# Neural Machine Reading for Domain-Specific Text Resources

## Sebastian Arnold sarnold@beuth-hochschule.de

Prof. Alexander Löser Prof. Felix A. Gers Prof. Philippe Cudré-Mauroux



Data Science and Text-based Information Systems (DATEXIS) Beuth University of Applied Sciences Berlin, Germany

## UNI FR

UNIVERSITÉ DE FRIBOURG UNIVERSITÄT FREIBURG

Department of Informatics Faculty of Science and Medicine University of Fribourg Fribourg, Switzerland

28.10.2020

# About Me

- 2006-2011Technische Universität BerlinBachelor of Science in Computer Science
- 2011–2015 **Technische Universität Berlin** Master of Science in Computer Science Student Assistant @ DIMA, Prof. Volker Markl
- 2015–2020 **Beuth University of Applied Sciences Berlin** Research Assistant @ DATEXIS, Prof. Alexander Löser
- 2016–2020 Université de Fribourg (CH) Doctoral Program in Computer Science PhD Student @ eXascale Infolab, Prof. Cudré-Mauroux

since 07/2020 Curalie GmbH, Berlin Machine Learning Expert



<u>@sebastianarnold</u>

# PhD Timeline & Main Contributions

# "Neural Machine Reading for Domain-Specific Text Resources" (since 06/2015)

Supervisors: Advisor: External examiner:	Prof. Alexander Löser Beuth University of Applied Sciences Berlin (D) Prof. Philippe Cudré-Mauroux eXascale Infolab, University of Fribourg (CH) Prof. Felix A. Gers DATEXIS group @ Beuth Prof. Laura Dietz Univ. of New Hampshire (US)	ACLE 2019 TOTAL AND ACLE 2019		
Publications:	DOLAP Workshop @ CIKM 2015 Melbourne, COLING I WWW Demos 2018 Lyon, TACL @ ACL2019 Florence, V	<b>Demos</b> 2016 Osaka, <b>WWW 2020</b> Taipei		
Reviewer for:	WWW'18/19/20, CIKM'20, ACL'18/19, AAAI'19, EMNL	.P'19, ESSLLI'18		
Talks:	German-Canadian Concourse @ Charité, Quo Vadis 2015, Zalando Meetup, Lange Nacht der Wissenschaften 2016, Holtzbrink Al Day 2017 @ SpringerNature, Dagstuhl Seminar 2019, Graduate School Data Science @ Einstein Center Berlin			
Industry Projects:	Siemens, SAP, Zalando SE, SpringerNature, Mobile.de, C	harité, BMWi, H2020		
Open Source:	https://github.com/sebastianarnold/			



- 1. Introduction
- 2. Background
- 3. TASTY A Robust Model for Efficient Entity Linking
- 4. SECTOR Coherent Topic Segmentation and Classification
- 5. CDV Contextual Document Representations for Answer Retrieval
- 6. Systems
- 7. Conclusion and Future Work

"Die Natur muß gefühlt werden, wer nur sieht und abstrahirt, kann ein Menschenalter, im Lebensgedränge der glühenden Tropenwelt, Pflanzen und Thiere zergliedern, er wird die Natur zu beschreiben glauben, ihr selbst aber ewig fremd sein."

- Alexander von Humboldt, an Johann Wolfgang von Goethe, Paris 3. Januar 1810.

# 1. Introduction

# Example: Information Search

"What are symptoms of Cystic fibrosis (CF)?"

- Query intent classification [Rose+Levinson 2004]
- Natural Language Processing [Jurafsky 2008]
- Information Extraction (IE) [Sarawagi 2008]
- Semantic representation [Moro 2014]



# **Knowledge Processing and Retrieval**

"What are symptoms of Cystic fibrosis (CF)?"

- Question Answering (QA)
- Information Retrieval (IR) [Manning 2008]
- Feedback Loop
- Exploration and refinement q = "Cystic fibrosis"~0.8 AND "symptoms"~0.8 [Marchionini 2006] Cystic fibrosis (CF) is a genetic disorder that causes mucus to build up and damage orga

**Cystic fibrosis** (CF) is a genetic disorder that causes mucus to build up and damage organs in the body, particularly the lungs and pancreas. Signs and symptoms may include salty-tasting skin; persistent coughing; frequent lung infections; wheezing or shortness of breath; poor growth; weight loss; greasy, bulky stools; difficulty with bowel movements; and in males, infertility. Over time, mucus buildup and infections can lead to permanent lung damage, including the formation of scar tissue (fibrosis) and cysts in the lungs. CF is caused by mutations in the *CFTR* gene and inheritance is autosomal recessive.<sup>[1][2][3]</sup> Treatment aims to relieve symptoms and usually includes respiratory therapies, inhaled medicines, pancreatic enzyme supplement, and nutritional supplements. Newer medications such as CFTR modulators have been approved for use in the United States. Ongoing research is focused on finding a cure for the disease.<sup>[2]</sup>



