

Multi-documentation Summarization of Science Articles

Justin To, Luka Liu, Milan Dean

W 266: Natural language Processing
UC Berkeley School of Information

April 15, 2023

Abstract

In today's world, with an abundance of information available, the skill of generating accurate and concise summaries from extensive texts is more crucial than ever before. This study primarily focuses on the task of abstractive multi-document summarization (MDS), which involves generating a summary based on multiple interrelated documents. To build an MDS model tailored to Multi-XScience, a dataset that the pre-trained models have not encountered before, we explored different strategies such as fine-tuning and model-stacking, using two state-of-the-art pre-trained large language models (LLMs), namely Centrum and LED. In evaluating the generated summaries, we employed both quantitative ROUGE-based scoring and qualitative analysis of the model outputs to gain insights into how the models adapted to specific elements of the task and dataset. Our objective is to contribute to research on adapting pre-trained models to new and domain-specific MDS tasks and datasets using different approaches, such as fine-tuning and model stacking, and identify elements that may affect the success of such approaches.

1 Introduction

MDS is the task of summarizing multiple texts into a concise and informative summary. Compared to single document summarization, it is challenging to keep the summary coherent and comprehensive, given the documents are longer¹ and of a more complex structure². Despite the challenges, MDS has wider application potential as real-world tasks often involve collating and summarizing information from multiple sources. An important question therefore is, given the advances in NLP and the availability of pre-trained LLMs, how should we build a customized abstractive MDS model if we are given a domain-specific dataset previously unseen by the models?

In this paper, we employ the Multi-XScience dataset, a relatively recent dataset that has not been widely used for pre-training/fine-tuning LLMs, to simulate the situation where we are given an MDS task for a broad domain, e.g. academic writing from the sciences. The defining characteristics of the dataset, namely (a) the significant portion of samples with long inputs (up to ~ 6300 tokens) and (b) that different articles within each sample are related yet not covering the exact same event³, led us to pick two models for experimentation. The first is the Longformer Encoder Decoder (LED) model, which was proposed by Beltagy et al. [1] and introduced the concept of local vs global attention so that attention-based transformers can still effectively handle the computation of long inputs. The second is the Centrum model, proposed by Puduppully and Steedman [6], for its capability of handling long inputs and inclusion of a centroid-based approach for dealing with multiple source documents.

Leveraging on the available pre-trained checkpoints, we first tested the off-the-shelf LED and

¹For instance, the widely used CNN/Daily Mail dataset has an average input length of around 780 tokens. In contrast, around 30% of the inputs in the Multi-XScience dataset has over 1024 tokens, while the longest sample has over 6300 tokens.

²In MDS, the source articles are written by multiple authors with varying writing styles and article structures.

³Further discussed in Section 3.1

Centrum using the Multi-XScience dataset and then proceeded to fine-tune them for improved performance. Then, we explored the use of a two-step setup whereby individual source articles of each sample are first shortened through summarization, and then combined again for MDS processing. For evaluation, we employed the common ROUGE score metric, showing the improved performance from fine-tuning (ROUGE-2 from 5.2 to 6.9; ROUGE-L from 14.6 to 17.8, which are close to state-of-the-art (SOTA)). Further analyses were conducted on the variation of performance across inputs of different lengths, and the amount of “copying” from the main article. Lastly, a qualitative analysis was done on 25 samples for each model tested, and the fine-tuned and two-step models were found to demonstrate capabilities in extracting information from and contrasting different sources, while adapting its writing style to one suited for the dataset.

2 Related Works

In recent years, there has been significant research in the field of MDS, with a particular focus on the use of transformer-based models for abstractive models and through approaches such as sparse attention mechanisms, hierarchical sentence representations, and document-level clustering. Much progress has also been made in providing datasets for MDS training and evaluation.

On model development, the LED model proposed by Beltagy et al. [1] addresses the issue of processing long documents by introducing a sparse attention mechanism that allows the model to scale to documents of up to 16384 tokens. This approach has outperformed previous models, including advanced ones such as PEGASUS, on long-document tasks such as summarization.

