UNSILO White Paper

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

Machine learning and artificial intelligence tools are promoted as solutions to some of mankind'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 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, 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 four 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, by 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 one-gram ontology terms like "KNN" and "Vitamin D" were classified as "Relevant broad Concepts", which contribute one point to the 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, including names and geolocations of authors and sponsoring organizations, which should be provided as metadata properties.

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 the these services since the documents used were scholarly articles starting with an abstract of approximately 5K, where almost every noun phrase is 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 33.0 points per article compared to 13.4 points for all other services across all domains. This corresponds to an average score 2.5 times higher than the competition. The

second best score was obtained by the Microsoft Cognitive Services API, which averaged 22.5 points across all domains. The Performance Summary and the Service Output and Classifications can be viewed in detail in Table 1 and Table 2.

One of the criteria was that the extracted concepts be **specific** and **unambiguous**. Most of the competing services return broad ontology concepts like "HIV" or "Ceramics" or ambiguous concepts like "Study" or "Feature" which have low descriptive value and are less suited for classification or fingerprinting of documents. Nearly all terms judged to be relevant and precise were multi-word terms, but only Microsoft returns multi-word phrases of a quality comparable to UNSILO, and this may be an important factor of their relative success. Google, on the other hand, almost exclusively returns ambiguous single word terms.

Detecting and giving precedence to multi-word terms is the 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 precise. 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. 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. Microsoft appears to have a slightly more risky strategy than UNSILO resulting in some longer phrases, like "Significant percent of HIV infected individuals" but also some noise, such as "nonobese subjects aged". Microsoft has the highest percentage of noise ratio at 10% compared to zero for UNSILO. Amazon includes more function words (AKA stop words) in the the phrases than any of the others. Sometimes this does not matter very much, e.g. "the cosine similarity", but in other cases it makes the concept less self-contained, e.g. "other perovskite compounds"

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 does extract some good phrases, but seems to have an issue with its scoring algorithm, as it presents very many key phrases with a confidence score of over 0.99, but few of these are truly relevant. IBM appears to favour named entities, which in scientific articles unfortunately tends to put names of authors and their affiliations among the top concepts.

To avoid redundancy and meet the **uniqueness** criterion, it is helpful to apply some normalization that maps similar phrases to a common form. Google not only fails to do basic lemmatization (e.g. collapsing "feature" and "features"), but tends to return the exact same phrase multiple times. Even in the best performing systems there is still some room for improvement in recognizing variations of the same concept, e.g. involving synonyms.

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, and Microsoft comes in as a decent second. The other competitors do not perform very well.

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 favors multi-word phrases. Such phrases, provided that their boundaries are correctly detected, are much more likely to capture precise meaning than single word terms are, 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 informs 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

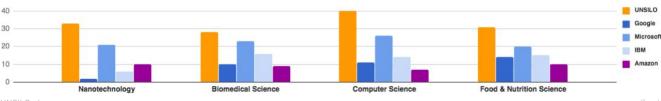


Table 2: Service Output and Classification of Concepts found in Scholarly Articles

A: Enhanced piezoelectric properties in vanadium-modified lead-free (K0.485Na0.5Li0.015)(Nb0.88Ta0.1V0.02)O3 ceramics prepared from nanopowders

UNSILO Concept Extraction API		
Common form	Score	Eva
KNN	1.00	1
KNN Ceramic	0.58	2
KNN System	0.46	2
PZT	0.44	1
Room Temperature Dielectric Constant	0.42	1
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	2
Perovskite	0.34	2
KNbO3	0.34	2
Electromechanical Coupling Factor	0.34	2
Dense KNN	0.34	2
Piezoelectric Ceramic	0.34	2
O3 Ceramic	0.33	2
Good Piezoelectric Property	0.33	(
Electromechanical Coupling	0.33	1

Intos //cloud.google.com/natural-lan Google Cloud Natural Lang	uage API	
Common form	Score	Eva
ceramics	1.00	
KNNV3	1.00	2
KNN	0.80	8
KNNV2	0.60	-
phase	0.40	1
properties	0.40	1
ceramics	0.40	
ceramics	0.40	1
К0	0.40	8
Vx) 03	0.40	-
properties	0.20	
microstrain	0.20	9
study	0.20	9
O3 ceramic	0.20	1
stress	0.20	1
ceramics	0.20	-
materials	0.20	
balls	0.20	-
system	0.20	
PZT	0.20	8