Sepastian Arnold 8

# GoOLAP: Answering Analytical Queries with Tables [2012]

COC Search. Explor	Pe. Decide.	cur	rrent White H	louse e	employees		1	Bearch			Start Tour
EXPLORE X	LIST X	← s	elect your res	sult							
Organization White House W	on/employments of /hite House										
		1	position/em	ployme	nts of White H	ouse			sho	wing 1 to 10 of 3680 e	ntries prev   next
	No States		country co	mpare v	city compare	person compare v	organization comp	oare v company compare v	provine	ceorstate compare v	position compa
	and second difference in the		1			Robert Gibbs	White House				press secretar
- Little Little			2			Jay Carney	White House				press secretar
			3			Rahm Emanuel	White House				chief of staff
and the second			4			Scott McClellan	White House				press secretar
			5			Robert Gibbs	White House				spokesman
We found 17.060 facts about the		:::	6			Jay Carney	White House				spokesman
Organization White House in	the World		7			Alberto Gonzales	White House				counsel
Wide Web.		:::	8			John Dean	White House				counsel
			9			Scott McClellan	White House				spokesman
Similar objects to White Ho	ouse	1	0			James Brady	White House				press secretar
Whitehouse (Company, 9 fa Whitehouse (Person, 8 facts	icts) s)		communica	tions of	White House				sh	owing 1 to 10 of 219 e	ntries previ pext
Whitehouse (City, 3 facts)		141	communica	cions of	write riouse				an	owing 1 to 10 01210 c	initias previntext
White House (Company, 1 f	act)		date	perso	on1 compare v	persondescription1	person compare v	organizationorcompany con	npare 🔻	persondescription	facility compare '
			1 1993-01-2	27				White House		small delegation	
Relations in this result ra	nk sort	:::	2				Eleanor Statues	White House			
	[Ŧ]		3				Lloyd Blankfein	White House			
generic relations (9371)	1.4.1	***	4					White House		CIA Director	
<ul> <li>position/employments (7</li> </ul>	7256)		5				John S. McCain , Jr.	White House			
<ul> <li>communications (288)</li> </ul>	L#.)	:::	6				Ehud Olmert	White House			
✓ diplomatic relations (131)	1) I		7				Douglas	White House			

A. Löser, <u>S. Arnold</u>, T. Fiehn (2012): **The GoOLAP Fact Retrieval Framework.** Business Intelligence. Springer <u>S. Arnold</u>, A. Löser, T. Kilias (2015):. **Resolving Common Analytical Tasks in Text Databases.** ACM DOLAP 2015

# Vision: Machine Reading [Etzioni 2006]

MR is "the automatic, unsupervised understanding of text [achieved by] the formation of a coherent set of beliefs based on a textual corpus and a background theory."



**Our goal:** Support the human information-seeking process with Machine Reading.

# Challenges for Domain-specific Language Understanding

# Traditional IE pipelines depend on task-specific training data or handcrafted rules.

- Domain-specific language
- Variations and noise
- Heterogenous document structure
- High task variance
- Insufficient training data
- Error propagation

- → adaptiveness and transferability
- $\rightarrow$  robustness, focus on recall
- $\rightarrow$  coherence, local context sensitivity
- $\rightarrow$  generalization, multi-task, zero-shot
- → efficiency, self-supervision, background knowledge
- → end-to-end design, differentiable models, dense representations

# **Hypothesis: Neural Machine Reading**

- 1. **Deep Neural Networks** enable **self-supervised training** of language models based on distributional information.
- 2. Machine Reading increases error tolerance and reduces adaptation cost for reading domain-specific text.
- 3. End-to-end models combine general training objectives with background knowledge to fulfill task-specific information needs.

# **Research Questions**



**RQ1:** Identify domain-specific **named entities**.



RQ2: Detect topics and structure in long documents.



RQ3: Embed

- discourse structure
- into document
- representations.



RQ4: Retrieve answer passages using vector representations.

# **Focus of the Thesis**



# Main Contributions and Publications

# **TASTY** – Named Entity Recognition and Linking

<u>S. Arnold</u>, F. A. Gers, T. Kilias and A. Löser (2016). Robust Named Entity Recognition in Idiosyncratic Domains. **arXiv:1608.06757 [cs.CL]** <u>S. Arnold</u>, R. Dziuba and A. Löser (2016).

TASTY: Interactive Entity Linking As-You-Type. COLING 2016 (Demos)

# **SECTOR** – Topic Segmentation and Classification

R. Schneider, <u>S. Arnold</u>, T. Oberhauser, T. Klatt, T. Steffek and A. Löser (2018). Smart-MD: Neural Paragraph Retrieval of Medical Topics. **WWW 2018 (Companion)**: 203–206

<u>S. Arnold</u>, R. Schneider, P. Cudré-Mauroux, F. A. Gers and A. Löser (2019). SECTOR: A Neural Model for Coherent Topic Segmentation and Classification. **TACL Vol. 7**: 169–184 (presented at ACL)

# **CDV** – Contextual Discourse Vectors for Answer Retrieval

<u>S. Arnold</u>, B. van Aken, P. Grundmann, F. A. Gers and A. Löser (2020). Learning Contextualized Document Representations for Healthcare Answer Retrieval. **WWW 2020**: 1332–1343

J.-M. Papaioannou, <u>S. Arnold</u>, F. A. Gers, A. Löser, M. Mayrdorfer, and K. Budde (2021).

