

Bringing Structure into Summaries

- A Faceted Summarization Dataset for Long Scientific Documents

Rui Meng, Khushboo Thaker, Lei Zhang, Yue Dong, Xingdi Yuan, Tong Wang, Daqing He



- Divide document into a linear sequence of subtopics
- Present ideas and arguments progress in a logical and orderly manner

- Divide document into a linear sequence of subtopics
- Present ideas and arguments progress in a logical and orderly manner
- Examples
 - Academic: IMRaD Introduction, Methods, Results, and Discussion



- Divide document into a linear sequence of subtopics
- Present ideas and arguments progress in a logical and orderly manner
- Examples
 - Academic: IMRaD Introduction, Methods, Results, and Discussion
 - Newspapers: **5W1H**, Inverted Pyramid Style

	The Lead
	The most important information about an event
/	Who? What? Where? When? Why? How?
1	The Body
	The crucial information expanding the topic
	Argument, Controversy, Story, Evidence,
	Background details
	The Tail
	Extra information
	Interesting, related
	items.
	Journalist
	Assessment
	\checkmark

- Divide document into a linear sequence of subtopics
- Present ideas and arguments progress in a logical and orderly manner
- Examples
 - Academic: IMRaD Introduction, Methods, Results, and Discussion
 - Newspapers: **5W1H**, Inverted Pyramid Style
 - Patient Report: **SOAP** subjective, objective, assessment, plan

COPD/pr	eumonia
	GOALS
1. Pt. will 2. Pt. will surfaces.	demonstrate productive cough in seated position, 3/4 triats. ambulate 150ft with supervision, no assistive device, on level indoor
S	Pt. reports not feeling well today, "I'm very tired".
0	Auscultation findings: scattered rhonchi all lung fields. Chest PT was performed in sitting (ant. and post.). Techniques included percussion, vibration, and shaking. Pt. performed a weak combined abdominal and upper costal cough that was nonbronchospastic, congested, and non-productive. The cough/huff was performed with VC. Pectoral stretch/thoracic cage mobilizations performed in seated position Pt. given towel roll placed in back of seat to open up ant. chest wall. Strengthening exercises in standing - pt. performed hip flexion, extension, and abduction; knee flexion 10 reps x 1 set B. Pt. performs HEP with supervision (in evenings with wife). Pt. instructed to hold tissue over trach when speaking to prevent infection and explained importance of drinking enough water.
A	Pt. continues to present with congestion and limitations in coughing productivity. Pt. has been compliant with evening exercise program, which has results in increased to to therapeutic exercise regime and an increase in LE strength. Amb. not attempted to 20 to pt. report of fatigue Pt. should be able to tolerate short distance ambulation within the next few days.
Ρ	Cont. current exercise plan including CPT; emphasize productive coughing techniques; increase strengthening exer reps to 15; attempt

Summary in Academic Writing

- Traditional abstracts
 - A concise summary of a research paper
 - Single paragraph of approximately 200~300 words in length
 - Highlight key content aspects: purpose/importance of the research, main outcomes etc.

Summary in Academic Writing

- Traditional abstracts
 - A concise summary of a research paper
 - Single paragraph of approximately 200~300 words in length
 - Highlight key content aspects: purpose/importance of the research, main outcomes etc.
- Can we bring structure into abstracts?
 - Structured abstracts
 - Emerald (this study)
 - PubMed (Gidiotis and Tsoumakas, 2019)
 - Benefits
 - Faster and easier access by the reader
 - Greater clarity for both reader and writer

Structured Abstract of Emerald

- FacetSum Dataset
 - 60,532 scientific papers w/ fulltext
 - Covering 25 different academic fields

# documents Train: 46,289 / Dev: 6,000 / Test: 6,000 / OA-Test: 2,243										
# words in abstracts										
Full Purpose Method Findings Value										
mean	290.4	54.1	52.0	68.6	47.3					
std	± 82.8	± 28.4	± 27.8	± 32.4	± 24.2					
# words in paper sections										
	Full Intro. Method Result Conc.									
recall%	-	84.3%	67.0%	72.4%	79.0%					
mean	6,827	885	1,194	2,371	747					
std	±2,704	± 557	±861	±1,466	±567					



- Marketing
- Human Resource Management
- Information and Knowledge Management
- Industry and Public Sector Management
- Education
- Managing Quality
- Economics
- Learning and Development
- Operations and Logistics Management
- Library and Information Science
- Built Environment
- Health Care Management / Healthcare