Branching off to models more specialized in MDS, PRIMERA, proposed by Xiao et al. [cutexiao2022primera](#), is a pre-training method for transformer-based models that leverages hierarchical sentence representations to improve performance on MDS tasks. The approach involves

constructing a sentence-level pyramid structure and applying a masked language modeling objective to the pyramid, leading to improved results over BART, PEGASUS, and LED.

Advancing the research on MDS further, Puduppully and Steedman [6] proposed a centroid-based pre-training method for MDS that leverages document-level clustering to capture document-level semantics. This approach showed improved performance on summarization tasks over previous models including PRIMERA, particularly for MDS datasets with a large number of documents.

On the dataset front, a commonly used dataset in the MDS domain is perhaps the Multi-news dataset, which was introduced by Fabbri et al. [2] and contains news articles from multiple sources. Other MDS datasets include Wikisum, which was generated by Liu et al. [4] from Wikipedia articles, WCEP compiled by Ghalandari et al. [3] based on news summaries from the Wikipedia Current Events Portal, and the Multi-XScience dataset by Lu et al. [5] used in this study.

3 Methods

3.1 Dataset

We used the Multi-XScience MDS dataset [5] in this study, which is a collection of scientific articles from scientific fields, including computer science, physics, and biology. Some key parameters of the dataset are set out as follows, while Appendix A summarizes other findings and visualizations from our EDA, as well as the data preprocessing steps adopted for the experiments.

The dataset is chosen for two reasons. First, Multi-XScience is among the few datasets previously unseen by MDS models during pre-training⁴. Picking datasets already seen by the models could pose difficulties in evaluation, not to mention that it would go against the overall goal to explore MDS model-building strategies for unseen data.

⁴LED’s pre-training data included Multi-News while Centrum’s pre-training included both Multi-News and WikiSum.

	Training	Validation	Test
Sample size	30369	5066	5093
Input	899	898	886
Token Length	565	567	547
	82-4694	79-4183	131-6348
Samples with inputs > 1024 tokens	31.21%	30.71%	30.87%
Label	142	141	142
Token Length	58	58	58
	24-735	25-324	26-418
Articles per sample	4.43	4.43	4.39
	2.62	2.63	2.55
	2-21	2-20	2-19

Table 1: Dataset description. The cells with three columns refer to mean/standard deviation/range

The second reason is that Multi-XScience is a more challenging task that is more comparable to real-world data collation and summarization. Many of the datasets (e.g. Multi-news, WCEP) are news-based, with the articles within a sample covering the same news event and key elements to be included in the summary appearing in multiple sources. The Multi-XScience dataset however, expects the model to write the related works section of a journal paper based on (a) the abstract of that paper which serves to provide context; and (b) the abstracts of the other journal articles that the main paper referenced. While these additional abstracts are usually⁵ closely related, they often concern different aspects of an issue, and the model needs to be able to compare and contrast the similarities and differences⁶. This difficulty of the dataset also shows up empirically⁷ in [8].

⁵Some exceptions exist, e.g. sample 845 (refer to Appendix D for sample review for sample review)) of the test set includes an irrelevant article on the design attributes of luggage carriers when the main article is on a MDS over online product reviews. See Appendix C for details.

⁶In contrast, the commonly used Multi-news dataset has, in each sample, multiple news reports surrounding a single event, and the overall theme to be covered in the summary is generally present in every source article.

⁷For instance, the zero- and few-short evaluation of BART, PEGASUS, LED and PRIMERA yielded a ROUGE-2 score of only 1.9 to 4.6 and a ROUGE-L score of 9.9 to 15.7 for Multi-XScience, while the score ranges for Multi-news is 3.7 to 13.6 for ROUGE-2 and 10.4 to

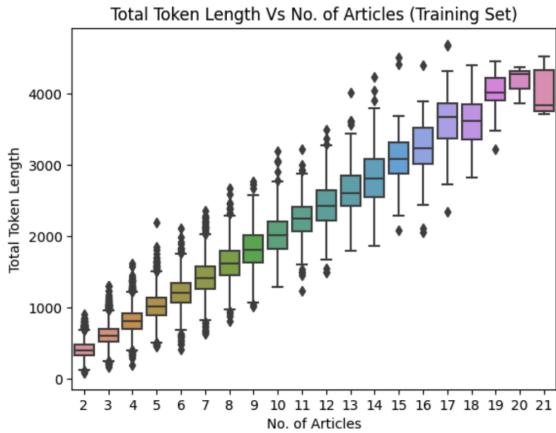


Figure 1: Total token length as a function of the number of articles in the training data set. The relationship between token length and the number of articles is linear.