Aspect-Based Passage Retrieval with Contextualized Discourse Vectors. Submitted to: ECIR 2021 System Demonstrations

TeXoo - A Zoo of Text Extractors. Apache V2 License: <u>https://github.com/sebastianarnold/TeXoo</u>

# 2. Background

# Distributional Hypothesis [Harris 1954]

Syntagmatic relations Combinations: "x and y and..."

- Neural Probabilistic Language Model: Learn to predict words based on their context [Bengio 2003]  $p(w_1, \dots, w_n)$ 
  - $p(w_1, \dots, w_T) = \prod_{t=1}^T p(w_t \mid w_1, \dots, w_{t-1})$

paint

colour

 $W_{+}$ 

dye

[Sahlgren 2008]

Paradigmatic relations Selections: "x or y or..."

green

blue

red

adores

likes

love

she

he

thev

• The **Skip-gram model** solves this problem more efficiently [Mikolov 2013]

 $\sum \sum \log p(w_c | w_t)$  $t=1 \ c \in \mathcal{C}_t$ 

Harris (1954): **Distributional Structure.** WORD Vol. 10, No. 2–3 Sahlgren (2008): **The Distributional Hypothesis.** Italian Journal of Linguistics 20, 33–54. Bengio et al. (2003): **A Neural Probabilistic Language Model.** Journal of Machine Learning Research 3. Mikolov et al. (2013): **Efficient Estimation of Word Representations in Vector Space.** arXiv:1301.3781

Sebastian Arnold 17

# 3. TASTY – A Robust Model for Efficient Entity Linking

<u>S. Arnold</u>, F. A. Gers, T. Kilias and A. Löser (2016). Robust Named Entity Recognition in Idiosyncratic Domains. **arXiv:1608.06757 [cs.CL]** 

<u>S. Arnold</u>, R. Dziuba and A. Löser (2016). TASTY: Interactive Entity Linking As-You-Type. **COLING 2016 (Demos)** 

# Challenge: Recognize and Link Named Entities

RQ1: What are general solutions to identify named entities in domain-specific text?

- → broad coverage
- → domain and languageindependent architecture
- → focus on **recall**
- → high robustness
- → efficient training

Expression of the chemokine receptor BLR2/EBI1 is specifically Epstein-Barr virus nuclear antigen 2. In our attempt to identify ch are related to Burkitt's lymphoma receptor 1 (SLR1) and are ex lymphocytes we used RT-PCR resulting in the isolation of a cDN transmembrane receptor termed BLR2. The protein shows signif similarities to the family of G-protein coupled chemokine receptor identical to the recently described receptor EBI1. Northen, blot a BLR2 mRNA could be highly stimulated in mitogen- and anti-Cr blood lymphocytes. BLR2-specific mRNA could be detected in a positive B cell lines. We show that transcription of the BLR2 gen induced in Epstein-Barr virus negative BL 41 cells via estrogen-r Epstein-Barr virus nuclear antigen 2, a key regulator of viral and immortalized B cells. Our data suggest an involvement of BLR2 migration in activated lymphocytes and in viral pathogenesis.

\* ° ° ` `

## Epstein–Barr virus nuclear antigen 2

The Epstein–Barr virus nuclear antigen 2 (EBNA-2) is a viral protein associated with the Epstein–Barr virus. EBNA-2 is the main viral

# Extracting entities requires local, contextual and global features.

# **Encoding Local Sub-word Information with Letter-trigrams**



Character-based word representations [Huang 2013] are a key component for robustness.

# **Capturing Sentence Context with Sequence Learning**



Bidirectional LSTM [Hochr. 1997, Graves 2012] effectively captures long-range dependencies.

Sebastian Arnold 21

# Key Results: TASTY

- TASTY is a **robust end-to-end model** for Entity Linking (NER+NEL)
- **Pushes down the costs** for a domainspecific model to labeling 4K sentences
- State-of-the-art performance for 20 English and German NER: 91.1% F1 on 0 CoNLL 2003 beats all models before 2016



## Next steps on our journey towards Neural Machine Reading:

• We envision a general document understanding beyond isolated entities

F1

• Model training should require as little supervision as possible

# 4. SECTOR – Coherent Topic Segmentation and Classification

<u>S. Arnold</u>, R. Schneider, P. Cudré-Mauroux, F. A. Gers and A. Löser (2019). SECTOR: A Neural Model for Coherent Topic Segmentation and Classification. **TACL Vol. 7**: 169–184 (presented at ACL)

# **Challenge: Understand Topics and Structure of a Document**

**RQ2:** How can Machine Reading models **detect topics and structure** in long documents?

- → topical information
- → structural information
- → coherent predictions
- $\rightarrow$  sentence granularity

Distributional Hypothesis needs to be extended with long-range structural information.