- Accounting and Finance
- Enterprise and Innovation
- Tourism and Hospitality
- Organization Studies
- Strategy
- Business Ethics and Law
- International Business
- Performance Management and Measurement
- Environmental Management / Environment
- Health and Social Care
- Management Science / Management Studies
- Mechanical and Materials Engineering

Structured Abstract of Emerald

- Summarize each study in four facets
 - Purpose
 - Design/methodology/approach
 - Findings
 - Originality/value

emerald insight

Discover Journals, Books & Case Studies

Home / Journals / Journal of Industry - University Collaboration / Volume 1 Issue 1 / End-to-end learning via a convolutional neural network for cancer cell line classification

End-to-end learning via a convolutional neural network for cancer cell line classification

Darlington A. Akogo , Xavier-Lewis Palmer Journal of Industry - University Collaboration ISSN: 2631-357X Article publication date: 12 April 2019 Issue publication date: 12 April 2019

Abstract

Purpose – Computer vision for automated analysis of cells and tissues usually include extracting features from images before analyzing such features via various machine learning and machine vision algorithms. The purpose of this work is to explore and demonstrate the ability of a Convolutional Neural Network (CNN) to classify cells pictured via brightfield microscopy without the need of any feature extraction, using a minimum of images, improving work-flows that involve cancer cell identification. **Design/methodology/approach** – The methodology involved a quantitative measure of the performance of a Convolutional Neural Network in distinguishing between two cancer lines. In their approach, they trained, validated and tested their 6-layer CNN on 1,241 images of MDA-MB-468 and MCF7 breast cancer cell lypes **Findings** – They obtained a 99% accuracy, providing a foundation for more comprehensive systems.

Originality/value – Value can be found in that systems based on this design can be used to assist cell identification in a variety of contexts, whereas a practical implication can be found that these systems can be deployed to assist biomedical workflows quickly and at low cost. In conclusion, this system demonstrates the potentials of end-to-end learning systems for faster and more accurate automated cell analysis.

Keywords End-to-end learning, Convolutional neural network, Cancer cell line classification Paper type Research paper

Structured Abstract of PubMed

- Five facets
 - \circ Introduction
 - Objections
 - Methods
 - Results
 - Conclusion
- Dataset
 - 19.7m+ scientific papers, 2.8m w/ fulltext
 - #unstructured : #structured = 4.5m : 15.1m

National Institutes of Health	M RSS	Save search Advanced	
Display Settings: 🕑 Abstra	ct	Send t	<u>o:</u> 🕑
Clin Toxicol (Phila), 2014 Jun;52(5	i):525-30. doi: 10.3109/15563650.20	14.913175. Epub 2014 May 5.	
Evaluation of dexmed toxicological events a	etomidine therapy for s t an academic medical	edation in patients with center.	
Mohorn PL ¹ , Vakkalanka JP, R	ushton W, Hardison L, Woloszyn /	A, Holstege C, Corbett SM.	
Author information			
Abstract INTRODUCTION: Although (agonist, has increased, its ro sequelae has not been well (linical use of dexmedetomidine le in patients admitted to intens astablished.	e (DEX), an alpha2-adrenergic rece sive care units secondary to toxicolo	ptor ogical
OBJECTIVES: The primary observed in poisoned patient	objective of this study was to de ts receiving DEX for sedation.	escribe clinical and adverse effects	
METHODS: This was an obs patients who received DEX for incidence of adverse effects arrhythmias. For comparison and every hour during DEX to duration; time within target R analgesia, and vasopressor	ervational case series with retr or sedation at an academic me of DEX therapy including brady , vital signs were collected hou herapy until the therapy ended, ichmond Agitation Sedation Sc requirements.	rospective chart review of poisoned dical center. The primary endpoint v vcardia, hypotension, seizures, and irly for the 5 h preceding the DEX th . Additional endpoints included thera core (RASS); and concomitant seda	was erapy apy tion,
RESULTS: Twenty-two patie similar to the commonly used vs. 93 beats/minute, p < 0.05 (111 vs. 109 mmHg, p = 0.74 during therapy. No additional duration of therapy was 6.5 a of other sedation and/or anal DEX initiation. Seven patient patients requiring vasopress	nts were included. Median initi I rates for sedation. Median here i). Median systolic blood pressu (5). Five patients experienced a l adverse effects were noted. M and 44.5 h, respectively. Seven gesia with four (23%) of these (32%) had concomitant vaso or support after DEX initiation.	al and median DEX infusion rates w art rate was lower during the therap ure before and during therapy was s an adverse effect per study definitio fledian time within target RASS and teen patients (77%) had concomita patients requiring additional agents pressor support with four (57%) of th	/ere y (82 similar ns nt use after hese
CONCLUSION: Common ad vasopressor support during t patients. Larger, comparative recommended in poisoned p	verse effects of DEX were note herapy warrants further investig a studies need to be performed atients	ed in this study. The requirement for gation into the safety of DEX in pois I before the use of DEX can be routi	oned

