

# KNOWLEDGE ACQUISITION THROUGH NATURAL LANGUAGE CONVERSATION AND CROWDSOURCING

Luka Bradeško

**Doctoral Dissertation**  
**Jožef Stefan International Postgraduate School**  
**Ljubljana, Slovenia**

**Supervisor:** Doc. Dunja Mladenić, Jožef Stefan Institute, Ljubljana, Slovenia

**Evaluation Board:**

Dr. Michael Witbrock, Chair, IBM, New York, New York

Prof. Erjavec, Member, Jožef Stefan Institute, Ljubljana, Slovenia

Prof. Iztok Savič, Member, Univerza v Novi Gorici, Nova Gorica, Slovenia

MEDNARODNA PODIPLOMSKA ŠOLA JOŽEFA STEFANA  
JOŽEF STEFAN INTERNATIONAL POSTGRADUATE SCHOOL



Luka Bradeško

## KNOWLEDGE ACQUISITION THROUGH NATURAL LANGUAGE CONVERSATION AND CROWDSOURCING

**Doctoral Dissertation**

## PRIDOBIVANJE STRUKTURIRANEGA ZNANJA SKOZI POGOVOR TER S POMOČJO MNOŽIČENJA

**Doktorska disertacija**

**Supervisor:** Doc. Dunja Mladenić

Ljubljana, Slovenia, April 2017



*To the world...*



# Acknowledgments

Thank everyone who contributed to the thesis: - EU Projects - Cyc - Dave - Michael - Vanessa - Dunja - Coworkers





# Abstract

The English abstract should not take up more than one page.



# Povzetek

Povzetek v slovenščini naj ne bo daljši od ene strani.



# Contents

<b>List of Figures</b>	<b>xv</b>
<b>List of Tables</b>	<b>xvii</b>
<b>List of Algorithms</b>	<b>xix</b>
<b>Abbreviations</b>	<b>xxi</b>
<b>Symbols</b>	<b>xxiii</b>
<b>Glossary</b>	<b>xxv</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Scientific Contributions . . . . .	1
1.1.1 Novel Approach Towards Knowledge Acquisition . . . . .	1
1.1.2 Knowledge Acquisition Platform Implementation as Technical Contribution . . . . .	2
1.1.3 A Shift From NL Patterns to Logical Knowledge Representation in Conversational Agents . . . . .	2
1.2 Thesis structure . . . . .	2
<b>2 Background and Related Work</b>	<b>3</b>
2.1 Labour Acquisition . . . . .	4
2.2 Interaction Acquisition . . . . .	4
2.2.1 Games . . . . .	5
2.2.2 Natural Language Conversation . . . . .	5
2.3 Reasoning Acquisition . . . . .	5
2.4 Mining Acquisition . . . . .	5
2.5 Acquisition with the help of existing knowledge . . . . .	5
2.6 Crowdsourcing Acquisition . . . . .	5
2.7 Acquisition of Geospatial Context . . . . .	5
<b>3 Knowledge Acquisition Approach</b>	<b>7</b>
<b>4 Real World Knowledge Acquisition Implementation</b>	<b>9</b>
4.1 Cyc . . . . .	9
<b>5 Evaluation</b>	<b>11</b>
<b>6 Conclusions</b>	<b>13</b>
<b>Appendix A Proofs of Theorems</b>	<b>15</b>
A.1 Proof of the Pythagorean Theorem . . . . .	15

<b>References</b>	<b>17</b>
<b>Bibliography</b>	<b>19</b>
<b>Biography</b>	<b>21</b>

# List of Figures

Figure A.1: Similar triangles used in the proof of the Pythagorean theorem. . . . 15





# List of Tables

Table 2.1: Structured overview of related KA systems . . . . .	4
--	---



# List of Algorithms



# Abbreviations

CC ... Curious Cat (a name of the knowledge acquisition application and platform that is a side result of this thesis)  
CYC ... An AI system (Inference Engine and Ontology), developed by Cycorp Inc.  
CycKB.. Cyc Knowledge Base (Ontology part of Cyc system)  
CycL ... Cyc Lanugage  
JSI ... Jožef Stefan Institute  
KA ... Knowledge Acquisition  
KDML.. Knowledge SDatabase Mark-up Language  
NL ... Natural Language