24

# Latent Topics can be Learned from Wikipedia Headings

signs

cause

types

risk

## Contents [hide]



- 1.2 Complications
- 2 Causes
  - 2.1 Type 1
  - 2.2 Type 2
  - 2.3 Gestational diabetes
  - 2.4 Maturity onset diabetes of the young
  - 2.5 Other types
- 3 Pathophysiology
- 4 Diagnosis
- 5 Prevention
- 6 Management
  - 6.1 Lifestyle
  - 6.2 Medications
  - 6.3 Surgery
  - 6.4 Support
- 7 Epidemiology
- 8 History
- 8.1 Etymology 9 Society and culture
  - 9.1 Naming
- 10 Other animals
- 11 Research
- 12 References
- 13 Further reading
- 14 External links



## preprocessing and pruning

(from 8.5k 'noisy' labels, bag-of-words) en disease (1.5k) diagnosis treatment symptoms causes epidemiology management prognosis history pathophysiology classification prevention genetics research culture factors mechanism society presentation differential disease surgery therapy clinical pathogenesis

de disease (1.0k) therapie symptome behandlung diagnose ursachen klinische ursache diagnostik verbreitung erscheinungen einteilung geschichte epidemiologie differentialdiagi prognose klassifikation formen verlauf haeufigkeit vorbeugung aetiologie pathogenese klinisches pathologie bild klinik entstehung symptomatik



clustering / normalization via BabelNet [Navigli 2012]

(85-95%) coverage)

de disease (25) en disease (27) treatment therapie symptom diagnose diagnosis symptom ursache cause kategorisierung classification epidemiology verlauf history geschichte prognosis management prognose pathophysiology praevalenz mechanism prevention fauna research genetics pathologie tomography definition klinik culture etymology infection genetik infektion fauna risk risiko pathology forschung geographie surgery screening mensch

medication

geography

other

complication

epidemiologie vorbeugung terminologie komplikation organe sonstiges

# Capturing Entire Documents with Sentence Embeddings + BLSTM

**Objective: predict topics on sentence level** 

- Bloom filter sentence compression [Serra 2017]
- Single-label classification into 25–30 topics
- Multi-label classification up to 603 words
- Segmentation based on deviation of the hidden layer topic embedding





# SECTOR Prediction on Par with Wiki Authors for "Dermatitis"



Source: https://en.wikipedia.org/w/index.php?title=Atopic dermatitis&diff=786969806&oldid=772576326

# Key Results: SECTOR

- Extends distributional models (top) with **topic and structure information** (bottom)
- Classifies 25-30 local topics with up to 71.6% F1
- **Detects topic shifts** to segment long documents into coherent passages



## Next steps on our journey towards Neural Machine Reading:

- We need to cover complementary information from named entities
- Hard passage segmentation should be replaced by continuous representations

# 5. CDV – Contextual Document Representations for Answer Retrieval

<u>S. Arnold</u>, B. van Aken, P. Grundmann, F. A. Gers and A. Löser (2020). Learning Contextualized Document Representations for Healthcare Answer Retrieval. **WWW 2020**: 1332–1343

# **Challenge: Capture Document Discourse Structure**

# RQ3: How can we embed discourse structure into document representations?

Pathophysiology In these finese influentation and hypertrophy of the relineasian sheath reservoirshy position the motion of th of a locking digit is not unique to trigger finger, and can be associated with dislocation, Duplocal dystonia, flexor tendon/sheath tumor, sesamoid bone anomalies, post-traumatic tendon en arpal head, and even hysteria. The differential diagnosis of pain at the MCP joint includes de Quer (for trigger thumb only), ulnar collateral ligament injury of the thumb (gamekeeper's thumb), MCP extensor apparatus injury, and MCP osteoarthritis.

# Entities and aspects provide discourse information complementary to language models.



# **Conclusion Conclusion Co**

Trigger finger is a long recognized condition characterized by a sometimes painful locking of the digit on flexion and extension. It is caused by the inflammation and subsequent narrowing of the A1 pulley through which the flexor

"Cystic

fibrosis"

# **Covering Diseases and Health Problems in Vector Space**

We create **entity embeddings** for diseases by training a Fasttext [Bojan. 2017] and BLSTM [Palangi 2016] model to map entity descriptions to entity IDs:

- Descriptions from UMLS, GARD, DiseasesDB, Wikidata and Wikipedia
- Covers 27,000+ diseases
- Provides fallback encoding for unseen entity names



# **Distributed Representation of Clinical Aspects**

We extend SECTOR with distributed representations of clinical aspects, e.g. symptoms, treatment, causes, diagnosis, prognosis, prevention, risk factors, ...