Common form	Score	Eva
Enhanced piezoelectric properties	NA	2
good piezoelectric properties	NA	0
superior piezoelectric properties	NA	0
enhancement of piezoelectric properties	NA	0
piezoelectric properties of lead-free piezo	NA.	2
alternative lead-free piezoelectric cerami	NA.	2
O3 exhibits high piezoelectric properties	NA	-1
piezoelectric applications	NA	2
O3 ceramics	NA	2
study of KNN ceramics	NA	0
Pure KNN ceramics	NA	2
sinterability of KNN ceramics	NA	2
sensitivity of properties	NA	-1
properties comparable	NA	-1
Pb-free piezoelectric materials	NA	2
piezoelectric charge constant	NA	2
Na0	NA	1
conventional sintering of KNN	NA	2
K0	NA	1
maximum piezoelectric coefficient	NA	2

KNN	e	Eval
KNNV	0	1
KNNV	3	-1
PZT	9	-1
KNN 0. nancpowders 0. LI 0. LI 0. Gibbs 0. Ps 0. Williamson – Hall 0. Department of Science and Technology 0. Japan 0. India 0. Li Askimoto 0.	7	-1
Name	8	1
KNN 0.0 Li 0.0 Glibbs 0.0 Ps Williamson – Hall 0.0 Department of Science and Technology 0.0 Japan 0.0 Japan 0.0 Kakimoto 0.0	8	0
Li 0.0 edge technologies 0.0 (Gibbs 0.0 Ps 0.0 Williamson – Hall 0.0 Wayne Kerr 0.0 Japan 0.0 Kakimoto 0.0	7	1
edge technologies	5	0
Gibbe 0. Ps 0. Williamson – Hall 0. Department of Science and Technology 0. Japan 0. Japan 0. India 0. Jakakimoto 0.	3	1
Ps 0.0 Williamson – Hall 0.0 Wayne Kerr 0.0 Japan 0.0 Kakinoto 0.0 Kakinoto 0.0	2	0
Williamson – Hall 0. Department of Science and Technology o 0. Wayne Kerr 0. Japan 0. India 0. Kakimoto 0.	2	0
Department of Science and Technology o Uwayne Kerr Ugapan	0	-1
Wayne Kerr 0. Japan 0. India 0. Kakimoto 0.	8	-1
Japan 0. India 0. Kakimoto 0.	6	2
India 0. Kakimoto 0.	5	1
Kakimoto 0.	5	0
	5	1
V	5	1
X-ray 0.	4	1
non-stoichiometry 0.	4	1

Common form	Score	Eva
lead	0.99+	
sensitivity	0.99+	(
Introduction	0.99+	- (
industrial use	0.99+	
other perovskite compounds	0.99+	(
the lead-free alternatives	0.99+	
Ta	0.99+	
Pr	0.99+	
evaporation	0.99+	
lead-free piezoceramics	0.99+	- 2
piezoelectric properties	0.99+	2
our search	0.99+	-
Ec	0.99+	(
Nb	0.99+	
Pb	0.99+	
energy conversion	0.99+	
MPB	0.99+	(
The enhancement	0.99+	(
special handling	0.99+	
LiTaO3 [9]	0.99+	-

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

https://bmcresnotes.biomedcentral.com/articles/10.1186/1756-0500-5-204

Common form	Score E	
Vitamin D	1.00	
Bone Mineral Density	0.35	
Vitamin D Receptor	0.23	
Serum CTx Concentration	0.21	
HIV Positive Individual	0.19	
Bone Formation Marker	0.17	
Bone Resorption Marker	0.16	
Bone Formation	0.16	
Bone Formation Biomarker	0.16	
Serum OC Level	0.16	
Bone Biomarker	0.14	
Serum Vitamin	0.12	
HIV-infected Patient	0.12	
HIV Negative Individual	0.12	
Collagen	0.11	
Efavirenz	0.11	
Bone Resorption	0.11	
Osteocalcin Concentration	0.11	
25-OH Vitamin	0.10	
Hepatitis C	0.10	