# Symbols

$j^*$  ... black-body irradiance

$\sigma$  ... Stefan's (or Stefan-Boltzmann) constant





# Glossary

Glossary of terms, dada, bada



# Chapter 1

## Introduction

An intelligent being or machine solving any kind of a problem needs knowledge to which it can apply its intelligence while coming up with an appropriate solution. This is especially true for the knowledge-driven AI systems which constitute a significant fraction of general AI research. For these applications, getting and formalizing the right amount of knowledge is crucial. This knowledge is acquired by some sort of Knowledge Acquisition (KA) process, which can be manual, automatic or semi-automatic. Knowledge acquisition, using an appropriate representation and subsequent knowledge maintenance are two of the fundamental and as-yet unsolved challenges of AI. Knowledge is still expensive to retrieve and to maintain. This is becoming increasingly obvious, with the rise of chat-bots and other conversational agents and AI assistants. The most developed of these (Siri, Cortana, Google Now, Alexa), are backed by huge financial support from their producing companies, and the lesser-known ones still result from 7 or more person-years of effort by individuals

**Finish**

Knowledge acquisition and subsequent knowledge maintenance, are two of the fundamental and as-yet not-completely-solved challenges of Artificial Intelligence (AI).

### 1.1 Scientific Contributions

This section gives an overview of scientific and other contributions of this thesis to the knowledge acquisition approaches.

#### 1.1.1 Novel Approach Towards Knowledge Acquisition

Traditionally KA (knowledge acquisition) approach focuses on one type of acquisition process, which can be either Labor, Interaction, Mining or Reasoning(Zang, Cao, Cao, Wu, & CAO, 2013). In this thesis we propose a novel, previously untried approach that intervenes all aforementioned types with current user context and crowdsourcing into a coherent, collaborative and autonomous KA system. It uses existing knowledge and user context, to automatically deduce and detect missing or unconfirmed knowledge(reasoning) and uses this info to generate crowdsourcing tasks for the right audience at the right time(labor). These tasks are presented to users in natural language (NL) as part of the contextual conversation (interaction) and the answers parsed (mining) and placed into the KB after consistency checks(reasoning). The approach contribution can be summed up as a) definition of the framework for autonomous and collaborative knowledge acquisition with the help of contextual knowledge (chapter X), and b) demonstrate and evaluate the contributions of contextual knowledge and approach in general chapter X.

### 1.1.2 Knowledge Acquisition Platform Implementation as Technical Contribution

Implementation of the KA framework as a working real-world prototype which shows the feasibility of the approach and a way to connect many independent and complex sub-systems. Sensor data, natural language, inference engine, huge pre-existing knowledge base (Cyc)[CycRef](#), textual patterns and crowdsourcing mechanisms are connected and interlinked into a coherent interactive application ([Chapter X](#)).

### 1.1.3 A Shift From NL Patterns to Logical Knowledge Representation in Conversational Agents

Besides the main contributions presented above, one aspect of the approach introduces a shift in the way how conversational agents are being developed. Normally the approach is to use textual patterns and corresponding textual responses, sometimes based on some variables, and thus encode the rules for conversation. As a consequence of natural language interaction, the proposed KA framework is in some sense a conversational agent which is driven by the knowledge and inference rules and uses patterns only for conversion from NL to logic. This shows promise as an alternative approach to building non scripted conversational engines ([Chapter X](#)).

## 1.2 Thesis structure

The rest of the thesis is structured into chapters covering specific topics. [Chapter X](#) introduces

## Chapter 2

# Background and Related Work

In this chapter we will give an overview of approaches and related works on broader knowledge acquisition research field, information extraction, crowdsourcing and geo-spatial context mining.

Knowledge Acquisition has been addressed from different perspectives by many researchers in Artificial Intelligence over decades, starting already in 1970 as a sub-discipline of AI research (Feigenbaum-economicPhd), and since then resulting in a big number of types and implementations of approaches and technologies/algorithms. In more recent survey of KA approaches (Zang et al., 2013), authors categorize all of the KA approaches into four main groups, regarding the source of the data and the way knowledge is acquired:

- *Labour Acquisition.* This approach uses human minds as the knowledge source. This usually involves human (expert) ontologists manually entering and encoding the knowledge.
- *Interaction Acquisition.* As in Labour Acquisition, the source of the knowledge is coming from humans, but in this case the KA is wrapped in a facilitated interaction with the system, and is sometimes implicit rather than explicit.
- *Reasoning Acquisition.* In this approach, new knowledge is automatically inferred from the existing knowledge using logical rules and machine inference.
- *Mining Acquisition.* In this approach, the knowledge is extracted from some large textual corpus or corpora.

We believe this categorization most accurately reflects the current state of machine (computer) based knowledge acquisition, and we decided to use the same classification when structuring our related work, focusing more on closely related approaches and extending where necessary. According to this classification, our work presented in this thesis, fits into a hybrid approach combining all four groups, with main focus on interaction and reasoning. We address the problem by combining the labour and interaction acquisition (users answering questions as part of NL interaction aimed at some higher level goal, such as helping the user with various tasks), adding unique features of using user context and existing knowledge in combination with reasoning to produce a practically unlimited number of potential interaction acquisition tasks, going into the field of crowd-sourcing by sending these generated tasks to many users simultaneously.