• Covers 14,000+ aspects learned from headings of medical Wikipedia articles



# Training CDV with Multi-task Objective

• Goal: Learn **discourse embedding** to align entities  $\epsilon_{1.T}$  and aspects  $\alpha_{1.T}$  with all sentences  $s_{1.T}$  in a long document:  $\mathcal{L}_{cdv}(\Theta) = \frac{1}{T} \sum_{t=1}^{T} (\|\hat{\epsilon}_t - \epsilon_t\| + \|\hat{\alpha}_t - \alpha_t\|)$ 

<b>s</b> s01 s2 s35 s6 s78 s9 s10 s11 s12 s1323 s2426 s2728 s2931	E(s) Arteriosclerosis Arteriosclerosis Arteriosclerosis;Arteriolosclerosis Arteriosclerosis;Atherosclerosis Arteriosclerosis;Monckeberg's_arter. Arteriosclerosis;Monckeberg's_arter. Arteriosclerosis Arteriosclerosis;Hyaline_arterioscl. Arteriosclerosis Arteriosclerosis Arteriosclerosis Arteriosclerosis Arteriosclerosis Arteriosclerosis Arteriosclerosis Arteriosclerosis	A(s) description signs; symptoms pathophysiology pathophysiology pathophysiology pathophysiology pathophysiology diagnosis treatment epidemiology history society; culture
		collect, our care



# **Contextualized Document Representation**

- Sentence embeddings  $\sigma_t$  are preprocessed using BioBERT [Lee 2019] or Fasttext [Bojan. 2017]
- Distributional models cannot capture structural context from long documents [Arnold 2019]
- Therefore, we transform  $\sigma_{1.T}$  into contextual discourse vectors  $\delta_{1.T}$  by using BLSTM [Hochr. 1997]  $\vec{h}_t = \text{LSTM}_{\Theta}(\vec{h}_{t-1}, \sigma(s_t))$  $\tilde{h}_t = \text{LSTM}_{\Theta}(\vec{h}_{t+1}, \sigma(s_t))$  $\delta_t = \tanh(W_{he}(\vec{h}_t \oplus \vec{h}_t) + b_e)$ 
  - $\hat{\epsilon}_t = \tanh(W_{\delta\epsilon}\delta_t + b_{\epsilon})$  $\hat{\alpha}_t = \tanh(W_{\delta\alpha}\delta_t + b_{\alpha})$



# **Loss Function**

Our loss function needs to handle the large variance of documents.

We optimize model parameters  $\theta$  by minimizing Huber loss [Huber 1992]





# Challenge: Retrieve Answer Passages from Long Documents

RQ4: How effective are document representations for **retrieving answer passages**?



**Cystic fibrosis** (CF) is a genetic disorder that causes mucus to build up and damage organs in the body, particularly the lungs and pancreas. Signs and symptoms may include salty-tasting skin; persistent coughing; frequent lung infections; wheezing or shortness of breath; poor growth; weight loss; greasy, bulky stools; difficulty with bowel movements; and in males, infertility. Over time, mucus buildup and infections can lead to permanent lung damage, including the formation of scar tissue (fibrosis) and cysts in the lungs. CF is caused by mutations in the *CFTR* gene and inheritance is autosomal recessive.<sup>[1][2][3]</sup> Treatment aims to relieve symptoms and usually includes respiratory therapies, inhaled medicines, pancreatic enzyme supplement, and nutritional supplements. Newer medications such as CFTR modulators have been approved for use in the United States. Ongoing research is focused on finding a cure for the disease.<sup>[2]</sup>

# Answer Retrieval on Precomputed CDVs



score<sub>t</sub> 
$$(Q(E, A), C_D) = \operatorname{cosine}(\epsilon_Q \oplus \alpha_Q, \hat{\epsilon}_t \oplus \hat{\alpha}_t)$$

→ No query/answer pairs are required to train the model

 $\rightarrow$  CDV vectors  $C_D$  need to be precomputed only once

# **Results for Healthcare Answer Retrieval**



We evaluate the retrieval of answer passages from over 2K healthcare articles from Wikipedia, NIH (including rare diseases) and Patient.info

\* interaction-based models

# Key Results: Contextual Discourse Vectors

- CDV extends pre-trained Language Models with long-range discourse information
- Trained with self-supervision on Wikipedia data
- Pre-computed vectors enable efficient retrieval with high recall.
- Outperforms BM25 (Elastic Search), DSSM, Duet (Microsoft), MatchPyramid, Conv-KNRM, HAR and others with up to 65.2% R@1
- Potential errors caused by hierarchical, related and overlapping information

## **Our journey towards Neural Machine Reading:**

- CDV fulfills all properties we have defined for automatic language understanding.
- This model satisfies our hypothesis of Neural Machine Reading.



linical manifestations



# TASTY: Tag-as-you-type Entity Linking



S. Arnold, F. A. Gers, T. Kilias and A. Löser (2016). Robust Named Entity Recognition in Idiosyncratic Domains. arXiv:1608.06757 [cs.CL] S. Arnold, R. Dziuba and A. Löser (2016). TASTY: Interactive Entity Linking As-You-Type. COLING 2016 Demos