As can be seen from the Table 1, the input length of the Multi X-Science dataset shows great variability, with some inputs so short that they can fit in more traditional summarization models, while others exceed even 4096, which is the max input length for Centrum. The same variability can also be observed in terms of the number of articles per sample (see Figure 1), with the two distributions showing a close-to-linear pattern. This provides an opportunity for us to use the Multi X-Science dataset to explore possible differences in model performance between short and long samples.

3.2 Models

Due to resource constraints, we focused on two pre-trained LLMs, namely LED and Centrum for our experiments. Of the three recent and well-performing models we reviewed, we did not pick PRIMERA for two reasons - (a) Centrum is, based on literature review, the best-performing model while LED is a common starting point for both PRIMERA and Centrum, meaning it could be treated as a more advanced baseline; and (b) the authors of PRIMERA already provided performance results of PRIMERA on the Multi-XScience test set.

20.8 for ROUGE-L. A similar gap remains even after fine-tuning of the models.

We adopted two different baselines for performance comparison. The first one, dubbed the “baseline model” is simply a lead-based model, whereby the first 3 sentences are chosen as the summary and is a common primary strategy implemented in various papers (e.g. [5], [7] and [9]). This also provides a useful baseline for comparing the degree of “copying” (see Section 3.3) from the main article which is an undesirable trait for the chosen dataset. The second baseline is the publicly available checkpoint for LED, which serves as the base architecture for recent MDS LLMs.

On top of the 2 baseline models, we ran experiments on the performance of the following models, with details on their fine-tuning and inference settings provided in Appendix B:

1. an alternative checkpoint of LED fine-tuned on the arXiv dataset
2. publicly available checkpoint of Centrum;
3. a version of the baseline LED⁸ that we fine-tuned on the Multi-XScience training set;
4. a similarly fine-tuned version of Centrum; and
5. a two-step model stacking the fine-tuned LED and Centrum.

Of these experiments, models (c) and (d) stem from our plan to test out the performance of fine-tuning strategies for MDS. During the process, we noticed that despite the larger input length of Centrum (4096 tokens), there are still cases where the inputs have to be truncated and information omitted. We therefore explored under model (e) a two-step model stacking process whereby the source articles are first condensed by our fine-tuned LED individually followed by an MDS process by (a) the fine-tuned LED or

⁸In theory, it would be better to fine-tune the alternative LED checkpoint which (a) is a larger version of LED with double the parameters; and (b) has been fine-tuned on text from the scientific fields. However, the GPU we had was unable to train the larger model while maintaining the max input token length at a reasonable level.

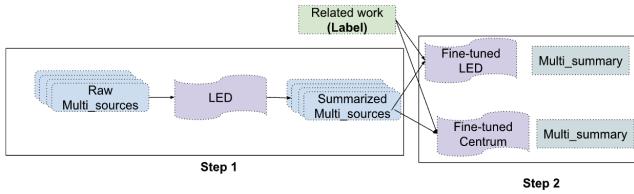


Figure 2: Two-steps model design

(b) the fine-tuned Centrum. The idea is to investigate if Centrum performs better when presented with more concise sources and without token truncation, see Figure 2. Refer to Appendix B for more details on the two step model.

3.3 Evaluation metrics

As in the standard practice for summarization, the primary metric we use is the ROUGE score. Among the different versions of ROUGE, we chose ROUGE-2 and ROUGE-L, with a particular focus on the latter given the prevalence of compound nouns in academic writings⁹, for gauging model performance. Compared to ROUGE-2, ROUGE-L is better suited for this task. In calculating these metrics, precision, recall, and f-measure are all considered¹⁰, and scores were computed for the overall test dataset, as well as the subsets of short, medium, and long token length samples¹¹.

