

Comparing UNSILO concept extraction to leading NLP cloud solutions

By Mario Juric, Head of R&D at UNSILO, Mads Rydahl, CVO at UNSILO, and Hilke Reckman, NLP specialist at UNSILO.ai (updated October 2019)

Machine learning and artificial intelligence tools are promoted as solutions to some of humanity’s hardest challenges. But Machine learning can be applied to the same problem in many ways, and service providers may apply the same methods and still return different results. How can we meaningfully compare the results of machine learning tools from different providers? In this paper, an update of our 2017 paper, we provide an overview of the machine learning techniques used by UNSILO, and compare the output of the UNSILO Concept Extraction Service to that of other leading concept extraction tools.

Background

Although machine learning and artificial intelligence tools can be used to solve a number of different tasks that were previously the exclusive domain of Subject Matter Experts (SMEs), they do not “understand” knowledge like a human expert. Like most natural-language analytics providers, UNSILO uses a combination of probabilistic Natural Language Processing (NLP), structured knowledge in the form of ontologies and thesauri, hard-coded rules, and adaptive machine learning to determine the most important elements in text, and power services like document similarity, reader interest profiles, recommendation of reviewers or publication venues, and trend analysis.

Methodology

For this White Paper, the UNSILO Concept Extraction API was compared with the most widely adopted concept extraction services available today; the Google Cloud Natural Language API, the Microsoft Cognitive Services Text Analytics API, the IBM Watson Alchemy Language API, and the Amazon Comprehend Keyphrase Extraction API. To test performance across a variety of different subjects and terminologies, we randomly selected scholarly articles from four domains: Nanotech, Biomedical Science, Computer Science, and Food & Nutrition Science.

The full text of each article was submitted to each of the designated API services, and from each service, the top 20 concepts were examined according to a set of qualitative criteria: a) Relevance to the subject matter of the article, b) Specificity and unambiguity, c) Syntactic completeness, and d) Uniqueness; whether a concept is a synonym of another concept in the same set. Based on these criteria, each concept was assigned to one of four classes, and a corresponding point score was awarded, resulting in an aggregated document evaluation score, calculated as the sum total of the class score of the top 20 concepts. For example, correctly identified ontology terms like “KNN” and “Vitamin D” were classified as “Relevant broad Concepts”, which contribute one point to the document evaluation score, while longer phrases with an unambiguous meaning that are in common use within the domain were classified as “Relevant Precise Concepts”. Duplicate concepts and concepts that were deemed synonymous with another concept in the same set were classified as “Irrelevant or Redundant”, as were concepts with no connection to the subject matter..

Relevant Precise Concept	2 points
Relevant Broad Concept	1 point
Irrelevant, Redundant, Ambiguous	0 point
Fragment, Error, Noise	-1 point

Caution: In contrast to the other services, the publicly available Microsoft Cognitive Services Text Analytics API and the Amazon Comprehend Keyphrase Extraction API only parse the first 5K of each document. Perhaps counterintuitively, this may favour these services since the documents used were scholarly articles starting with an abstract of approximately 5K, where most noun phrases are highly relevant to the subject matter of the whole article.

Results and Analysis

Results show that the UNSILO Concept Extraction API does a better job at identifying relevant concepts in every tested domain, scoring on average 30.0

points per article compared to 12.9 points for all other services across all domains. This corresponds to an average score 2.3 times higher than the competition. The second best score was obtained by the Microsoft Cognitive Services API, which averaged 16.5 points across all domains. The Performance Summary and the Service Output and Classifications can be viewed in detail in Table 1 and Table 2.

Detecting and giving precedence to *multi-word terms* is key to successful fingerprinting and classification, because single word terms tend to be ambiguous or imprecise, whereas multi-word terms typically are **unambiguous** and more **specific**. For example, “fiber” and “intake” can refer to many things, but “dietary fiber intake” represents a clear concept that helps a user understand what a document is about. Previously only UNSILO and Microsoft returned good quality multi-word phrases, but now both IBM and Google have moved towards returning more multi-word terms. The results of IBM have improved significantly and are now comparable to those of Microsoft. Google performs less consistently in this area than IBM and still returns many unhelpful single-word terms.