# **TraiNER: Bootstrapping Named Entity Recognition**



# Smart-MD: Clinical Decision Support System



R. Schneider, <u>S. Arnold</u>, T. Oberhauser, T. Klatt, T. Steffek and A. Löser (2018). **Smart-MD: Neural Paragraph Retrieval of Medical Topics.** WWW 2018 (Companion): 203–206

# **CDV Healthcare Retrieval**

## **CDV Search | Highlight**

Find passages 1 ss a specific	topic:
COVID-19 medication	Search

#### Similar Diseases

2019 novel coronavirus respiratory syndrome

0

- Canine infectious tracheobronchitis
- Adenoviridae infectious disease
- Avian infectious bronchitis
- Human parainfluenza
- SARS-CoV-2

Dataset: 11.1K articles from the COVID-19 Open Research Dataset (PMC Open Access subset) All articles are made available under a Creative Commons or similar license. Specific licensing information for individual articles can be found in the PMC source and CORD-19 metadata. More datasets: Wikipedia I CORD-19



BEUTH HOCHSCHULE FÜR TECHNIK BERLIN University of Applied Sciences

Made by DATEXIS (Data Science and Text-based Information Systems) at Beuth University of Applied Sciences Berlin

Deep Learning Technology: Sebastian Arnold, Betty van Aken, Paul Grundmann, Felix A. Gers and Alexander Löser. Learning Contextualized Document Representations for Healthcare Answer Retrieval. The Web Conference 2020 (WWW'20)

Funded by The Federal Ministry for Economic Affairs and Energy; Grant: 01MD19013D, Smart-MD Project. **Digital Technologies** Gefördert durch



Results for Query ( COVID-19 medication )

### Potential Treatments for COVID-19: a Narrative Literature Review - Conclusion PMC Open Access (CC BY-NC-SA) | 88.68%

Apparently, in addition to the drugs currently prescribed to treat COVID-19, Arbidol hydrochloride, interferon, and Thalidomide plus Methylprednisolone can also be used due to their effects reported in clinical studies. However, more studies are needed to confirm the use of corticosteroids, as there are conflicting reports regarding their efficacy. Also, potential drugs listed in Table 1, such as Remdesivir, Atazanavir, Saguinavir, and Formoterol, and Tocilizumab can be introduced as treatments for COVID-19 if they prove to be effective in animal and clinical studies.

## The origin, transmission and clinical therapies on coronavirus diseased 19 (COVID-19) outbreak - an update on the status - Antiviral treatments

#### PMC Open Access (CC BY-NC-SA) | 88.30%

Based on the experience of fighting the epidemic SARS-CoV and MERS-CoV previously, we may learn some lessons for some treat strategies against coronavirus. Antiviral drugs and systemic corticosteroid treatment practice previously, including neuraminidase inhibitors (oseltamivir, peramivir, commonly used in 2 zanamivir, etc), gand vclovir, and ribavirin, as well as methylprednisolone [46, 75] for influenza virus, are invalid for COVID-19 and not recommended. Remdesivir (GS-5734) is a 1'-cyano-substituted adenosine nucleotide analog prodrug and shows broad-spectrum antiviral activity against several RNA viruses. Based on

### The origin, transmission and clinical therapies on coronavirus disease 2019 (COVID-19) outbreak - an update on the status - Current therapies PMC Open Access (CC BY-NC-SA) | 87.86%

Given the lack of effective antiviral therapy against COVID-19, current treatments mainly focused on symptomatic and respiratory support according to the Diagnosis and Treatment of Pneumonia Caused by COVID-19 (updated to version 6) issued by National Health Commission of the People's Republic of China. Nearly all patients accepted oxygen therapy, and WHO recommended extracorporeal membrane oxygenation (ECMO) to patients with refractory hypoxemia. Rescue treatment with convalescent plasma and immunoglobulin G are delivered to some critical cases according to their conditions.

COVID-19: what has been learned and to be learned about the novel coronavirus disease - Tackling cytokine storms

#### PMC Open Access (CC BY-NC-SA) | 87.47%

It has been known that a cytokine storm results from an overreaction of the immune system in SARS and MERS patients 33. Cytokine storm is a form of systemic inflammatory response featured by the release of a series of cytokines including TNFa, IL-1B, IL-2, IL-6, IFNa, IFNB, IFNY, and MCP-1. These cytokines induce immune cells

## Highlight for Query (COVID-19 medication)

#### COVID-19: what has been learned and to be learned about the novel coronavirus disease

#### PMC Open Access (CC BY 4.0)

#### Introduction

The Spring Festival on January 25, 2020 has become an unprecedented and unforgettable memory to all Chinese who were urged to stay indoors for all the holiday and for many weeks after due to the outbreak of a novel viral disease. The virus is highly homologous to the coronavirus (CoV) that caused an outbreak of severe acute respiratory syndrome (SARS) in 2003; thus, it was named SARS-CoV-2 by the World Health Organization (WHO) on February 11, 2020, and the associated disease was named CoV Disease-19 (COVID-19) 1. The epidemic started in Wuhan, China, and guickly spread throughout the entire country and to near 50 others all over the world. As of March 2, 2020, the virus has resulted in over 80,000 confirmed cases of COVID-19, with more than 40,000 patients discharged and over 3,000 patients who died, WHO warns that COVID-19 is "public enemy number 1" and potentially more powerful than terrorism 2.