Unlike other MDS datasets where the key information is expected to be found in most, if not

⁹This is an empirical observation from viewing samples in the dataset, and is quite intuitive considering academic terms are often ones glued together from multiple words (e.g. the theme of “Synchronous Optical Network Ring Assignment Problem” in sample 4820 of the test set, and “fixed-size ordinally forgetting encoding” in sample 3157).

¹⁰We also look only at the mean score but not the lower and upper ranges due to the limited scope of our study, though we recognize that looking at the variation in performance across samples is in itself a valid future investigation direction.

¹¹We divided the 5093 test samples into three groups. “Short” samples consisted of 1273 samples with less than 486 tokens (lower quartile) for the inputs; “medium” with 2546 samples with input lengths between 486 and 735 (upper quartile) tokens; and the remaining 1274 samples with input lengths greater than 1150 tokens are grouped into the “long” samples.

all, of the source articles, the Multi X-Science dataset has a different structure whereby the main article (i.e. first source) is expected to provide general context while the other sources will each provide a related, yet different angle of the topic. It is desirable for models to extract information not just from the main article, and we therefore measured the degree of “copying” from the main article by repurposing standard ROUGE calculations¹².

Finally, we supplemented our evaluation with a qualitative analysis whereby 25 random samples were selected and we reviewed the summaries written by the baseline LED and the 5 models experimented, judging them based on fluency, ability to extract information from multiple sources, and contrast information from the main article and other sources, as well as factual accuracy.

4 Results and Discussion

4.1 Overall performance

In the following table, we report the ROUGE-2 (R-2) and ROUGE-L (R-L) scores of the 2 baselines plus 5 models we tested, with precision (P), recall (R), and f-measure (F) scores presented. Scores for overall, short, medium, and long samples are also presented.

As shown in the Table 2 and Table 3, the fine-tuned models (LED and Centrum) provided a markedly improved performance over the baselines and off-the-shelf models. In fact, even though the models were only fine-tuned for up to 2 epochs¹³, our figures quickly approached the SOTA figures reported in [8] (i.e. ROUGE-2 of 6.8; ROUGE-L of 18.2) demonstrating the viability of LED and Centrum to be adopted for previously unseen MDS data. When considering the breakdowns, the improvement for longer samples was much more visible, indicating the MDS-specific model architectures are perhaps more

¹²Implemented using the standard rouge function by replacing the target reference that the function takes in with the main article of the sample.

¹³Training details and considerations are at Appendix B.

suitable for tasks with longer input lengths.

This result is consistent with what our qualitative analysis (Appendix C). For example, in sample 485, the Centrum model was able to extract the key ideas of 3 relevant studies and contrast them with the main article, forming a coherent summary, while the baseline and off-the-shelf models were overwhelmed by the large number of articles (9 in total) and choose only to copy one of the articles in full. This is very different from shorter samples, e.g. sample 4160 with only 3 articles, where the off-the-shelf models were still able to prepare a summary covering elements, albeit in a more extractive manner without efforts to change the tone or merge the extracted parts together.

For the two-step model, the ROUGE scores do not show any advantage of the approach over simple fine-tuned models. However, as we observe from the qualitative analysis (e.g. samples 831 and 5068), there is a strong tendency for the two-step model to highlight differences and similarities between the main and other articles. However, the downside is that the two-step model made quite a number of factual fallacies in the process¹⁴.

4.2 Degree of Copying

To assess the degree of copying from the main article, we calculated a re-purposed ROUGE-L score. As shown in Table 4, the fine-tuned and two-step models exhibited a significantly lower degree of “copying”, which is supported by the qualitative analysis showing that these models were able to extract information from multiple articles. In contrast, the baseline LED model showed the highest degree of “copying”, often copying substantial portions of the main article to the point that the number of copied sentences exceeded even that of the lead-based baseline.

Although the ROUGE scores did not reveal

¹⁴One possible reason is that we need better tuning on the amount and manner of information fed into the 2nd step model, but resource limitations were a main obstacle considering the generation of the first step answers alone took over 14 hours.

