# Short overview on newest trends and development directions of Natural Language Processing in Haskell

*Alina Parascan*

*Alexandru Velicea*

## 1. What is NLP ( Natural Language Processing) ?

As defined by Wikipedia " **Natural language processing** (**NLP**) is a field of computer science, artificial intelligence and linguistics concerned with the interactions between computers and human (natural) languages. As such, NLP is related to the area of human–computer interaction. Many challenges in NLP involve natural language understanding, that is, enabling computers to derive meaning from human or natural language input."

A short insight of NLP starts with the very beginning of computers. In 1950 Alan Turing published an article called "Computing Machinery and Intelligence" which proposed what is now known as the "Turing Test" as a criterion of intelligence. In reality, he made a breakthrough in this field much earlier, during the second World War when he conceived and built a machine that helped decrypt the "Enigma" code used by the Germans to pass secret military messages. This was one of the first computers and of

course the first instance in which a computer was used for Natural Language Processing purposes. Its importance is proved by the fact that decrypting the "Enigma" code shortened the war with 2 to 4 years and was estimated to saving the lives of over 4 million people [1]

## 2. Uses of NLP libraries

There are numerous uses of NLP libraries: keyword generation, language identification, full-text search, artificial intelligence, auto completion, translation, topic detection and clustering, content summarization, parsing human-written dates and locations, etc.

We will mention some of the most notable Major Tasks in NLP:

**Automatic summarization**, which produces a short and readable summary of a single or multiple documents, preserving the relevant information.[2]

**Machine translation,** which means to automatically translate a text from one human language to another. Putting aside Google Translate, some notable software in this field are Babylon, Power Translator, Prompt[3]

---

[1] Jack Copeland, (18 June 2012). "Alan Turing: The codebreaker who saved 'millions of lives'". BBC News Technology. Retrieved 26 October 2014

[2] Dipanjan Das, André F.T. Martins, "A Survey on Automatic Text Summarization", Carnegie Mellon University, November 2007

[3] according to http://translation-software-review.toptenreviews.com/

**Named entity recognition (NER),** which is the task of correctly identification of proper names of people, organizations, locations, or other entities, determining in the same time the type of each such name is (e.g. person, location, organization). [4]

**Natural language generation,** which convert information from computer databases into readable human language, thus deliberately constructing a natural language text in order to meet specified communicative goals. [5]

**Part-of-speech tagging,** which is the process of classifying words into their **parts of speech** and labeling them accordingly. [6] As we know, many words, especially common ones, can serve as multiple parts of speech. For example, "color" can be a noun ("the color of the sky is blue") or verb ("to color a drawing"); "stone" can be a noun, verb (to stone) or adjective (a stone bridge); and "out" can be any of at least five different parts of speech. Some languages have more such ambiguity than others. English language has little inflectional morphology, therefore is particularly prone to such ambiguity. Chinese is prone to such ambiguity because it is a tonal language during verbalization. Such inflection is not readily conveyed via the entities employed within the orthography to convey intended meaning.

**Parsing -** determines the parse tree (grammatical analysis) of a given sentence.

**Question answering –** determines the answer of a human-language question. Point wise questions have a specific correct answer (such as "What is the capital of Romania?"), but

---

[4] Benjamin Bengfort, "An Introduction to Named Entity Recognition in Natural Language Processing - Part 1", 17 April 2013, www.datacommunnitydc.org

[5] Claire Gardent, "Natural Language Generation", CNRS/LORIA, Nancy, http://www.loria.fr/~gardent/teaching/nlg-bkk06.pdf

[6] Steven Bird, Ewan Klein, and Edward Loper, "Natural Language Processing", O'Reilly Media, 2009

sometimes the question considered is an open-ended one (such as "What is the definition of happiness?").

**Sentence breaking (also known as sentence boundary disambiguation) -** finds the sentence boundaries in a given piece of text. Sentence boundaries are often marked by period, commas or other punctuation marks, but these same characters can serve other purposes (e.g. marking abbreviations).

**Sentiment analysis - e**xtract subjective information usually from a set of documents, often using online reviews to determine the contextual polarity of specific objects in a text. In other words, it determines whether a piece of writing is positive, negative or neutral.[7] It is especially useful for identifying trends of public opinion in the social media, for the purpose of marketing.