#### Nuclear acid assays

The detection of SARS-CoV-2 RNA via reverse-transcriptase polymerase chain reaction (RT-PCR) was used as the major criteria for the diagnosis of COVID-19. However, due to the high false-negative rate, which may accelerate the epidemic, clinical manifestations started to be used for diagnosis (which no longer solely relied on RT-PCR) in China on February 13, 2020. A similar situation also occurred with the diagnosis of SARS 59. Therefore, a combination of disease history, clinical manifestations, laboratory tests, and radiological findings is essential and imperative for making an effective diagnosis. On February 14, 2020, the Feng Zhang group described a protocol of using the CRISPR-based SHERLOCK technique to detect SARS-CoV-2, which detects synthetic SARS-CoV-2 RNA fragments at 20 × 10-18 mol/L to 200 × 10-18 mol/L (10-100 copies per microliter of 5 nience if verified in clinical samples. input) using a dipstick in less than an hour without premiring elaborate instrumentation 60. Hopefully, the new technique can dramatically enhance the sensitivity Antiviral therapy

At the time of writing, no effective antiviral therapy confirmed. However, intravenous administration with remdesivir, a nucleotide analog, has been found to be efficacious in an American patient with COVID-19 64. Remdesivir is a novel antiviral drug developed by Gilead initially for the treatment of diseases caused by Ebola and Marlburg viruses 76. Later, remdesivir also demonstrated possible inhibition of other single stranded RNA viruses including MERS and SARS viruses 77.78. Based on these, Gilead has provided the compound to China to conduct a pair of trials on SARS-CoV-2-infected individuals 79, and the results are highly anticipated. In addition, baricitinb, interferon-q, lopinavir/ritonavir, and ribavirin have been suggested as potential for patients with acute respiratory symptoms 80,81. Diarrhea, nausea, vomiting, liver damage, and reactions can occur following combined therapy with lopinavir/ritonavir 80. The interaction of the with other drugs used in the patients should be monitored carefully.

#### Plasma from recovered patients and antibody generation

The collection of the blood from patients who recovered from a contagious disease to t suffering from the same disease or to protect healthy individuals from catching the disease ha Indeed, recovered patients often have a relatively high level of antibodies against the patho Antibodies are an immunoglobulin (Ig) produced by B lymphocytes to fight pathogens and oth and they recognize unique molecules in the pathogens and neutralize them directly 83 Based was collected from the blood of a group of patients who recovered from COVID-19 and was in. seriously ill patients. Their symptoms improved within 24 hours, accompanied by reduced inflammatic. loads and improved oxygen saturation in the blood. However, verification and clarification are necess propose the method for large-scale use before specific therapies are not yet developed.

**Online:** https://cord19. cdv.demo. datexis.com

J.-M. Papaioannou, S. Arnold, F. A. Gers, A. Löser, M. Mayrdorfer, and K. Budde (2020). Aspect-Based Passage Retrieval with Contextualized Discourse Vectors. Submitted to: ECIR 2021 System Demonstrations

# 7. Conclusion and Future Work

# Contributions

# TASTY, SECTOR and CDV are Neural Machine Reading architectures for domain-specific language understanding:

- Variety of languages and domains
- Robust against noise and spelling variations
- Document and discourse structure
- Entity Linking, Topic Modeling, Answer Retrieval tasks
- Efficient to train
- In-depth error analysis

Implementation available under Apache V2 license <a href="https://github.com/sebastianarnold/TeXoo">https://github.com/sebastianarnold/TeXoo</a>

# Limitations

## Writing systems

- All experiments in Latin script
- Sentence and token normalization necessary

# **Self-supervision**

- Preprocessing tailored to training data
- Holistic approach: Weak-supervision

# **Optimization objectives**

- We focused on task accuracy (Prec/Rec/F1)
- Reality: result diversification, freshness, trust, popularity, etc.
- Dynamically changing objectives

# **Future Work**

**Hierarchical Knowledge Representations** 

- Capture hypernym/hyponym relations
- Model **semantic hierarchies** in vector space

# **Uncertainty Modeling**

- **Data distribution** changes from training to inference time
- Report level of confidence

**Continual Learning** 

- Continuously learn from new data
- Prevent catastrophic forgetting

# Acknowledgements

Thanks to everyone who supported me during this time with inspiration, feedback and grounding.

