance, in the standardization box of Fig. 1, the annotator standardizes the selected occurrence of the verb *to roll* as "John rolls the sail". After writing the standardization, the annotator assigns the occurrence to a main "variation class" (1.2): the assignment is done by means of a competence based judgement, which is crucial at the end in determining the productive variations of the action verb in the ontology. The occurrence can be tagged as:

- PRIMARY, if and only if the verb refers to a physical action and the referred action is a proper instance of



Figure 1. Selection and standardization of corpus occurrences

the verb;
- MARKED, if the verb is used in a metaphorical sense or within a phraseology;
- SUBLEMMA, if it is a phrasal verb. These kinds of occurrences will be treated at the end of the annotation procedure (this aspect will not be considered here).

The decision concerning the status of the occurrence (PRIMARY or MARKED) makes use of an operational test roughly derived from Wittgenstein's work (Wittgenstein 1953). The occurrence is judged PRIMARY if it is possible to say to somebody who does not know the meaning of the verb V that "the referred action and similar events are what we intend with V", otherwise the occurrence is MARKED.
For instance the occurrence in "John rolls the sail" is assigned to PRIMARY variation, since indeed it is possible to point to it as a typical instance of the verb. The same can be said of "John rolls his sleeve up".
On the contrary, the instances standardized as "the registration rolls until [election] announcement" and "John rolls the words around in [his] mind" are not what you can point to in instantiating the meaning of *to roll* and therefore have been tagged as MARKED.
Only occurrences assigned to the PRIMARY variation class generate the set of productive action types stored in the ontology, and therefore they must be clustered (1.3). The workflow requires the observation of the full actual variation. When all the instances are annotated, the annotator identifies action types by means of a judgement based on the cognitive similarity among instances. Standardizations assigned to the same type should be similar for what regards:

- the involved body schema;
- the focal properties of action;
- the "equivalent verbs", i.e. the synset (Fellbaum 1998) that can be applied to the referred action (for instance, *to wind* for type 1 and *to rotate* for type 2 in the descriptions below).

Among the occurrences of a certain type, the annotator chooses the most representative one as a best example, creates the type and assigns each single standardization to it by dragging and dropping (1.4).
The overall criterion for type creation is to keep granularity to its minimal level, assigning instances to the same type as long as they can fit within the extension of a "best example".
In the case of *to roll*, occurrences have been gathered into seven types:
Type 1: "John rolls his sleeve up" contains all the standardized occurrences that instantiate the action in which the agent turns something over and over on itself (*to wind*);
Type 2: "John rolls onto his side" contains all the sentences in which the agent rotates himself along a surface (*to rotate oneself*);

Type3: "John rolls the barrel along the ground" contains all the standardized occurrences that instantiate the action in which the agent causes the object to rotate along a surface, accompanying it during its movement (*to transport*);
Type4: "John rolls the ball along the ground" contains all the standardized occurrences that instantiate the action in which the agent causes the object to rotate along a surface, by an impulse (to throw);
Type5: "John rolls his ankle around" contains all the standardized occurrences that instantiate the action in which an agent rotates a body part, typically a ball-and-socket joint (*to rotate a body part*);
Type6: "The ball rolls along the floor" contains all the standardized occurrences that instantiate the event in which an object rotates freely along a surface (unaccompanied) (*to travel freely*);
Type7: "John rolls the dough" contains all the standardized occurrences that instantiate the action in which the agent rotates a malleable material, giving it shape (*to shape*).