Common form	Score	Eval
HIV	1.00	1
drugs	0.12	0
individuals	0.09	(
patients	0.06	(
CTx	0.06	1
patients	0.03	(
vitamin d	0.03	1
concentrations	0.03	(
treatment	0.03	1
vitamin D	0.03	0
effects	0.03	0
study	0.03	. (
supplementation	0.03	. (
bone biomarkers	0.03	1
efavirenz.	0.03	1
oc	0.03	1
vitamin D supplementation	0.03	1
bone biomarkers	0.03	1
individuals	0.00	(
bone	0.00	1

Common form	Score	Eva
supplementation of vitamin D	NA	. 0
Vitamin D deficiency	NA	
bone biomarkers	NA	
vitamin D deficient HIV positive individual	NA	8 8
HIV positive patients	NA	1
HIV-infected patients	NA	a 8
catabolism of vitamin D	NA	9 8
IU vitamin D	NA	-
Effects of vitamin D supplementation	NA	
HIV Efavirenz Vitamin D Findings Backgr	NA	-
HIV negative individuals	NA	
duration of HIV infection	NA	8 5
bone disorders	NA	
months of vitamin D administration	NA	. 9
HIV viral load	NA	
Significant percent of HIV infected individ	NA	9
early diagnosis of HIV infection	NA	
bone specific biomarkers including osteo	NA	1
favorable bone formation	NA	2 1
indicator of bone resorption	NA	1

Common form	Score	Eval
vitamin D	1.00	1
Vitamin D deficiency	0.66	2
HIV infection	0.63	2
efavirenz	0.46	1
bone fracture	0.30	1
ALP	0.26	1
HIV virus	0.23	1
HIV Clinic of Iranian HIV/AIDS Research	0.22	0
Chronic hepatitis B	0.20	1
oc	0.19	1
Chronic hepatitis C	0.19	1
Tehran University of Medical Sciences	0.19	
hypogonadism	0.18	- 1
Chicago	0.18	0
Tehran	0.18	0
drug abuse	0.18	1
adrenal insufficiency	0.18	1
USA	0.17	0
Iran	0.17	0
Osteoblast	0.17	1

Common form	Score	Eval
Associations	0.99+	0
Results	0.99+	0
the activation	0.99+	0
Significant percent	0.99+	0
the catabolism	0.99+	1
treatment	0.99+	1
bone metabolism	0.99+	2
HIV negative individuals	0.99+	2
HIV	0.99+	1
the effects	0.99+	0
Vitamin D deficiency	0.99+	2
any study	0.99+	-1
HIV infection	0.99+	2
present study	0.99+	0
107	0.99+	-1
the world	0.99+	0
98	0.99+	-1
survival	0.99+	- 1
these individuals	0.99+	0
the patients	0.99+	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

1.00 0.55 0.55 0.48 0.42 0.37	2 2
0.55 0.48 0.42 0.37	2 2
0.48 0.42 0.37	2
0.42 0.37	2
0.37	
177.7	
0.36	2
	2
0.36	2
0.36	2
0.32	2
0.32	2
0.32	2
0.30	2
0.30	2
0.29	2
0.28	2
0.26	2
0.26	2
0.25	2
0.05	2
	0.26 0.26 0.25 0.25

Common form	Score	Eva
features	1.00	1
features	0.56	
clustering	0.44	1
model	0.22	(
SVM	0.22	1
corpora	0.11	1
model	0.11	(
features	0.11	(
feature	0.11	(
emotion classification	0.11	1
models	0.11	0
methods	0.11	1
words	0.11	1
corpus	0.11	(
clustering method	0.11	1
dialogue acts	0.11	1
results	0.11	(
word frequency distribution	0.11	1
words	0.00	
Korean	0.00	1
		4

Common form	Score	Eval
dialogue acts	NA	2
dialogue intention recognition	NA	2
word similarity	NA	2
word embedding features	NA	2
dialogue act classification	NA	2
dialogue systems	NA	2
word embedding model	NA	2
user intention analysis	NA	2
word embedding methods	NA.	1
understanding user utterances	NA	2
previous classification models	NA	0
emotion classification	NA	2
Various classification models	NA.	0
problem of data sparseness	NA	0
data sparseness problem	NA	0
word frequency distribution	NA	2
analysis of embedding features	NA	1
sufficient training data	NA.	2
emotion recognition	NA	0
extensive amounts of training data	NA.	0

