

# FactBook Semantic Technologies<sup>1</sup>

## Introduction

**FactBook** uses semantic technologies to understand users' current interests and collect appropriate facts on their behalf. The underlying deep structure learning algorithms presumably mimic the way children acquire their native language.<sup>2</sup> This Whitepaper presents the basics of our approach.

## Overview

Machine learning of natural languages was for a long time represented mostly by supervised learning. Maybe due to the vast bulk of knowledge already accumulated in linguistics. Learning made use of large linguistic corpuses, trying to find formal rules reproducing correct parsing provided by skilled linguists and using existing linguistic knowledge.

The resulting parsers were overcomplicated and could parse only several dozen sentences per second (~ 1 kB/s). Which was four orders of magnitude less than state-of-the-art search engines indexing speed (10 MB/s). Thus, practical semantic indexing at scale was prohibited. Situation changed recently, with parsing speed up to 100 and even 1000 sentences<sup>3,4</sup> per second.

However, most parsers still utilize supervised learning, with corpuses available only for some languages in a certain domains. Most scientific and blogosphere jargons are out of reach of such approach.<sup>5</sup> Thus, deep unsupervised learning without linguistic priors attract more and more attention in recent years.<sup>6</sup>

We seek to learn language from scratch in unsupervised fashion, without dictionaries or any other prior linguistic knowledge. The resulting "semantic processor" is capable of parsing syntactical and semantic structures of sentences as fast as 2000 sentences per second (~ 300 kB/s).<sup>7</sup>

Contrary to the mainstream approach based on distributed representations in deep and recurrent networks<sup>8,9</sup> our model uses extremely sparse representations, namely recurrent correlation analysis. Linguistic rules form a large (open) set of mutually competing linguistic patterns, represented as binary trees, based on correlation statistics. It is fully applicable to any language and any domain.

## Language model

Supervised approach describes language, using precompiled dictionaries and linguistic corpuses. We aim to extract the "language model" from scratch, making use of the abundance of raw texts.

---

<sup>1</sup> Sergey Shumsky, serge.shumsky@gmail.com

<sup>2</sup> S.A.Shumsky, "Brain and Language: How we Understand Speech", in "Approaches to Cognition Modeling", V.G.Red'ko (Editor), 2014, ISBN 978-5-9710-1049-4

<sup>3</sup> "Large-Scale Syntactic Processing: Parsing the Web", Final Report of the 2009 JHU CLSP Workshop

<sup>4</sup> Alexander Volokh and Günter Neumann (2012) Task-oriented Dependency Parsing Evaluation Methodology, IEEE 13th International Conference on Information Reuse and Integration, IEEE Systems, Man, and Cybernetics Society (SMC), 2012

<sup>5</sup> Say in biology various genes and proteins have a quite different biological "meaning".

<sup>6</sup> Schmidhuber J. "Deep Learning in Neural Networks: An Overview". arXiv:1404.7828 v3 [cs.NE], 2014

<sup>7</sup> Semantic indexing of English Wikipedia takes 15 hours on a single server.

<sup>8</sup> R.Collobert, et.al. "Natural Language Processing (Almost) from Scratch". J. of Machine Learning Research 12 (2011) 2493-2537

<sup>9</sup> T.Mikolov, I.Sutskever, K.Chen, G.Corrado, J.Dean "Distributed Representations of Words and Phrases and their Compositionality" (NIPS 2013)

Namely, we seek to:

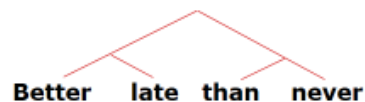
- Understanding **meanings**. Words and phrases have close meanings, if they are used in a similar way (e.g. names, surnames, dates, cities, countries, etc.). To that end, we exploit statistics of usage of linguistic variables in that or another context.
- Understanding **relations** - binding of words and phrases in a sentence. (E.g. name and surname, subject and its action, etc.). Each successive binding in our model increases the probability of similar bindings in the future.