- Supervisors, Mentors and Committee Alexander Löser, Philippe Cudré-Mauroux, Felix Gers, Peter Tröger, Amy Siu, Laura Dietz
- Coauthors and Team Members

Torsten Kilias, Robert Dziuba, Rudolf Schneider, Christopher Kümmel, Robin Mehlitz, Tom Oberhauser, Betty van Aken, Benjamin Winter, Paul Grundmann, Michalis Papaioannou DATEXIS Team @ BeuthHS eXascale Infolab Team @ UNIFR

• My Family and Friends



# **Questions & Discussion**

Our work is funded by the German Federal Ministry of Economic Affairs and Energy (BMWi) under grant agreements 01MD16011E (Medical Allround-Care Service Solutions) 01MD19003b (PLASS), 01MD19013D (Smart-MD) and H2020 ICT-2016-1 grant agreement 732328 (FashionBrain).

Gefördert durch:



Bundesministerium für Wirtschaft und Energie



aufgrund eines Beschlusses des Deutschen Bundestages





# **Neural Machine Reading for Domain-Specific Text Resources**

**Speaker: Sebastian Arnold** 

sarnold@beuth-hochschule.de @sebastianarnold

Data Science and Text-based Information Systems (DATEXIS) Beuth University of Applied Sciences

Berlin, Germany www.datexis.de



# Selected References (1)

Bengio, Ducharme, Vincent and Janvin (2003): A Neural Probabilistic Language Model. *Journal of Machine Learning Research* 3. Bojanowski, Grave, Joulin and Mikolov, (2017): Enriching Word Vectors with Subword Information. TACL Vol. 5 Devlin, Chang, Lee and Toutanova (2019): BEPT: Pre-Training of Deep Ridirectional Transformers for Language Understanding

Devlin, Chang, Lee and Toutanova (2019): BERT: Pre-Training of Deep Bidirectional Transformers for Language Understanding. NAACL-HLT 2019

Etzioni, Banko and Cafarella (2006): Machine Reading. AAAI 2006

Firth (1957): A Synopsis of Linguistic Theory, 1930–1955. Studies in Linguistic Analysis

Graves (2012): Supervised Sequence Labelling with Recurrent Neural Networks. Berlin Heidelberg: Springer

Harris (1954): Distributional Structure. WORD Vol. 10, No. 2-3

Hochreiter and Schmidhuber (1997): Long Short-term Memory. Neural Computation, Vol. 9, No. 8

Huang, He, Gao, Deng, Acero and Heck (2013): Learning Deep Structured Semantic Models for Web Search using Clickthrough Data. CIKM 2013

Huber (1992): Robust Estimation of a Location Parameter. Breakthroughs in Statistics, Springer

Jurafsky and Martin (2008): Speech and Language Processing. Vol. 2

Lee et al. (2019): BioBERT: A Pre-trained Biomedical Language Representation Model for Biomedical Text Mining. *Bioinformatics* Manning, Schütze and Raghavan (2008): Introduction to Information Retrieval. Cambridge University Press Marchionini (2006): Exploratory Search: From Finding to Understanding. *Communications of the ACM* Vol. 49, No. 4 Mikolov, Chen, Corrado and Dean (2013): Efficient Estimation of Word Representations in Vector Space. arXiv:1301.3781 Mitra, Diaz and Craswell (2017): Learning to Match using Local and Distributed Representations of Text for Web Search. WWW 2017 Moro, Raganato and Navigli (2014): Entity Linking meets Word Sense Disambiguation: A Unified Approach. *TACL* Vol. 2

# Selected References (2)

Navigli and Ponzetto (2012): BabelNet: The Automatic Construction, Evaluation and Application of a Wide-coverage Multilingual Semantic Network. Artificial Intelligence, Vol. 193

Palangi et al. (2016): Deep Sentence Embedding using Long Short-Term Memory Networks: Analysis and Application to Information Retrieval. TASLP Vol. 24, No. 4

Peters, Neumann, Iyyer, Gardner, Clark, Lee and Zettlemoyer (2018): Deep Contextualized Word Representations. NAACL-HLT 2018 Pink, Nothman and Curran (2014): Analysing Recall Loss in Named Entity Slot Filling. EMNLP 2014

Prokofyev, emartini and Cudré-Mauroux (2014): Effective Named Entity Recognition for Idiosyncratic Web Collections. WWW 2014 Rose and Levinson (2004): Understanding User Goals in Web Search. WWW 2004

Sahlgren (2008): The Distributional Hypothesis. Italian Journal of Linguistics 20

Sarawagi (2008): Information Extraction. Foundations and Trends in Databases 1, Nr. 3

Serrà and Karatzoglou (2017): Getting Deep Recommenders Fit: Bloom Embeddings for Sparse binary Input/Output Networks. RecSys 2017

Vaswani, Shazeer, Parmar, Uszkoreit, Jones, Gomez, Kaiser and Polosukhin (2017): Attention Is All You Need. NIPS 2017 Zhu, Ahuja, Wei and Reddy (2019): A Hierarchical Attention Retrieval Model for Healthcare Question Answering. WWW 2019