## 3.3. Annotation and Validation of types

One or more best examples can be added to a type in order to represent all possible argument structures: each best example has to contain the maximal argument projection, in order to represent the thematic structure of all standardized occurrences that belong to it.
The annotator edits every best example as shown in Fig.2 (step 2.2). The thematic grid must be filled, writing each argument in a separate cell and selecting the correct label from the adjacent combo-box; the tag-set for thematic role annotation consists of a restricted set of labels derived from traditional theoretical studies (Jackendoff, 1972) and current practices in computational lexicons (VerbNet). Up to three locally equivalent verbs can also be assigned.
An aspectual class (Aktionsart) is then assigned to each best example of a type by means of the Imperfective Paradox Test (Dowty, 1979). Aspect can assume three values: *event*, *process* or *state*.
In the validation of types, an excessive incidental granularity must be cut: a supervisor makes a comparison among the types created by the annotator in order to ensure that they do not refer to the same action.
Afterwards, the annotator assesses the internal consistency of the type, revising all of the occurrences previously tagged. To this end, each instance of a type is assessed with regard to its best example (2.3). During this operation, a thematic role is assigned to each argument of the standardized instance, selecting the corresponding text and then clicking on the respective button in the best example (Fig. 3).
The annotation procedure ends only when all occurrences have been validated in terms of the best examples and thematic roles have been assigned to the arguments. The annotator then produces a "script" for each type, briefly describing the action (2.4).
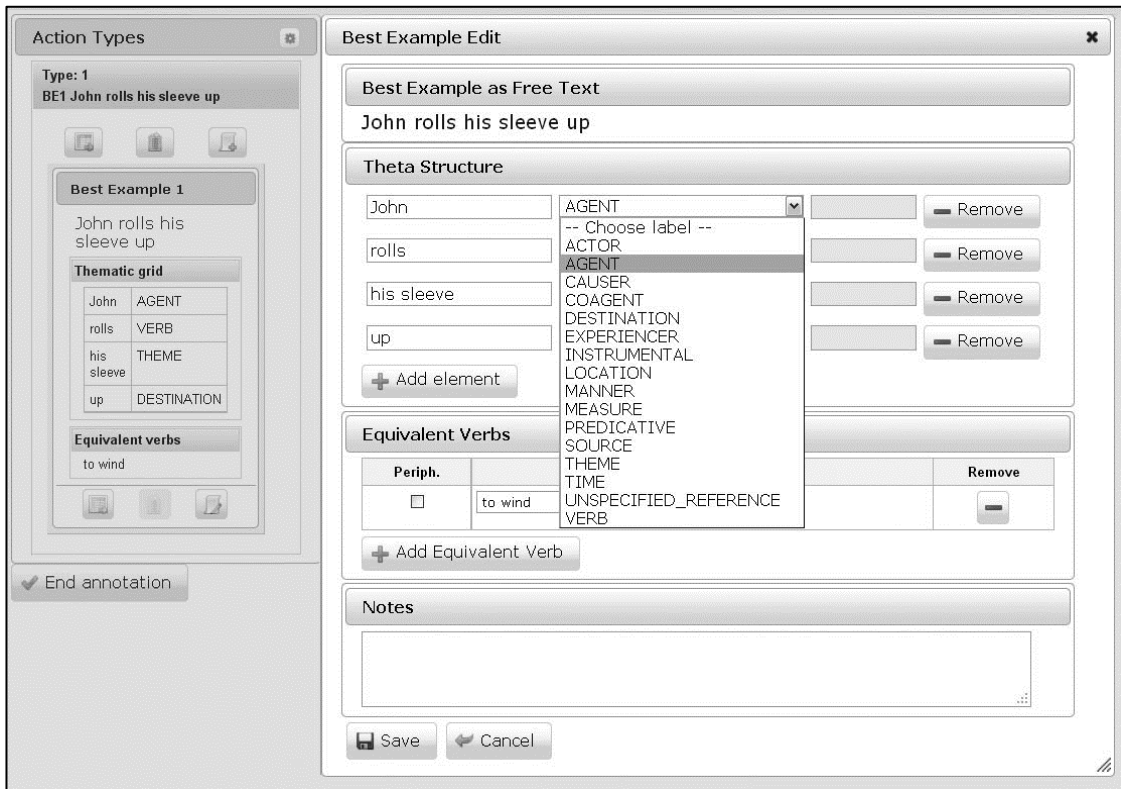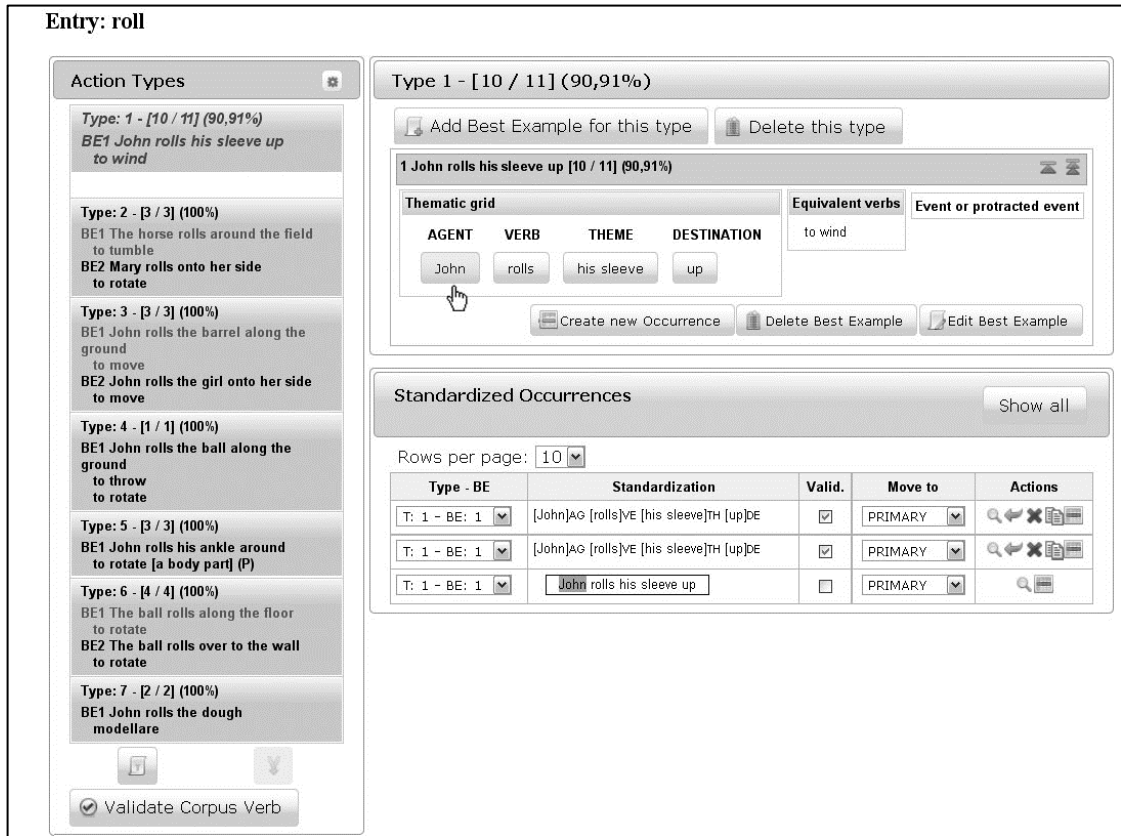
Figure 2. Semantic tagging of the Best examples