Common form	Score	Eval
natural language	1.00	2
Representative	0.80	0
Kim	0.67	0
SVM	0.62	1
f	0.55	0
application domain	0.54	2
core point	0.52	(
Paul Ekman	0.48	(
International Organization	0.47	
Sect.	0.46	0
Latent Semantic Analysis	0.46	2
SVM	0.45	1
knowledge base	0.43	2
modi	0.43	-
Twitter	0.43	- 31
Lee	0.43	0
Hasegawa	0.43	(
neural network	0.42	2
CRF	0.42	1
Xu	0.42	0

Common form	Score	Eval
the similarity	0.99+	0
the system	0.99+	0
This paper	0.99+	0
time	0.99+	0
a key summary	0.99+	1
a speaker	0.99+	1
WordNet	0.99+	1
this problem	0.99+	0
Another way	0.99+	0
the former pair	0.99+	-1
the cosine similarity	0.99+	2
these classification models	0.99+	1
Representative	0.99+	0
training data	0.99+	1
conversation	0.99+	1
such models	0.99+	-1
dialogue act classification	0.99+	2
words and emotions	0.99+	1
" and "bat"	0.99+	-1
the skewed nature	0.99+	-1

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

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

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

Common form	Score	Eva
group	1.00	(
PBMC	0.96	4
PBMC	0.30	(
activity	0.13	(
Ox-LDL	0.13	1
Lp-PLA2	0.09	
Lp-PLA2	0.09	(
grains	0.09	
PLS-DA	0.09	-
groups	0.04	(
blood cells	0.04	2
UPLC-LTQ-Orbitrap MS	0.04	3
components	0.00	(
intervention	0.00	(
epidemic metabolic disorder	0.00	:
plasma	0.00	
metabolites	0.00	
rice	0.00	8
T2D	0.00	
subjects	0.00	(
		14

Common form	Score	Ev
Lp-PLA2 activity	NA	
Intervention diet	NA	
Dietary intervention	NA	
release of Lp-PLA2	NA	
plasma Lp-PLA2 activities	NA	
blood cells	NA	
week intervention study	NA.	y - 1
B cells	NA.	
usual refined-rice diet	NA.	
usual diet	NA	
T cells	NA	
natural killer cells	NA	
PBMC metabolites	NA.	
week intervention phase	NA.	a 8
replacement of refined rice	NA.	
blood glucose	NA.	
cooked refined rice	NA.	
typical diet	NA	
Nonobese subjects aged	NA.	9 9
Assessment of dietary intake	NA	

Common form	Score	Eval
PBMC	1.00	1
MDA	0.50	1
oleamide	0.45	1
blood glucose	0.44	1
IFG	0.43	1
Korea	0.43	0
PBMC	0.43	1
ВМІ	0.41	1
IGT	0.40	1
Lp-PLA2	0.40	1
PBMC cleamide	0.40	1
Hitachi Ltd.	0.38	0
Diabetes	0.37	1
ox-LDL	0.37	1
Thermo Fisher Scientific	0.36	0
Tokyo	0.35	0
Japan	0.35	0
fatty acids	0.35	1
immune system	0.34	1
metabolic rate	0.34	1
		15

Common form	Score	Eval
All participants	0.99+	0
inflammatory reactions	0.99+	2
2 week days	0.99+	0
all participants	0.99+	-1
replacement	0.99+	0
Diabetes	0.99+	1
The present study	0.99+	0
two subjects	0.99+	0
An intervention study	0.99+	1
the immune system	0.99+	1
T cells	0.99+	2
this intervention diet	0.99+	0
adults	0.99+	- 1
The remaining 80 subjects	0.99+	0
These blood cells	0.99+	1
metabolites	0.99+	1
a 12-week intervention phase	0.99+	1
the innate immune system	0.99+	1
an inverse association	0.99+	0
2013	0.99+	-1

UNSILO.ai www.unsilo.ai