Among other notable and recent uses of NLP libraries we would like to mention the work of Hakan Burden, Rogardt Heldal and Peter Ljungl  from University of Gothenburg, Sweden, "*Opportunities for Agile Documentation Using Natural Language Generation*", that discuss the opportunities for Agile using Natural Language Generation in order to prevent the general miscommunication that occurs between the client and the actual developers and to realize a more accurate approach to requirements specifications. They propose Natural Language Generation as a suitable framework for automatically retrieving relevant information from the source code of certain modules within the application that will impact future tasks. This information will be then displayed in an appropriate text format (i.e. adding the necessary HTML tags or in a LATEX format for the presentation). Their conclusion was:

> "Using the code for generating documents that describe different system properties enables a single source of information in the development process, aligning documentation with implementation. Different views could then be generated from the code during implementation, with content and structure

---

[7] www.semantria.com "What is Sentiment Analysis?"

depending on who the consumer of the document and what the purpose of the text is. This could be done together with the relevant stakeholders, enabling concise documents that are easy to keep consistent with the implementation. The opportunities for collaborations between industry and academia in text generation as a means for Agile documentation are abundant! "

These are a few of the most important and recent uses of NLP, most of which have already been implemented and used, but considering the continuous growth of the volume of information present on the internet and the increasing importance of social media in day to day life, we are convinced that NLP will be used in effective ways to sort, search and interpret vast volumes of data/information.

## 3. NLP and Haskell

Currently there are over 70 NLP libraries available on Hackage and this speaks volumes in regards to the direction and evolution of functional programming in general, Haskell in particular and of course the increasing attention that Natural Language Processing is getting recently. Even though the basic algorithms and principles used for NLP were developed a while ago, some of them have evolved and have been improved through the years and are still used. One of those is the Porter Stemmer Algorithm that was originally written in 1979 at Cambridge University. The purpose of this algorithm is to reduce words to common roots; for example it should identify the string "cats" or "catlike" as based on the root "cat" and "stemmer", "stemming", "stemmed" as based on "stem". This is particularly useful when building a search engine, or a language translation tool, or for any Natural Language Processing purpose as reducing words to their root makes processing them easier.

The Haskell version of this algorithm was written by Dimitry       and it is still used present day.

Another important tool in the Natural Language Processing world is LOLITA. The name is an acronym for "**L**arge-scale, **O**bject-based, **L**inguistic **I**nteractor, **T**ranslator and **A**nalyzer". It is an NLP system that provides core capabilities such as:

1. " conversion of English language text to an internal (semantic network based) representation of the text's meaning
2. inferences in the semantic network
3. conversion of portions of the semantic network to English "[8]

"LOLITA was an early example of a substantial application written in a functional language: it consisted of around 50,000 lines of Haskell, with around 6000 lines of C. It is also a complex and demanding application, in which many aspects of Haskell were invaluable in development.

LOLITA was designed to handle unrestricted text, so that ambiguity at various levels was unavoidable and significant. Laziness was essential in handling the explosion of syntactic ambiguity resulting from a large grammar, and it was much used with semantic ambiguity too." [9]

---

[8] http://homepages.inf.ed.ac.uk/wadler/realworld/natlangproc.html

[9] http://en.wikipedia.org/wiki/LOLITA

## 4. WordNet

Of course, in order to utilize NLP capabilities to the fullest one needs different helping tools; for example a database that contains words of a specific language in order to be able to process text and information. One of the most significant databases of this sort is WordNet.  It is a lexical database for the English language and it was originally developed in  1985 in the Cognitive Science Laboratory at Princeton University under the direction of psychology professor George Armitage Miller.

WordNet groups English words into sets of synonyms (also called synsets) and it provides some short definitions and also records relations among these synonyms sets and their members. The database includes lexical categories such as: nouns, verbs, adjectives and adverbs but it disregards prepositions and other function words. It is primarily used in artificial intelligence and text analysis and the most recent version contains 155,287 words that are organized in 117,659 synsets and thus holding a total of 206,941 word-sentence pairs.

WordNet has been used for a number of different purposes such as automatic text classification, text summarization, machine translation, and even automatic crosswords puzzle generation.

Even though it was first developed in 1985 it is still currently used and interfaces in different programming languages have been developed for WordNet throughout the years. For example, one of the first Haskell interfaces for WordNet has been written in 2003 by Hal Daumé and after that two more have been written; one in 2008: it is available on Hackage at *http://hackage.haskell.org/package/WordNet-0.1.2* and the most recent Haskell Interface for this database was written in 2013 (also available on Hackage - *https://hackage.haskell.org/package/WordNet*).