Previous works that can compares to our solution is divided into the systems that exploit existing knowledge (generated anew during acquisition or pre-existing from before in other sources) (Forbus2007; Kvo2010; Mitchel2015; Singh et al., 2002; Witbrock et al., 2003; Sharma & Forbus, 2010), reasoning (Witbrock et al., 2003; Speer, 2007; Speer,

Fix this, refer to chapters i to specific w

Lieberman, & Havasi, 2008; Kuo & Hsu, 2010), crowdsourcing (**Singh2002**; **Pedro2012a**; Speer et al., 2009; Kuo & Hsu, 2010; Pedro, Appel, & Jr, 2013), acquisition through interaction (**Pedro2012**; Speer et al., 2009; Pedro et al., 2013), acquisition through labour(**add, probably rather refer to subsections**) () and natural language conversation(**Pedro2012**; Speer, 2007; Speer et al., 2009; Witbrock et al., 2003; Kuo & Hsu, 2010).

**Test referencing table** (see Table 2.1).

Table 2.1: Structured overview of related KA systems

System	Ref.	Cat.	Source	Repr.	PK	CS	C
Cyc project (Cycorp)	(Lenat, 1995)	Labour	KE	CycL	/	/	/
ThoughtTrasure(Signiform)	( <b>Mueller2003</b> )	Labour	KE	LAGS	/	/	/
HowNet (Keen.)	(Lenat, 1995)	Labour	KE	KDML	/	/	/

## 2.1 Labour Acquisition

This category consists of KA approaches which rely on explicit human work to collect the knowledge. A number of expert (or also untrained) ontologists or knowledge engineers is employed to codify the knowledge by hand into the given knowledge representation (formal language).

*Cyc.* The most famous and also most comprehensive and expensive knowledge acquired this way, is Cyc KB, which is part of Cyc AI system (Lenat, 1995). It started in 1984 as a research project, with a premise that in order to be able to think like humans do, the computer needs to have knowledge about the world and the language like humans do, and there is no other way than to teach them, one concept at a time, by hand. Since 1994, the project continued through Cycorp Inc. company, which is still continuing the effort. Through the years Cyc Inc. employed computer scientists, knowledge engineers, philosophers, ontologists, linguists and domain experts, to codify the knowledge in the formal higher order logic language CycL (**Cyc Language**). As of 2006 (**Matuszek2006**), the effort of making Cyc was 900 non-crowdsourced human years which resulted in 7 million assertions connecting 500,000 terms and 17,000 predicates/relations (Zang et al., 2013), structured into consistent sub-theories (Microtheories) and connected to the Cyc Inference engine and Natural Language generation. Since the implementation of our approach is based on Cyc, we give a more detailed description of the KB and its connected systems in section 4.1 on page 9. Cyc Project is still work in progress and continues to live and expand through various research and commercial projects.

*ThoughtTreasure.* Approximately at the same time as Cyc Inc. company was formed (1994), Eric Mueller started to work on a similar system, which was inspired by Cyc and is similar in having a combination of common sense knowledge concepts connected to their natural language presentations. The main differentiator from Cyc is, that it tries to use simpler representation compared to first-order logic as is used in Cyc. Additionally, some parts of ThoughtTreasure knowledge can be presented also with finite automata, grids and scripts (**Mueller1999**; **Mueller2003**). In 2003 the knowledge of this system consisted of 25,000 concepts and 50,000 assertions. ThoughtTreasure was not so successful as Cyc and ceased all developments in 2000 and was open-sourced on Github in 2015.

*ThoughtTreasure.* Dada

## **2.2 Interaction Acquisition**

adada

### **2.2.1 Games**

adada

### **2.2.2 Natural Language Conversation**

dada

## **2.3 Reasoning Acquisition**

adad dada

## **2.4 Mining Acquisition**

adad

## **2.5 Acquisition with the help of existing knowledge**

adad

## **2.6 Crowdsourcing Acquisition**

adad

## **2.7 Acquisition of Geospatial Context**

adad





## Chapter 3

# Knowledge Acquisition Approach



## Chapter 4

# Real World Knowledge Acquisition Implementation

### 4.1 Cyc

TBW



## Chapter 5

# Evaluation

TBW



## Chapter 6

# Conclusions

We came to the following conclusions . . .





## Appendix A

# Proofs of Theorems

### A.1 Proof of the Pythagorean Theorem

Let us prove the Pythagorean Theorem from page ??.

*Proof.* This proof is based on the proportionality of the sides of two similar triangles, that is, upon the fact that the ratio of any two corresponding sides of similar triangles is the same regardless of the size of the triangles.