Model	R-2	Overall	Short	Medium	Long
Baseline (Lead)	P	4.50	4.78	4.68	3.87
	R	8.13	6.12	7.87	10.7
	F	5.40	<i>5.02</i>	5.56	5.48
Baseline (LED)	P	4.52	4.12	4.75	4.47
	R	6.99	8.17	7.07	5.63
	F	<i>5.21</i>	5.2	5.40	<i>4.83</i>
LED (Alt Ver)	P	5.25	5.32	5.30	5.09
	R	6.31	6.97	6.34	5.57
	F	5.39	5.58	5.44	5.11
Centrum	P	4.14	4.02	4.23	4.09
	R	8.12	9.48	8.14	6.70
	F	5.24	5.36	<i>5.34</i>	4.93
Our LED	P	8.75	8.03	8.67	9.63
	R	5.89	5.21	5.69	6.94
	F	6.56	5.83	6.41	7.62
Our Centrum	P	7.62	7.10	7.58	8.28
	R	7.27	6.38	7.13	8.44
	F	6.88	6.14	6.78	7.82
2-Step (Centrum)	P	5.87	5.66	5.92	5.99
	R	6.67	6.15	6.62	7.26
	F	5.86	5.42	5.87	6.28

Table 2: Rouge-2 Scores for Models Tested (Best in **bold**; worst is *italicized*)

Model	R-2	Overall	Short	Medium	Long
Baseline (Lead)	P	12.4	14.9	12.7	9.17
	R	21.7	18.8	21.3	25.2
	F	14.6	15.5	15	<i>12.9</i>
Baseline (LED)	P	12.9	11.6	13.3	13.5
	R	20.5	23.8	20.4	17.3
	F	14.9	14.8	15.1	14.6
LED (Alt Ver)	P	15.1	15.4	15.1	14.3
	R	18.6	20.3	18.7	16.7
	F	15.7	16.2	15.7	15.0
Centrum	P	11.3	10.8	11.4	11.7
	R	22.7	25.5	22.8	19.7
	F	<i>4.4</i>	<i>14.4</i>	<i>14.6</i>	14.2
Our LED	P	23.6	24.2	23.5	23.2
	R	15.9	15.7	15.5	16.8
	F	17.7	17.6	17.4	18.4
Our Centrum	P	19.8	20.9	19.6	19.1
	R	18.9	18.4	18.7	19.7
	F	17.8	17.9	17.6	18.1
2-Step (Centrum)	P	17.0	18.0	17.0	15.9
	R	19.6	19.7	19.5	19.6
	F	17.0	17.4	17.0	16.7

Table 3: Rouge-L Scores for Models Tested (Best in **bold**; worst is *italicized*)

Model	R-2	Overall	Short	Medium	Long
Baseline (Lead)	P	39.7	56.6	39.7	22.7
	R	40.7	41.1	39.0	43.9
	F	36.8	44.9	36.8	28.5
Baseline (LED)	P	76.9	83.1	86.9	50.7
	R	75.5	92.8	81.5	46.3
	F	74.2	85.6	81.9	47.2
LED (Alt Ver)	P	45.3	58.8	43.9	34.7
	R	33.6	44.3	31.5	27.0
	F	37.4	48.9	35.6	29.5
Centrum	P	48.2	56.5	48.6	39.2
	R	57.5	73.7	55.7	45.1
	F	51.2	62.4	50.8	41.0
Our LED	P	33.6	35.5	34.7	29.5
	R	13.5	13.0	13.2	14.7
	F	18.3	18.0	18.2	18.7
Our Centrum	P	29.4	32.9	29.4	26.1
	R	17.3	16.4	16.9	19.0
	F	20.5	20.7	20.1	20.9
2-Step (Centrum)	P	24.7	27.9	25.0	21.0
	R	17.2	17.1	16.8	18.0
	F	19.1	20.1	18.9	18.6

Table 4: Rouge-L Scores for Measuring “Copying” (Best in **bold**; worst is *italicized*)

any significant “copying” for the off-the-shelf LED and Centrum models, it is noteworthy that these models occasionally copied from a single reference article rather than the main article (e.g. see the output of the base Centrum for sample 4371). This indicates that there is still room for improving the “copying” metric, such as by taking the minimum of the ROUGE-L scores compared to each input article (rather than just the main article) to discourage “copying” from a single source.