### Key technology: semantic indexing of phrases

#### □ Understanding *meaning*

Who: **George** ~ **Bill** ~ **Donald**  
When: **yesterday** ~ **recently** ~ **April 18**  
What: **said** ~ **reported** ~ **communicated**

#### □ Understanding *relations*



### Golem's language model

Children can easily acquire any language, since all human languages have adapted themselves to our mental capabilities. Thus, it is natural to construct language model based on what we already know about the algorithms of human brain.<sup>10</sup>

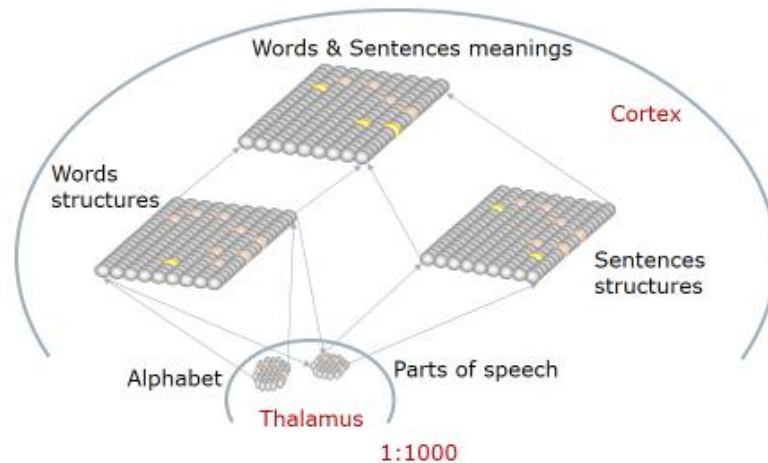
In particular, we know that all our experience is stored in the neocortex in a set of interconnecting 2D self-organizing feature maps. We built our model of human "language faculty" having in mind recursive self-organizing maps, as possible correlation analysis machinery in the neocortex. We also supposed different parts of neocortex to process morphological, syntactical and semantic information.

The resulting language model "Golem" is to our best knowledge the first biologically viable model of unsupervised language learning "from scratch". It has several modules, using the same set of basic learning algorithms in a hierarchical manner - to find out and encode typical linguistic patterns.

---

<sup>10</sup> S.A.Shumsky, "Reverse Engineering the Brain Architecture: Role and Co-operation of the Main Subsystems", in: Lectures on Neuroinformatics, pp.13-45. (Neuroinformatics-2015), Moscow, 2015.

## Our approach: Brain-inspired model of «language faculty»



### Golem's morphology module

Morphology module learns to merge letters iteratively into morphemes and words. The resulting morphological structure of the words is used by syntactic and semantic modules, e.g. for stemming and part of speech analysis.

Clear

Parse

СТЕММИНГ	Словоформы
<i>(morpholog</i>	<i>morphology 1.0000</i>
	<i>morphological 0.5631</i>
	<i>morphologically 0.1395</i>
	<i>morphologies 0.0370</i>
	<i>morphologic 0.0242</i>
	<i>morphological 0.0206</i>
	<i>morphologie 0.0122</i>
	<i>morphologist 0.0054</i>
	<i>morphologically 0.0022</i>
	<i>morphologists 0.0017</i>
	<i>morphologische 0.0017</i>
	<i>morphology_ 0.0012</i>
<i>(morphology)</i>	
<i>(morph ology)</i>	
<i>(mor ph olog y)</i>	
<i>(m or p h ol og y )</i>	
<i>( m o r p h o l o g y )</i>	

The above figure illustrates the merging procedure, stemming and various word forms of the word ‘*morphology*’, learned in unsupervised manner.

## Golem’s syntactic module

In a similar fashion, the syntactic module learns various patterns in parts of speech usage in sentences. The figures below illustrate the iterative procedure of corresponding bottom up sentences parsing.