This proves the increasing interest in recent years for text processing and Natural Language Processing and also that functional programming is one of the main programming styles that are and will be used for these types of endeavors as more and more libraries and tools are developed in order to support the growing need and interest for such subjects.

## 5. Practical examples of NLP libraries and applications written in Haskell

There are numerous libraries for Natural Language Processing available on Hackage and quite a few recent applications. We have chosen to start with a short presentation of the Chatter NLP Library. [10]

It is of course open source and available for download and usage and it comprises NLP algorithms that are developed for the purpose of tagging parts of speech according to the Brown Corpus. [11]

After installing it with the aid of `cabal install` command the user can input a text and the tagger function will output a list of every word that was previously input and tag each word into a category according with Brown Corpus list. For example the text "The best jokes have no punchline." will be analyzed as follows:

```
[TaggedSent [
POS {posTag = DT, posToken = Token "The"},      //DT - means singular
                                                determiner/quantifier (this, that)
POS {posTag = JJS, posToken = Token "best"},     //JJS - semantically superlative
                                                adjective (chief, top)
POS {posTag = NNS, posToken = Token "jokes"},    //NNS - plural noun
```

[10]  http://hackage.haskell.org/package/chatter

[11]  http://en.wikipedia.org/wiki/Brown_Corpus#Part-of-speech_tags_used.

POS {posTag = VBP, posToken = Token "have"},
POS {posTag = DT, posToken = Token "no"},
POS {posTag = NN, posToken = Token "punchline"},
POS {posTag = Term, posToken = Token "."}]]

We will describe a tutorial on how to install and utilize chatter:

**Step 1 : Install cabal**

Open a cmd



and type:

>cabal install

>cabal update

>cabal install cabal-install

Step two: Install chatter from Hackage:



Step three : Create a Haskell file named Main.hs and paste the following code :

1. Import the parts of speech library and the pack function:

```
import NLP.POS

import Data.Text (pack)
```

2. Obtain the default tagger provided by the library:

```
main = do

tagger <- defaultTagger
```

3. Feed the tag function a tagger and a text to see the corresponding parts of speech per each word:

```
let text = pack "The best jokes have no punchline."

print $ tag tagger text
```

4. The output will be an association list of the word to its part of speech:

```
[[ ("The", Tag "at"),

   ("best", Tag "jjt"),

   ("jokes", Tag "nns"),

   ("have", Tag "hv"),

  ("no", Tag "at"),

   ("punchline",Tag "nn"),

   (".",Tag ".") ]]
```