4.3 Other observations

The qualitative analysis also allowed us to make the following observations on the models, including:

1. The fine-tuned and two-step models showed clear signs of learning common writing styles (using phrases such as “there is a large body of work...”) which is desirable in real world applications;
2. For samples containing irrelevant articles,

the fine-tuned models were able to safely navigate away from them (e.g. sample 845)

3. Despite having been fine-tuned on journal articles, the models showed greatest difficulties when presented with highly technical papers that use common words in uncommon meanings (e.g. in mathematics in sample 4858). This shows the potential to enhance performance by adding another fine-tuning step on domain-specific texts. Alternatively, the two-step model has continued to perform well in such cases which shows the potential of such an approach.

5 Conclusion

In conclusion, we demonstrated the viability of building custom-based MDS models through fine-tuning which is able to quickly approach SOTA scores even with simple and limited tuning. Through breakdown and qualitative analyses of the results, we highlighted the importance of considering the token length of input texts when selecting a text summarization model. Our results suggest that the fine-tuned LED and Centrum models are much more adept at MDS for longer texts, while the two-step model shows potential especially for highly technical texts despite needing more tuning. Future work could focus on further optimizing the tuning techniques (e.g. fine-tuning with domain-specific texts first before the MDS data), improving on the two-step approach, as well as looking at the performance stability of the models across data samples.

Appendix

A EDA and Data Preprocessing

A.1 Introduction

This appendix provides further findings from the EDA on the Multi-XScience dataset, as well as the data preprocessing procedures undertaken in the experiments. In particular, the distribution in the length of inputs and labels respectively guided us in determining (a) the maximum number of tokens fed into the models; and (b) the maximum output of length of generated summaries.

A.2 Distribution of Input and Output Lengths

To determine the input length, we looked at the average length of the samples (i.e. each with multiple articles) in the Multi-XScience dataset. Based on the parameters set out in Table 1 of the paper, and taking into account model limitations¹⁵ and GPU capabilities, we set the input numbers of tokens as 4096 in the fine-tuning process. As for the baseline and off-the-shelf models which do not require fine-tuning, we experimented on other input token lengths including 512, 1024 and 16384.

To determine the max output length¹⁶, we looked at the distribution of the token length of the standard summaries (label) that came with the dataset. A sample distribution for the training dataset, which closely resembles the distributions for the validation and test datasets, is provided below.

Considering the desired level of detail in the summary and the distribution of standard sum-

¹⁵LED can take in at most 16384 tokens and Centrum can take in 4096 tokens.

¹⁶LED and Centrum both do not come with options to set minimum output lengths. One can force such a parameter into the generation function, but initial experimentations showed minimal impact on the output.

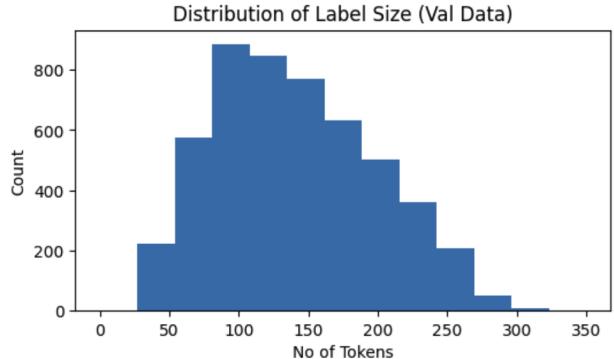


Figure 3: Distribution of the number of tokens in the label of training data. The number of tokens is right-skewed and centered around 100.

maries (label) tokens, we decided to set 256 as the maximum number of tokens¹⁷ at inference time.

A.3 Data Preprocessing

Multi-XScience includes a main abstract, which is from the main article; a standard summary of the related works with a separation of “@cite”, which is the summary of multi-source documentations, label datasets; and reference abstracts, which are the multiple abstracts from a different source of works. We processed the dataset by defining the separator of documents as “|||||”, a document separator, then concatenating reference abstracts using the document separator to replace citation references.

B Model

B.1 Baseline Model

Copying the first 3 sentences from each reference abstract was used as our first baseline. It is a simple and easy-to-implement approach that provides a starting point for summarizing multiple documents. It assumes that the first few sentences of a document contain important information that should be included in the summary.