*[ (deep)(learning)(is)(part)(of) ][ (a)(broader)(family)(of)(machine)(learning)(methods) ]*  
*[ (deep)(learning) ][ (is)(part)(of) ][ (a)(broader)(family) ][ (of)(machine)(learning)(methods) ]*  
*[ (deep) ][ (learning) ][ (is) ][ (part)(of) ][ (a)(broader) ][ (family) ][ (of)(machine)(learning) ][ (methods) ]*

*[ (recent)(research)(has)(increasingly)(focused)(on) ][ (unsupervised)(learning)(algorithms) ]*  
*[ (recent)(research) ][ (has)(increasingly)(focused)(on) ][ (unsupervised)(learning) ][ (algorithms) ]*  
*[ (recent) ][ (research) ][ (has)(increasingly) ][ (focused)(on) ][ (unsupervised) ][ (learning) ][ (algorithms) ]*

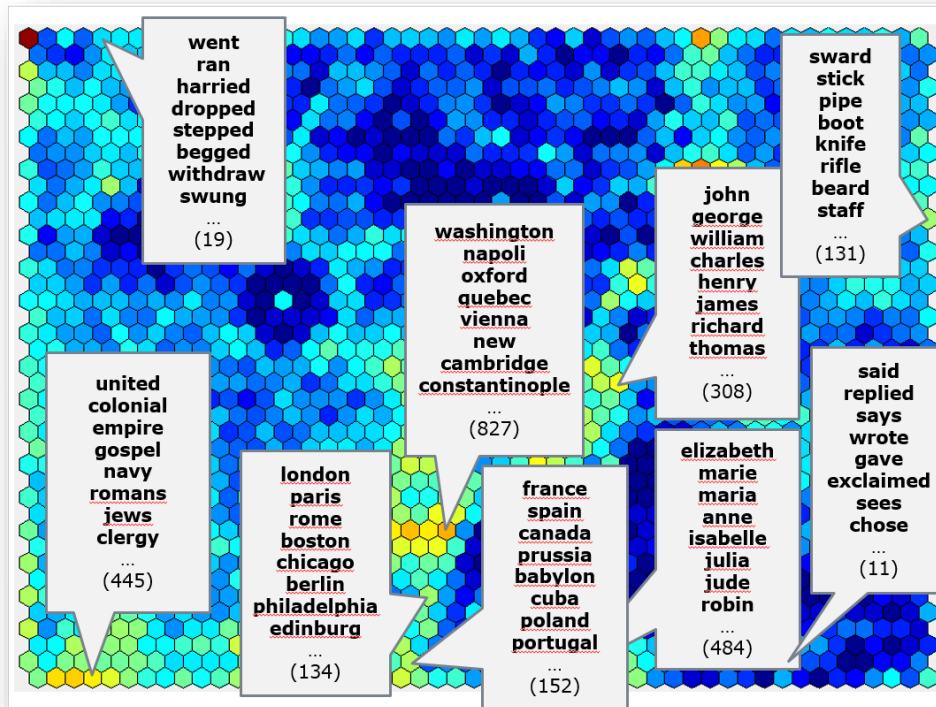
*[ (a)(recurrent)(neural)(network) ][ (is)(a)(class)(of)(artificial)(neural)(network) ]*  
*[ (a)(recurrent) ][ (neural)(network) ][ (is)(a)(class)(of) ][ (artificial)(neural)(network) ]*  
*[ (a) ][ (recurrent) ][ (neural) ][ (network) ][ (is)(a) ][ (class)(of) ][ (artificial) ][ (neural)(network) ]*

## Golem’s semantic module

Semantic module encodes the meanings of the words, phrases and sentences. The starting point is the semantic space for words, illustrated in the next figure. Each «semantic cell» on that map represents one of the basic meanings. For some cells the words with the closest meanings are shown. (The words in Golem may have several meanings, i.e. belong to several cells.)