An important challenge in extracting key phrases is to correctly detect phrase boundaries, to meet the criterion of **syntactic completeness**. A phrase should be coherent and self-contained. All systems occasionally make errors here, except for Amazon, which seems to extract whole noun phrases only (but including stop words like ‘the’). UNSILO returns much fewer incomplete concepts than its closest competitors.

To avoid redundancy and meet the **uniqueness** criterion, it is helpful to apply some normalization that maps similar phrases to a common form. Systems tend to return near-duplicates, for example ‘Dialogue acts’ and ‘dialog acts’ (different spelling and capitalization), ‘HIV’ and ‘Human immunodeficiency virus’ (abbreviation and expanded form). Even in the best performing systems there is still some room for improvement in recognizing variations of the same concepts.

Identifying suitable phrases is important, but not enough. These phrases also need to be scored for **relevance**, to correctly reflect the topics of a document. The tested systems differ quite a bit from each other in terms of relevance scoring. Amazon Comprehend presents very many key phrases with a confidence score of over 0.99, although few of these are truly relevant. The fact that they return and score each occurrence separately, rather than aggregation statistics about a concept severely limits the information they have available for scoring. IBM appears to have reduced its previous preference for returning named entities, which improves its performance, at least on scientific articles.

Discussion and Conclusions

We have sought to evaluate the UNSILO Concept Extraction API by comparing the output to that of competing systems. We have shown that such a comparison can be quite insightful. The primary limitation encountered was that several of the competing public APIs could not process a complete article. The manual evaluation may be somewhat subjective, and the number of points assigned to individual phrases may be subject to debate, but the overall picture is nevertheless rather clear: UNSILO performs best, with Microsoft coming in second, now closely followed by IBM.

Why then, does the UNSILO Concept Extraction work so much better than the competing services? Part of the explanation has to do with an extraction strategy that favours multi-word phrases. Such phrases, provided that their boundaries are correctly detected, are much more likely to capture precise meaning than single-word terms, and they are much less likely to be ambiguous. Another key factor that explains UNSILO’s success is a relevant background corpus. UNSILO Concept Extraction was trained on a corpus of scholarly articles. The patterns learned from this corpus inform the decisions on phrase boundaries and relevance. One could argue that this is an unfair advantage. However, one could also take it to mean that a general purpose keyphrase extraction API, which cannot be trained and adapted to a specific corpus, will necessarily show limited performance.

Table 1: Concept Quality Across Academic Domains

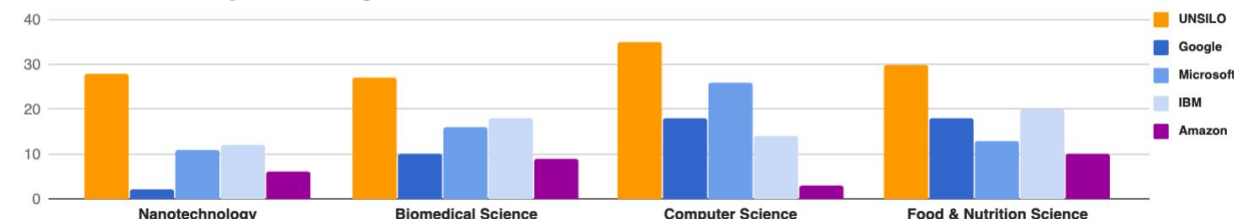


Table 2: Service Output and Classification of Concepts found in Scholarly Articles

A: Enhanced piezoelectric properties in vanadium-modified lead-free (K_{0.485}Na_{0.51}Li_{0.015})(Nb_{0.88}Ta_{0.1}V_{0.02})O₃ ceramics prepared from nanopowders

<http://www.sciencedirect.com/science/article/pii/S092538814027820>

<https://file.unsilio.com/unsilo-api-semantic-suite/>

UNSILO Concept Extraction API	Score	Eval
Common form		
KNN	1.00	1
KNN Ceramic	0.58	2
KNN System	0.46	2
PZT	0.44	1
piezoelectric ceramics	0.41	1
Room Temperature Dielectric Constant	0.42	2
Piezoelectric Property	0.41	1
Ferroelectric	0.40	1
Piezoelectric	0.40	1
Pure KNN	0.39	2
KNN Crystal	0.37	2
High Tetragonality	0.36	2
Li-doped KNN	0.35	1
Perovskite	0.34	2
KNbO ₃	0.34	1
Electromechanical Coupling Factor	0.34	2
Dense KNN	0.34	2
Piezoelectric Ceramic	0.34	2
O ₃ Ceramic	0.33	-1
Good Piezoelectric Property	0.33	1
Electromechanical Coupling	0.33	1
28		