¹⁷During a number of early stage experiments with LED, the maximum was set at 200 instead.

Main abstract	Related_work (Label)	Multi_sources	Data Preprocessing
Abstract: We give a purely topological definition of the perturbative quantum invariants of links and 3-manifolds associated with Chern-Simons field theory. Our definition is as close as possible to one given by Kontsevich. We will also establish some basic properties of these invariants, in particular that they are universally finite type with respect to algebraically split surgery and with respect to Torelli surgery. Torelli surgery is a mutual generalization of blink surgery of Garoufalidis and Levine and clasper surgery of Habiro.'	'related_work': 'Two other generalizations that can be considered are invariants of graphs in 3-manifolds, and invariants associated to other flat connections @cite_16 . We will analyze these in future work. Among other things, there should be a general relation between flat bundles and links in 3-manifolds on the one hand and finite covers and branched covers on the other hand @cite_26 .'	Source 1: 'This note is a sequel to our earlier paper of the same title [4] and describes invariants of rational homology 3-spheres associated to acyclic orthogonal local systems. Our work is in the spirit of the Axelrod–Singer papers [1], generalizes some of their results, and furnishes a new setting for the purely topological implications of their work.' Source 2: 'Recently, Mullins calculated the Casson-Walker invariant of the 2-fold cyclic branched cover of an oriented link in S^3 in terms of its Jones polynomial and its signature, under the assumption that the 2-fold branched cover is a rational homology 3-sphere. Using elementary principles, we provide a similar calculation for the general case. In addition, we calculate the LMO invariant of the p-fold branched cover of twisted knots in S^3 in terms of the Kontsevich integral of the knot.'	related_work': Two other generalizations that can be considered are invariants of graphs in 3-manifolds, and invariants associated to other flat connections We will analyze these in future work. Among other things, there should be a general relation between flat bundles and links in 3-manifolds on the one hand and finite covers and branched covers on the other hand

Figure 4: Multi-XScience data sample

This approach can serve as a useful baseline for several reasons:

1. It is easy to implement and does not require a lot of computational resources. This makes it a good starting point for developing more advanced summarization models.
2. It is straightforward to understand and explain. This makes it a good choice for initial experiments and evaluations.
3. It can provide a quick estimate of the performance of a summarization system. By comparing the summaries generated by this approach to human-written summaries or other machine-generated summaries, it is possible to get a rough idea of how well the system is performing.

B.2 Baseline Base LED (1K and 16K) Model

The base LED model has (“allenai/led-base-16384”), with max input token length of 16384 & ~200M parameters and consists of 12 transformer encoder layers and 2 transformer decoder

layers. This model is designed to be computationally efficient and suitable for low-resource environments. The baseline LED was generated based on a pre-trained LED model by Beltagy, L (2020) [1] on Multi-XScience test datasets with input tokens varying from 1k and 6k.

B.3 Off-the-shelf LED (1K and 16K) Model

The off-the-shelf LED is referred to the “large LED” during our experiments to indicate that it has more parameters and is larger than the “baseline” LED model. The model has (“allenai/led-large-16384-arxiv”), with max input token length of 16384 & ~512M parameters, and consists of 24 transformer encoder layers and 4 transformer decoder layers. This model has been pre-trained for general summarization tasks. The baseline LED was generated based on a pre-trained large LED and evaluated on Multi-XScience test datasets.

B.4 Off-the-shelf Centrum Model

The “ratishsp/Centrum” model with max input token length of 4096 & ~192M parameters and

consists of 12 transformer encoder layers and 12 transformer decoder layers. It is smaller than some of the larger transformer models, such as T5 and GPT-3, but larger than some of the smaller transformer models, such as the BART base. The publicly available Centrum checkpoint was built upon the LED model architecture using the 4096 token version. This means Centrum is not able to take in 16384 tokens.