Let  $ABC$  represent a right triangle, with the right angle located at  $C$ , as shown in Figure A.1. We draw the altitude from point  $C$ , and call  $H$  its intersection with the hypotenuse  $AB$ . Point  $H$  divides the length of the hypotenuse into two parts.

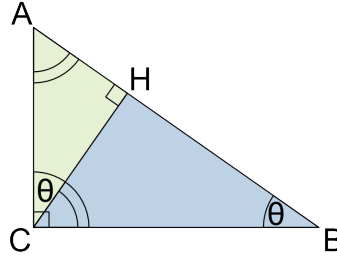


Figure A.1: Similar triangles used in the proof of the Pythagorean theorem.

The new triangle  $ACH$  is similar to triangle  $ABC$ , because they both have a right angle (by definition of the altitude), and they share the angle at  $A$ , meaning that the third angle will be the same in both triangles as well, marked as  $\theta$  in Figure A.1. By a similar reasoning, the triangle  $CBH$  is also similar to  $ABC$ .

Similarity of the triangles leads to the equality of ratios of corresponding sides:

$$\frac{BC}{AB} = \frac{BH}{BC} \text{ and } \frac{AC}{AB} = \frac{AH}{AC}. \quad (\text{A.1})$$

The first result equates  $\cos \theta$  and the second result equates  $\sin \theta$ .

These ratios can be written as:

$$BC^2 = AB \times BH \text{ and } AC^2 = AB \times AH. \quad (\text{A.2})$$

Summing these two equalities, we obtain:

$$BC^2 + AC^2 = AB \times BH + AB \times AH = AB \times (AH + BH) = AB^2, \quad (\text{A.3})$$

which, tidying up, is the Pythagorean theorem:

$$BC^2 + AC^2 = AB^2. \tag{A.4}$$

□

# References

- Kuo, Y. & Hsu, J. (2010). Goal-Oriented Knowledge Collection. *AAAI Fall Symposium: Commonsense Knowledge*, 64–69. Retrieved from <http://www.aaai.org/ocs/index.php/FSS/FSS10/paper/viewPDFInterstitial/2278/2605>
- Lenat, D. B. (1995). Cyc: A Large-Scale Investment in Knowledge Infrastructure. *Communications of the ACM*, 38(22).
- Pedro, S. D. S., Appel, A. P., & Jr, E. R. H. (2013). Autonomously reviewing and validating the knowledge base of a never-ending learning system. In *Proceedings of the 22nd ...* (pp. 1195–1203). Retrieved from <http://dl.acm.org/citation.cfm?id=2488149>
- Sharma, A. & Forbus, K. (2010). Graph-Based Reasoning and Reinforcement Learning for Improving Q/A Performance in Large Knowledge-Based Systems. In *2010 aaai fall symposium series* (pp. 96–101). Retrieved from <http://www.aaai.org/ocs/index.php/FSS/FSS10/paper/download/2246/2596>
- Singh, P., Lin, T., Mueller, E., Lim, G., Perkins, T., & Zhu, W. (2002). Open Mind Common Sense: Knowledge acquisition from the general public. *Cooperative Information Systems Oct. 30-Nov. 1 2002*, 1223–1237. doi:10.1007/3-540-36124-3\_77
- Speer, R. (2007). Open mind commons: An inquisitive approach to learning common sense. *Workshop on Common Sense and Intelligent User ...* Retrieved from <http://www.fatih.edu.tr/%7B~%7Dhugur/inquisitive/Open%20Mind%20Commons%20An%20Inquisitive%20Approach%20to.PDF>
- Speer, R., Krishnamurthy, J., Havasi, C., Smith, D., Lieberman, H., & Arnold, K. (2009). An interface for targeted collection of common sense knowledge using a mixture model. *Proceedings of the 14th International Conference on Intelligent User Interfaces*, 137–146. doi:10.1145/1502650.1502672
- Speer, R., Lieberman, H., & Havasi, C. (2008). AnalogySpace : Reducing the Dimensionality of Common Sense Knowledge. *Aaai*, 548–553.
- Witbrock, M., Baxter, D., Curtis, J., Schneider, D., Kahlert, R. C., Miraglia, P., ... Vizedom, A. (2003). An Interactive Dialogue System for Knowledge Acquisition in Cyc. In *Proceedings of the eighteenth international joint conference on artificial intelligence*. Acapulco, Mexico.
- Zang, L.-J., Cao, C., Cao, Y.-N., Wu, Y.-M., & CAO, C.-G. (2013). A Survey of Commonsense Knowledge Acquisition. *Journal of Computer Science and Technology*, 28(4), 689–719. doi:10.1007/s11390-013-1369-6



# Bibliography

## **Publications Related to the Thesis**

All publications related to the thesis should be referenced in the text.

## **Other Publications (optional)**

...



# Biography

The author of this thesis . . .