<https://cloud.google.com/natural-language/>

Google Cloud Natural Language API	Score	Eval
Common form		
KNNV0	0.03	0
Vxj O3	0.02	-1
K0	0.01	-1
piezoelectric ceramics	0.01	2
PZT	0.01	0
balls	0.01	0
xTaO	0.00	-1
Ec	0.00	0
scanning tunneling microscopes	0.00	1
Curie temperature	0.00	1
niobate system	0.00	0
piezoceramics system	0.00	1
materials	0.00	0
V + content	0.00	1
work	0.00	0
lead	0.00	1
state reaction method	0.00	-1
Fig	0.00	-1
µC	0.00	-1
titanate	0.00	0
2		

<https://azure.microsoft.com/en-us/services/cognitive-services/text-analytics/>

Microsoft Cognitive Services Text Analytics API	Score	Eval
Common form		
high piezoelectric properties	NA	0
enhanced piezoelectric properties	NA	0
good piezoelectric properties	NA	1
better piezoelectric properties	NA	0
superior piezoelectric properties	NA	0
enhancement of piezoelectric properties	NA	2
O ₃ ceramics	NA	-1
alternative lead-free piezoelectric ceramic	NA	2
piezoelectric properties of lead-free piezo	NA	2
piezoelectric applications	NA	1
study of KNN ceramics	NA	0
sinterability of KNN ceramics	NA	2
Pure KNN ceramics	NA	2
NaO	NA	-1
K0	NA	-1
piezoelectric charge constant	NA	2
corresponding ceramics	NA	0
textured ceramics	NA	1
properties comparable	NA	-1
sensitivity of properties	NA	0
11		

<https://file.unsilio.com/unsilo-api-semantic-suite/>

IBM Watson AlchemyLanguage API	Score	Eval
Common form		
high energy ball milling	0.60	2
piezoelectric properties	0.59	2
KNNV0	0.59	0
Nb0.88Ta0.1V	0.58	-1
crystallite size	0.57	2
K0.485Na0.51	0.57	-1
morphotropic phase boundary	0.57	1
higher values of density	0.56	0
gradual increase	0.56	0
dielectric constant	0.56	2
industrial use	0.55	1
room temperature	0.55	0
grain growth	0.55	1
x increases	0.55	-1
piezoelectric properties of lead-free piezo	0.54	2
remnant polarization	0.54	1
optimum c	0.54	0
V5	0.54	0
KNNV1	0.54	0
lead-free	0.54	1
12		

<https://aws.amazon.com/comprehend/features/>

Amazon Comprehend Keyphrase Extraction API	Score	Eval
Common form		
The study	1.00	0
global restrictions	1.00	0
these techniques	1.00	0
terms	1.00	0
the formation	1.00	0
the field	1.00	0
This composition	1.00	0
the sinterability	1.00	1
an optimum introduction	1.00	0
its usage	1.00	0
their lead-based counterparts	1.00	0
a new solid solution	1.00	0
our search	1.00	0
recent years	1.00	0
human beings	1.00	0
A study	1.00	0
lead-free piezoceramics	1.00	2
the piezoelectric applications	1.00	1
high densities	1.00	2
The enhancement	1.00	0
6		

B: Effects of vitamin D supplementation on the bone specific biomarkers in HIV infected individuals under treatment with efavirenz

<https://bmcsresnotes.biomedcentral.com/articles/10.1186/1756-0500-5-204>

UNSILO Concept Extraction API	Score	Eval
Common form		
Vitamin D	1.00	1
Bone Mineral Density	0.35	2
Vitamin D Receptor	0.23	2
Serum CTX Concentration	0.21	2
HIV Positive Individual	0.19	2
Bone Formation Marker	0.17	2
Bone Resorption Marker	0.16	2
Bone Formation	0.16	1
Bone Formation Biomarker	0.16	0
Serum OC Level	0.16	2
Bone Biomarker	0.14	1
Serum Vitamin	0.12	-1
HIV-infected Patient	0.12	2
HIV Negative Individual	0.12	2
Collagen	0.11	1
Efavirenz	0.11	2
Bone Resorption	0.11	2
Osteocalcin Concentration	0.11	2
25-OH Vitamin	0.10	-1
Hepatitis C	0.10	1
27		