Figure 3.  Validation of action types

## 4. Defining the cross-linguistic ontology of action in a Wittgenstein-like scenario

Working with more than one language, IMAGACT will produce a language independent type inventory. Experience in ontology building has shown that the level of consensus that can be reached in defining entities referred to by language expressions is very low, since the identification of such entities relies on underdetermined definitions. The key innovation of IMAGACT is to provide a methodology which exploits the language independent capacity to appreciate similarities among *scenes*, distinguishing the *Identification* of action types from their *Definition*. Crucially, only the identification (and not the active writing of a definition) is required to set up the cross-linguistic relations.

In IMAGACT the ontology building makes use of the universal language of images, which allows the reconciliation, in a unique ontology, of the types derived from the annotation of different language corpora.

For instance, the distinction among types 1-4 in Table 1 is relevant for foreseeing the cross linguistic variation of action concepts. The difference among types is easily recognized by humans and does not require the definition of a set of differential features, which are radically underdetermined.

In Wittgenstein's terms, how can you explain to somebody what a *game* is? Just point out a play and say "this and similar things are games" (Wittgenstein, 1953). This Wittgenstein-like scenario will be exploited to identify action types at a cross-linguistic level, avoiding a direct comparison of descriptions derived from corpus annotation.

In this scenario the annotation of the English verb *to roll* will lead to a mapping of the types extracted onto scenes which will represent them, as in Fig. 4.

Then, when setting up the cross-linguistic ontology, we will discover that scene B is also extended by the Italian verb *arrotolare*, and that the variation of the English verb
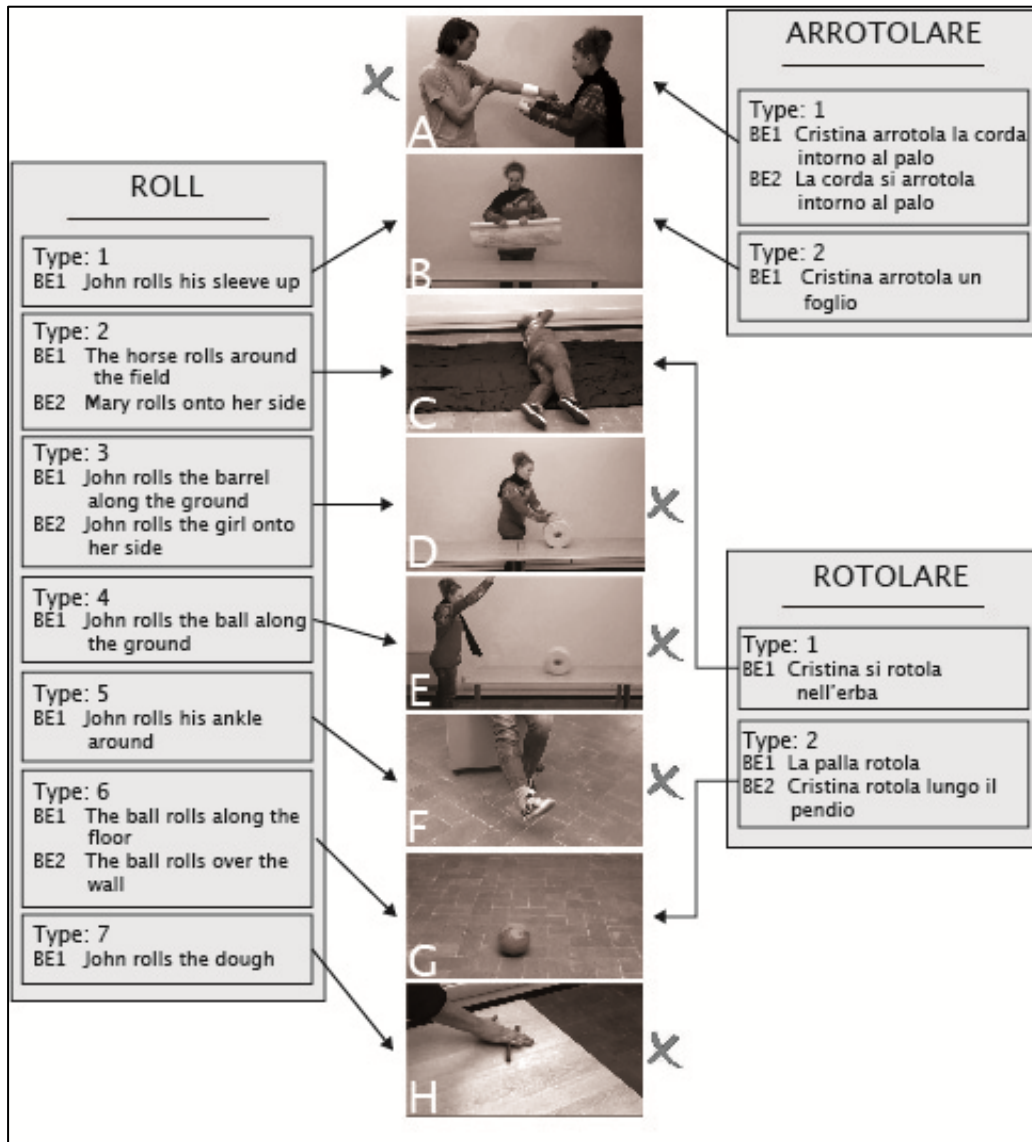


Figure 2. Cross-linguistic Gallery of Scenes representing the variation of *to roll*, *arrotolare* and *rotolare*

*to roll* is greater with respect to its Italian counterpart, since two corresponding Italian verbs (*arrotolare* and *rotolare*) find application only in a subset of the action types extended by *to roll*. Moreover the differential in meaning of the Italian verbs will be further highlighted, since it will become evident that at least one type extended by *arrotolare* (Fig.4 A) is not a possible extension of *to roll*.