# Tutorial 2 KNLP Web application by Mark Watson

\>cabal install yesod



\>cabal install hsparql

>cabal install json

\>cabal install

encoding



cabal install encoding

cabal install split

>cabal install split



```
[10 of 42] Compiling Text.XML.HaXml.Parse    ( src\Text\XML\HaXml\Parse.hs, dist\bu
ild\Text\XML\HaXml\Parse.o )
[11 of 42] Compiling Text.XML.HaXml.Pretty   ( src\Text\XML\HaXml\Pretty.hs, dist\
build\Text\XML\HaXml\Pretty.o )
[12 of 42] Compiling Text.XML.HaXml.Html.Generate ( src\Text\XML\HaXml\Html\Gene
rate.hs, dist\build\Text\XML\HaXml\Html\Generate.o )
[13 of 42] Compiling Text.XML.HaXml.Html.Parse  ( src\Text\XML\HaXml\Html\Parse.h
s, dist\build\Text\XML\HaXml\Html\Parse.o )
[14 of 42] Compiling Text.XML.HaXml.Wrappers  ( src\Text\XML\HaXml\Wrappers.hs, d
ist\build\Text\XML\HaXml\Wrappers.o )

src\Text\XML\HaXml\Wrappers.hs:34:47:
    lexical error in string/character literal at character '\r'
cabal: Error: some packages failed to install:
HaXml-1.24.1 failed during the building phase. The exception was:
ExitFailure 1
encoding-0.7.0.2 depends on HaXml-1.24.1 which failed to install.

C:\Users\Alina>cabal install split
Resolving dependencies...
All the requested packages are already installed:
split-0.2.2
Use --reinstall if you want to reinstall anyway.

C:\Users\Alina>
```
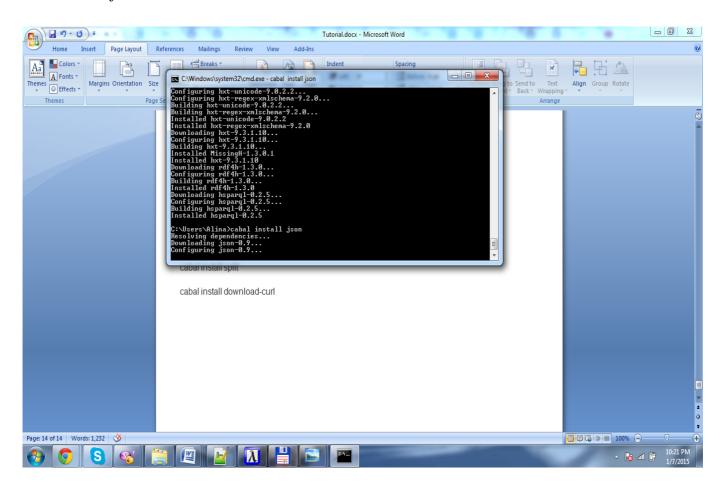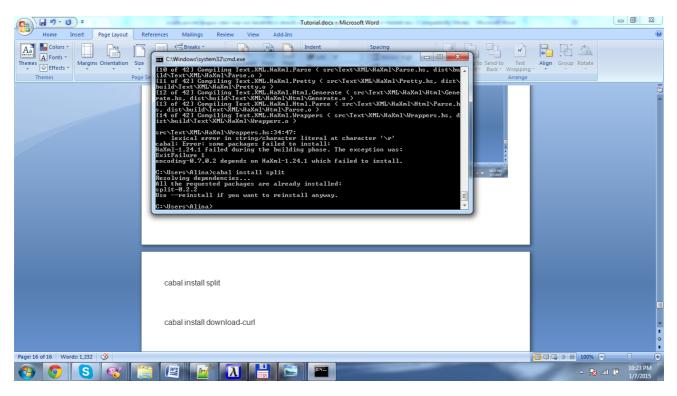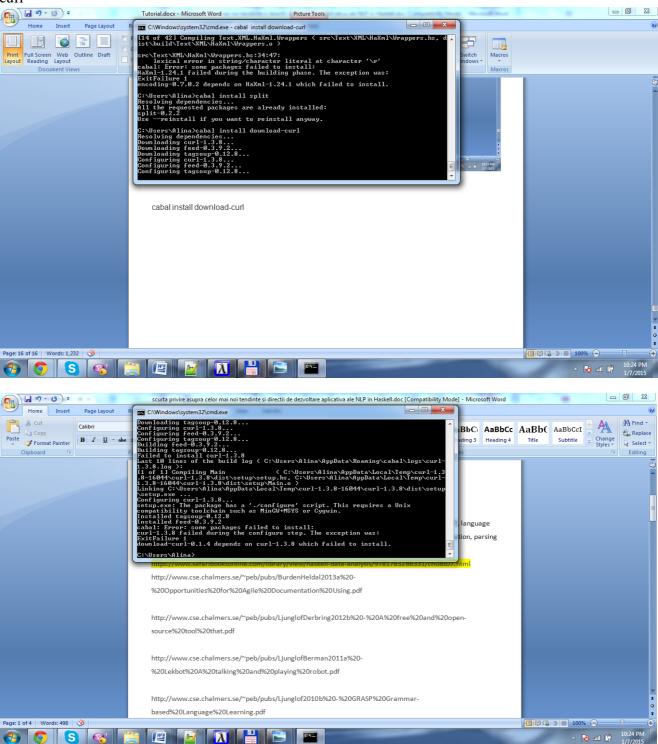
cabal install split

cabal install download-curl

\>cabal install download-

curl



cabal install download-curl



https://www.safaribooksonline.com/library/view/haskell-data-analysis/9781783286331/ch08s07.html

http://www.cse.chalmers.se/~peb/pubs/BurdenHeldal2013a%20-
%20Opportunities%20for%20Agile%20Documentation%20Using.pdf

http://www.cse.chalmers.se/~peb/pubs/LjunglofDerbring2012b%20-%20A%20free%20and%20open-
source%20tool%20that.pdf

http://www.cse.chalmers.se/~peb/pubs/LjunglofBerman2011a%20-
%20Lekbot%20A%20talking%20and%20playing%20robot.pdf