Google Cloud Natural Language API	Score	Eval
Common form		
both	0.17	1
vitamin d	0.02	-1
Serum CTX	0.01	1
health care problem	0.01	1
supplementation	0.01	1
treatment	0.01	1
88.4 %	0.01	-1
patient	0.00	0
HIV Efavirenz	0.00	-1
Vitamin D deficiency	0.00	2
infection	0.00	1
Baseline osteocalcin	0.00	2
collagen telopeptide	0.00	2
Findings	0.00	0
Human immunodeficiency virus	0.00	0
level range	0.00	0
Background	0.00	0
Vitamin D Findings Background	0.00	-1
diagnosis	0.00	1
10		

Microsoft Cognitive Services Text Analytics API	Score	Eval
Common form		
supplementation of vitamin D	NA	2
Vitamin D deficiency	NA	2
catabolism of vitamin D	NA	2
IU vitamin D	NA	-1
vitamin D deficient HIV positive individual	NA	0
Effects of vitamin D supplementation	NA	2
HIV Efavirenz Vitamin D Findings Backg	NA	-1
HIV-infected patients	NA	0
HIV positive patients	NA	2
HIV negative individuals	NA	2
HIV diagnosis	NA	1
Significant percent of HIV infected indivi	NA	0
duration of HIV infection	NA	1
bone disorders	NA	1
bone specific biomarkers	NA	2
HIV viral load	NA	1
multis individuals	NA	-1
early diagnosis of HIV infection	NA	0
bone biomarkers' concentrations	NA	2
patients' efavirenz	NA	-1
16		

IBM Watson AlchemyLanguage API	Score	Eval
Common form		
Effects of vitamin D supplementation	0.77	2
serum levels	0.67	1
positive patients	0.66	-1
catabolism of vitamin D	0.64	2
bone biomarkers' concentrations	0.60	0
Vitamin D deficiency	0.58	2
specific biomarkers	0.58	-1
P value Serum	0.56	-1
antiretroviral drugs	0.56	2
supplementation of vitamin D	0.55	2
Baseline osteocalcin	0.55	2
bone	0.54	1
enough duration of follow-up	0.54	-1
clinical data	0.54	1
Serum levels of vitamin D	0.54	2
normal range	0.54	0
single dose	0.54	0
HIV	0.53	1
bone formation	0.53	2
activation of osteoblasts	0.53	2
18		

Amazon Comprehend Keyphrase Extraction API	Score	Eval
Common form		
a fundamental element	1.00	0
the effects	1.00	0
supplementation	1.00	1
this population	1.00	1
the study	1.00	0
comparison	1.00	0
bone disorders	1.00	1
treatment	1.00	1
regulation	1.00	0
this infection	1.00	0
the catabolism	1.00	1
indicator	1.00	0
antiretroviral drugs	1.00	2
Vitamin D deficiency	1.00	2
a risk factor	1.00	0
Nine patients	1.00	0
the activation	1.00	0
any interventional study	1.00	1
any adverse event	1.00	0
duration	1.00	0
9		

C: A novel density-based clustering method using word embedding features for dialogue intention recognition

<https://rd.springer.com/article/10.1007/s10586-016-0649-7>

UNSILO Concept Extraction API	Score	Eval
Common form		
Dialogue Act	1.00	2
Latent Dirichlet Allocation	0.55	1
Deep Neural Network	0.55	1
Word Embedding	0.48	2
Latent Semantic Analysis	0.42	1
Lexical Feature	0.37	1
Emotion Classification	0.36	2
Support Vector Machine	0.36	2
Smart Home System	0.36	2
Density-based Cluster Method	0.32	2
Non-convex Cluster	0.32	2
Density-based Cluster	0.32	2
Feature Selection Method	0.30	2
Stochastic Neighbor Embedding	0.30	2
Word Frequency Distribution	0.29	2
Natural Language Understand	0.28	2
User Utterance	0.26	2
Cluster Feature	0.26	2
Maximum Entropy Model	0.25	2
Social Network Analysis	0.25	1
35		