-1 24702790[uid1]

PubMed

Pub Med

Analysis 1 - Sentence Position Distribution

- Find best-match sentences (extractive oracle) in source texts by Rouge-1
- Plot the relative position of each oracle sentence
 - **Purpose** is clearly skewed towards the beginning
 - Method distribution is generally uniform
 - **Findings** are mostly positioned towards the end
 - Value is located at the beginning and end



Findings

Analysis 1 - Sentence Position Distribution

- Find best-match sentences (extractive oracle) in source texts by Rouge-1
- Plot the relative position of each oracle sentence
 - **Purpose** is clearly skewed towards the beginning
 - Method distribution is generally uniform
 - **Findings** are mostly positioned towards the end
 - Value is located at the beginning and end



Analysis 2 - Section Alignment

- How different abstract facets align with different sections in an article
- Categorize paper sections by keyword matching of section titles
 - Introduction: intro, purpose
 - **Method**: design, method, approach
 - **Result**: result, find, discuss, analy
 - **Conclusion**: conclu, future
- Oracle scores show a strong correlation between summary facets and sections

Paper Section	n		1.1sior			
Abstract	Full	1+C	Introduc	Method	Result (Conclus
Full	62.09	56.47	48.47	43.32	49.73	50.42
Purpose	49.76	47.06	44.23	30.12	33.87	36.23
Method	45.36	34.23	30.82	37.53	29.07	28.46
Findings	52.09	45.28	33.65	29.49	42.80	42.35
Value	45.98	42.37	35.29	26.68	32.52	36.85

Automatic Faceted Summarization

- Method
 - Utilize a pre-trained encoder-decoder model BART
 - Fine-tune it on Emerald paper data
 - Input: (1) full paper (2) I+C (introduction+conclusion)



Automatic Faceted Summarization

- Method
 - Utilize a pre-trained encoder-decoder model BART
 - Fine-tune it on Emerald paper data
 - Input: (1) full paper (2) I+C (introduction+conclusion)
 - Output:
 - (1) **BART-CAT:** a long summary by concatenating all facets
 - Pros: simple to implement
 - Cons: target text might be very long, model can have difficulty handling long-distance dependency



Automatic Faceted Summarization

- Method
 - Utilize a pre-trained encoder-decoder model BART
 - Fine-tune it on Emerald paper data
 - Input: (1) full paper (2) I+C (introduction+conclusion)
 - Output:
 - (1) **BART-CAT**
 - (1) **BART-FACET:** use a control token in input specifying the facet to summarize
 - Pros: shorter output sequence
 - Cons: may output duplicate information among multiple facets



Results

- Results
 - BART-CAT << BART-Facet
 - The quality of outputs deteriorates rapidly as the decoding proceeds

Model	Source Text	Full	Purpose	Method	Findings	Value
BART (Lewis et al., 2020)	I+C	44.36	41.14	20.75	14.72	5.85
BART-Facet	I+C	47.09	43.47	29.07	30.97	28.90
BART	full paper	42.74	41.21	20.53	14.33	5.07
BART-Facet	full paper	45.76	42.55	28.07	28.98	28.70

Results

• Results

- BART-CAT << BART-Facet
 - The quality of outputs deteriorates rapidly as the decoding proceeds
- I+C > full paper
 - Difficulty in handling long inputs (avg 6.8k tokens).

Model	Source Text	Full	Purpose	Method	Findings	Value
BART (Lewis et al., 2020)	I+C	44.36	41.14	20.75	14.72	5.85
BART-Facet	I+C	47.09	43.47	29.07	30.97	28.90
BART	full paper	42.74	41.21	20.53	14.33	5.07
BART-Facet	full paper	45.76	42.55	28.07	28.98	28.70

Future Work

- Incorporating methods for long-document processing
 - Reducing input length by first extracting key sentences or segments
 - Architectures for long text, such as Longformer, BigBird
- Automatically structuring traditional abstracts

Thank You!

Email: rui.meng@pitt.edu

arXiv: arXiv: arxiv.org/pdf/2106.00130.pdf

Data & Code: <u>https://github.com/hfthair/emerald_crawler</u>

- Crawler for Emerald data
- Paper links and dataset splits