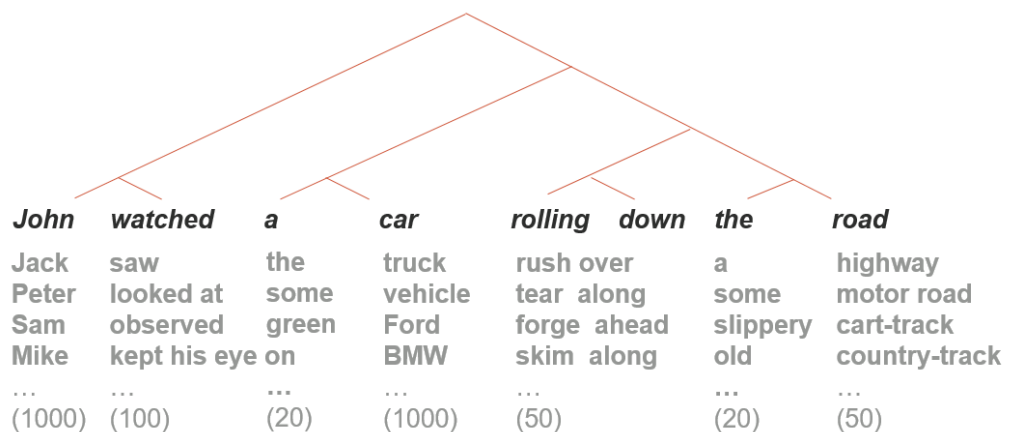
To represent the meaning of phrases and sentences we extend this initial semantic space, adding new meanings in the course of learning as most frequently used combinations of already known ones.

Since new meanings don’t change the existing ones, such unsupervised semantic learning may proceed in the course of semantic indexing. Thus, Golem gradually increases its semantic knowledge, reading new volumes of documents in various domains.



## Semantic search

Semantic indexing may be used, e.g. for semantic search, which takes into account deep linguistic structures, describing the general 'scenarios': who makes what with whom, etc. There are zillions of ways to express one and the same deep structure using different words. Golem provides compact sparse representation for the meaning of sentences, allowing fast semantic search, using inverted semantic index.



**1 deep structure:**  $1000 \times 100 \times 20 \times 1000 \times 50 \times 20 \times 50 = 10^{12}$  **surface variants**

In a way, semantics describe typical 'scenarios', where each 'role' may be played by that or another word. Our inverted index contains both words and their roles in scenarios.

## Collecting facts with semantic search

FactBook uses semantic search (together with the keyword search) to collect facts on a given topic from numerous documents. The topic of interest is automatically detected, based on the facts, already selected by the user.

Understanding the meaning of facts helps the system to grasp the interest profile from a handful of examples. The more facts user saves in his archive, the better are the new facts, extracted from the documents collection.

The following figure illustrates the facts, selected for a user, interested in the history of Newton's great discovery.

The screenshot shows the Factbook web application interface. At the top, the 'Factbook' logo is on the left, a search bar with the text 'when Newton discovered gravitation law' and a 'SEARCH' button is in the center, and 'My profiles' with a user icon is on the right. Below the header, there are three steps: 'STEP 1: DEFINE A TOPIC', 'STEP 2: COLLECT FACTS' (which is highlighted with a dark background), and 'STEP 3: EXPORT FACTS (30)'. The main content area has a blue header with the topic 'when Newton discovered gravitation law'. Below this, there are two fact cards. The first card is titled 'Sir Isaac Newton Discovered Gravity When An Apple Fell ... on Sep.11, 2011' and has an 'HTML' link. It contains three facts, each with a left arrow, a right arrow, and a 'save' icon. The facts are: 'By 1666 **Newton** had early versions of his three **laws** of motion.', 'Newton's novel idea of 1666 was to imagine that the Earth 's **gravity** influenced the Moon, counter-balancing its centrifugal force. From his law of centrifugal force and Kepler's third law of planetary motion, Newton deduced the inverse-square law.', and 'After his 1679 correspondence with Hooke , **Newton**, by his own account, found a proof that Kepler's areal **law** was a consequence of centripetal forces, and he also showed that **if** the orbital curve is an ellipse under the action of central forces then the radial dependence of the force is inverse square with the **distance** from the centre.' The second card is titled 'Isaac Newton, the discoverer of the law of gravitation on Apr.24, 2016' and also has an 'HTML' link. It contains one fact: 'In 1665, while on a visit in his native village, he saw an apple fall from a tree and began wondering what force made the apple fall.'

FactBook thus provides easy to use interactive tool for collecting facts from the Web or other large documents collections.

## Conclusion

In conclusion, we developed biologically inspired unsupervised learning algorithms, capable of learning language model from scratch. This model is nonparametric, i.e. it can continuously learn, extending its semantic knowledge, and improving the quality of semantic search.