Google Cloud Natural Language API	Score	Eval
Common form		
knowledge base	0.01	1
corpora	0.01	0
SVM	0.01	2
synonyms	0.01	2
i-neighbourhood value	0.01	2
acts	0.01	-1
word frequency distribution	0.00	2
tasks	0.00	0
dialogue system	0.00	2
DBSCAN	0.00	2
similarities	0.00	0
Korean	0.00	1
Word2Vec	0.00	2
efficiency problems	0.00	-1
degree	0.00	0
embedding features	0.00	1
corpus data	0.00	1
proportion	0.00	0
dialogue acts	0.00	0
Dialogue acts	0.00	2
16		

Microsoft Cognitive Services Text Analytics API	Score	Eval
Common form		
dialogue acts	NA	2
dialogue intention recognition	NA	2
dialogue act classification	NA	2
word similarity	NA	2
dialogue systems	NA	2
emotion classification	NA	2
previous classification models	NA	0
Various classification models	NA	0
embedding space	NA	1
classification performance	NA	1
user intention analysis	NA	2
user utterances	NA	2
Supervised learning-based classification	NA	1
similarity values	NA	0
emotion recognition	NA	2
problem of data sparseness	NA	0
data sparseness problem	NA	2
original lexical features	NA	1
word frequency distribution	NA	2
sufficient training data	NA	0
26		

IBM Watson AlchemyLanguage API	Score	Eval
Common form		
excessive use of lexical features	0.69	0
knowledge base	0.64	1
Various classification models	0.63	0
novel density	0.61	-1
original lexical features	0.60	1
dialogue intention recognition	0.60	2
efficient use of lexical features	0.59	0
problem of data sparseness	0.57	2
emotion recognition	0.57	2
Experimental results	0.56	0
understanding user utterances	0.56	2
direct use of word	0.56	-1
dialogue acts	0.56	2
nature of the word frequency distribution	0.56	0
previous classification models	0.56	0
convex clusters	0.56	0
dialogue systems	0.55	2
training set	0.55	1
different domains	0.54	1
clusters	0.54	0
14		

Amazon Comprehend Keyphrase Extraction API	Score	Eval
Common form		
words	1.00	1
the efficiency	1.00	0
our proposed model	1.00	0
the distances	1.00	0
a significant amount	1.00	0
the system	1.00	0
the problem	1.00	0
the lesser representation	1.00	0
This paper	1.00	0
lexical features	1.00	1
a user	1.00	1
different domains	1.00	0
a speaker	1.00	1
a firm	1.00	0
Recent research	1.00	0
the improvement	1.00	0
the actual degree	1.00	0
the objectives	1.00	0
this problem	1.00	0
lexical feature	1.00	0
3		

D: Replacing carbohydrate with protein and fat in prediabetes or type-2 diabetes: greater effect on metabolites in BMC than plasma

<https://rd.springer.com/article/10.1186/s12986-016-0063-4>

UNSILO Concept Extraction API	Score	Eval
Common form		
Plasma Lp-PLA2 Activity	1.00	2
Lp-PLA2 Activity	0.74	1
High Lp-PLA2 Activity	0.40	0
Plasma Lp-PLA2	0.39	1
Plasma ox-LDL	0.35	2
12-week Dietary Intervention	0.32	2
Impaired Fasting Glucose	0.31	2
Score Scatter Plot	0.28	0
Basal Metabolic Rate	0.25	2
PBMC Gene Expression	0.22	2
Glycemic Control	0.21	2
Type-2 Diabetes	0.20	2
High Lp-PLA2	0.20	0
ox-LDL Level	0.19	1
Total Energy Expenditure	0.19	2
lysPC Level	0.19	1
THP-1 Monocyte	0.18	2
Plasma Metabolite	0.18	2
Wallac Victor2 Multilabel Counter	0.16	2
Dietary Fiber Intake	0.15	2
30		

Google Cloud Natural Language API	Score	Eval
Common form		
Table 2	0.08	-1
LDL	0.06	2
Lp-PLA2	0.01	2
PLS-DA	0.01	2
Diabetes	0.0	