Therefore, the correspondence between types derived from different language corpora will follow from their references to the same gallery of scenes. This result is obtained without the comparison among definitions given by different annotators, thus identifying the cross-linguistic mapping of action verbs in a language independent ontology, skipping the underdeterminacy of definitions.

After the corpus annotation procedure, IMAGACT will deliver a database of action types with their language encoding in English and Italian. The set of sentences derived from corpora will instantiate each represented type in connection to a prototypic scene. The gallery will work for action concepts as ImageNet does for objects.

## 5. Competence-based extension to languages and Ontology implementation

IMAGACT will deliver a database of Action types represented by stereotypic scenes in 3D. Each scene will be associated with English and Italian verbs and with the set of sentences instantiating the type in the corpora. In the second stage of the project, the database will be further exploited making use of the imagery stored therein. More specifically, the relations stored in the corpus-extracted database will constitute the starting point for a competence-based extension to other languages.

The direct representation of actions through scenes will allow the mapping of different language lexicons onto the same cross-linguistic ontology. On the basis of this outcome it will be possible to ask informants with a different language competence what verb(s) is applied in his language for each type, identified by a scene and by a set of English sentences assigned to that scene, derived from corpus occurrences. The translation relation for the lexical entries in the respective language and the validated set of equivalences in IMAGACT will follow.

This work exploits linguistic diversity to implement the action typology. For instance, contrary to English and Italian which record a lot of General Verbs, Danish has a very specific verbal lexicon (Korzen, 2005). Therefore, we expect that action types which are relevant for Danish are not identified by corpus based work on other languages. For instance type 1 of *to put* in Table 1 will record a lot of occurrences instantiating this type. Many languages will move in parallel with English, however this will not be the case when a Danish mother tongue informant will go through the same instances of the type. The informant will apply *at sætte* looking to the scene in type 1 and will verify the consistency of this verb through the occurrences of the type. The translation will run in parallel with the same general verb *at sætte* when the argument will be a *glass* [glasset] or *a pot* [gryden], as in (3) and (4), but not when the argument is a *dress* [tøjet] as in (5):

(3)    Marco har sat [stillet] glasset på bordet
(4)    Konen har sat [stillet] gryden over ilden
(5)    Moderen har lagt tøjet på sengen

In the event described by (5) a different verb is strictly required. Danish, which is a language encoding *mood* in its action verbs (Talmy, 1985), applies *at lægge* for the instances where the object lies on its destination, as in Figure 5 B.

Therefore, a new type will arise in the database as a result of this language-specific categorization. The new prototypic scene in Figure 5B will be generated.
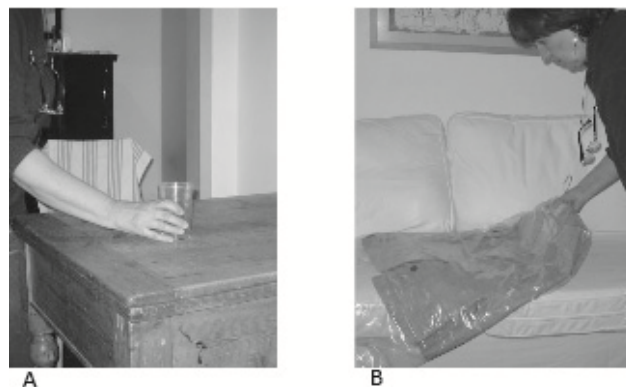


Figure 5. Competence based Extension of Action ontology

In IMAGACT the action ontology will provide equivalences for languages with high global impact but with strong diversity in cultural tradition and linguistic tendencies (Spanish and Chinese Mandarin). After its first delivery the IMAGACT infrastructure will grow freely as a function of its competence-based implementation in an open set of languages. We expect a huge amount of data from this task, which will ground the traditional concept of "Language specific categorization".