http://www.cse.chalmers.se/~peb/pubs/Ljunglof2010b%20-%20GRASP%20Grammar-
based%20Language%20Learning.pdf
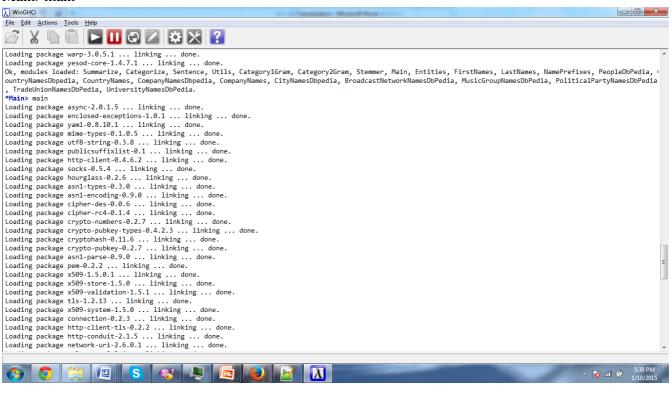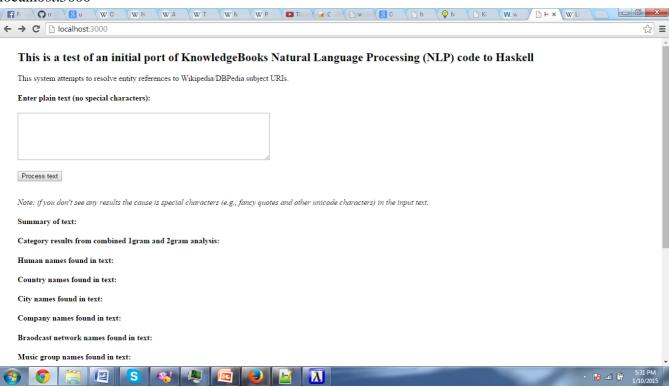
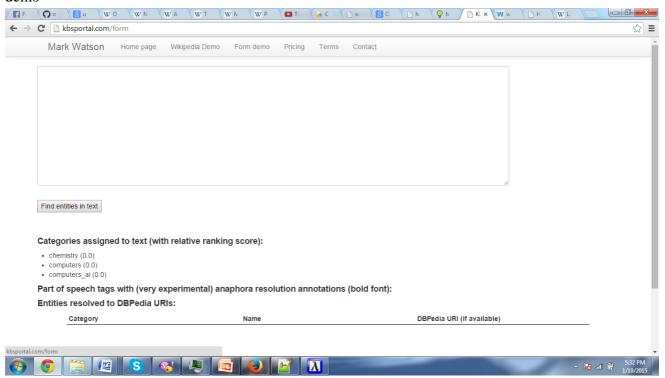Prelude>:load WebApp.hs



**start yesod server**

Main>main



localhost:3000

author's website:

http://markwatson.com/

demo



documentation for the app:

http://kbsportal.com/demo

**Bibliography:**

- *Natural Language Processing for the Working Programmer*
  2010, 2011 Daniël de Kok, Harm Brouwer

- Jack Copeland, (18 June 2012). "Alan Turing: The codebreaker who saved 'millions of lives'". BBC News Technology. Retrieved 26 October 2014

- Dipanjan Das, André F.T. Martins, "A Survey on Automatic Text Summarization", Carnegie Mellon University, November 2007

- http://translation-software-review.toptenreviews.com/

- Benjamin Bengfort, "An Introduction to Named Entity Recognition in Natural Language Processing - Part 1", 17 April 2013, www.datacommunnitydc.org

- Claire Gardent, "*Natural Language Generation*", CNRS/LORIA, Nancy, http://www.loria.fr/~gardent/teaching/nlg-bkk06.pdf

  Steven Bird, Ewan Klein, and Edward Loper, "*Natural Language Processing*", O'Reilly Media, 2009

- www.semantria.com "*What is Sentiment Analysis?*"

- http://homepages.inf.ed.ac.uk/wadler/realworld/natlangproc.html

- http://en.wikipedia.org/wiki/LOLITA

-  http://hackage.haskell.org/package/chatter

- http://en.wikipedia.org/wiki/Brown_Corpus#Part-of-speech_tags_used.

- *Grammar-Based Automated Music Composition in Haskell*
  Donya Quick and Paul Hudak, Yale University, 2013

- *The Haskell School of Music — From Signals to Symphonies —*
  Paul Hudak, Yale University, Department of Computer Science,  January 2014

- http://www.ibm.com/developerworks/library/j-ft20/

- http://www.smashingmagazine.com/2014/07/02/dont-be-scared-of-functional-programming/

- http://www.kurzweilai.net/psychopathic-killers-computerized-text-analysis-uncovers-the-word-patterns-of-a-predator
- http://haskell.cs.yale.edu/