## 6. Acknowledgements

## 7. References

British National Corpus, version 3 (BNC XML Edition). 2007. Distributed by Oxford University Computing Services URL: http://www.natcorp.ox.ac.uk/

CLIPS Corpus. URL: http://www.clips.unina.it

C-ORAL-ROM Corpus. Distributed by ELDA http://catalog.elra.info/product_info.php?products_id=757

De Mauro T., Mancini F., Vedovelli M., Voghera M. (1993). *Lessico di frequenza dell'italiano parlato* (LIP). Milano: ETASLIBRI.

Dowty, D. (1979). *Word meaning and Montague grammar*. Dordrecht: Reidel.

Fellbaum, Ch. (1998). *WordNet: An Electronic Lexical Database*. Cambridge (MA): MIT Press.

Givon, T. (1986). *Prototypes: Between Plato and Wittgenstein*. In C. Craig (Ed.), *Noun Classes and Categorization*. Amsterdam: Beniamin, pp. 77-102.

IMAGACT. http://www.imagact.it/

ImageNet http://www.image-net.org/

Jackendoff, R. (1972). *Semantic Interpretation in Generative Grammar*. Cambridge (MA): MIT Press.

Korzen, I. (2005). Endocentric and esocentric languages in translation. *Perspectives: Studies in translatology*, 13(1), pp. 21-37.

LABLITA Corpus of Spontaneous Spoken Italian. URL: http://lablita.dit.unifi.it/corpora/

Majid, A., Boster, J.S., Bowerman, M. (2008). The cross-linguistic categorization of everyday events: A study of cutting and breaking. *Cognition*, 109,(2), pp. 235-250.

Moneglia, M. (2011). Natural Language Ontology of Action. A gap with huge consequences for Natural Language Understanding and Machine Translation. In Z. Vetulani (Ed.) *Proceedings of the 5th Language & Technology Conference: Human Language Technologies as a Challenge for Computer Science and Linguistics*. Poznań: Fundacja Uniwersytetu im. A. Mickiewicza, pp. 95-100.

Moneglia, M., Panunzi, A. (2007). Action Predicates and the Ontology of Action across Spoken Language Corpora. The Basic Issue of the SEMACT Project. In M. Alcántara, T. Declerck (Eds.), *Proceeding of the International Workshop on the Semantic Representation of Spoken Language* (SRSL7). Salamanca: Universidad de Salamanca, pp.51-58.

Moneglia & Panunzi (2011). Specification for the annotation of verb occurrences in the IMAGACT project. Technical Report Draft. http://lablita.dit.unifi.it/projects/IMAGACT/folder.201 0-11-25.7365875310/

Moneglia, M., Monachini, M., Panunzi, A., Frontini, F., Gagliardi, G, Russo I. (2012). Mapping a corpus-induced ontology of action verbs on ItalWordNet. In *Proceedings of the sixth Global WordNet Conference (GWC 2012)*.

Ng, H.T., Lim, C.Y., Foo, S.K. (1999). A Case Study on Inter-Annotator Agreement for Word Sense Disambiguation. In *Proceedings of the ACL SIGLEX Workshop on Standardizing Lexical Resources (SIGLEX99)*. College Park (MD), pp. 9-13.

Rosch, E. (1978). Principles of Categorization. In E. Rosch & B.B. Lloyd (Eds.), *Cognition and Categorization*. Hillsdale: Lawrence Erlbaum Associates, 27–48.

Talmy, L. (1985). Lexicalization patterns: Semantic structure in lexical form. In Shopen T. (Ed.), *Language typology and syntactic description, Vol. III: Grammatical categories and the lexicon.* Cambridge (UK): Cambridge University Press.

Wittgenstein, L. (1953). *Philosophical Investigations*. Oxford: Blackwell.