B.5 Fine-tuned LED Model

Fine-tuning the “LED-large-16384-arxiv” model to the Multi-XScience train (size 30369) and validation (size 5066) datasets specifically for the related work summarization. We used 4096 max number of tokens when feeding the training and validation sets during fine-tuning, and used 200 as the max output tokens with a batch size of 2 (due to limited GPU resources) at inference time. The model is then used as a cross-entropy loss function for 2 epochs with 40 checkpoints generated for every 5

B.6 Fine-tuned Centrum Model

Fine-tuning the “ratishsp/Centrum” model to the X-science train and validation datasets specifically for the related works summarization. We used 4096 max number of tokens when feeding the training and validation sets during fine-tuning, and used 256 as the max output tokens with a batch size of 1 (due to limited GPU resources) at inference time. The model is then used as a cross-entropy loss function for 2 epochs, and evaluated on test datasets (size 5093) with ROUGE metrics.

B.7 Two-step LED/Centrum Model

In the first step, the LED model is used to generate a summary of each related source document. By summarizing each source document, the model can condense the information and reduce redundancy, making it easier to process and

analyze. For the first step, none of the individual source articles in the Multi-XScience dataset contain more than 4096 tokens, so there is no issue of the first step model (i.e. the fine-tuned LED) receiving truncated inputs.

In the second step, generate the summary based on the fine-tuned LED model to refine the summaries generated in the first step. In this way, we were hoping that the second model can learn to generate even more accurate and relevant summaries. Two different models were tested in the second step, namely our fine-tuned LED and the fine-tuned Centrum. Only the results of the latter, which proved superior, were reported in the paper.

B.8 Two-step Model results

However, from our evaluation and observation of this approach, the performance of the two-step model does not exceed the fine-tuned LED/Centrum model. One potential issue is that the large amount of information contained in the source documents may make it difficult for the model to generate concise and informative summaries. Ultimately, the performance is heavily dependent on the performance of step 1.

Overall, our results demonstrate the importance of utilizing advanced techniques like fine-tuning to unlock the full potential of language models in real-world applications.

C Qualitative Analysis

We also investigated how the model performed when summarizing short, medium, and long documents. The study aimed to evaluate the quality of the generated summaries and identify any trends or patterns in the model’s performance based on the length of the source documents. A total of 25 randomly drawn samples (from the 5093 samples in the X-Science test dataset) are analyzed. These 25 samples include:

- 10 short samples, i.e. with total token lengths below the lower quartile

- 10 medium samples, i.e. with total token lengths between the lower and upper quartiles
- 5 long samples, i.e. with token lengths above the upper quartile

Our analysis of the review (refer to the Analysis Appendix for more detailed review) results showed that the fine-tuned LED and Centrum learned the structure for multi-documents summary for X-science dataset, with particularly good results for longer documents (ref. Sample 485) while for shorter samples the advantage over the baseline and off-the-shelf models appear to be less as those models are able to extract information from multiple sources to some degree, albeit with limited ability to merge them into a single coherent summary (ref. Sample 4160).

Additionally, we observed that LED and Centrum seem to have difficulties in digesting some of the more technical topics, where the use of common language for unusual meanings might have confused the models (refer to no 4858). The two-step model somehow helps for these samples, but the downside is that the two-step model hallucinates quite a bit. In general, LED and Centrum never really hallucinate. To future improve this issue, we would suggest that further research in this area should consider:

- Add domain-specific knowledge: This can include adding specialized math dictionaries or knowledge bases to your training data.
- Use special tokens for math symbols: This can help the model recognize and differentiate between mathematical expressions and regular language. For example, use a special token such as “[math]” to indicate the start of a math expression and “[/math]” to indicate the end.
- Use more training data: Consider adding more math-specific data to the training set to help the model learn to recognize and understand math concepts.

- Fine-tune your model: fine-tuning it on a smaller dataset of math-specific documents. This can help the model learn to recognize and summarize math-related information more effectively.

Overall, our study highlights the importance of considering the length of source documents, and the training datasets topics when evaluating the performance of multi-document summarization models.

Future research in this area should focus on developing models that can effectively summarize long documents; and provide the model with more specialized knowledge to help the model recognize and differentiate between regular language and math expressions, improving its ability to summarize math-related documents effectively. Refer to the manual review documentation in more detail link.

D Detailed Review Samples

Refer to the following pages.

References

- [1] Iz Beltagy, Matthew E. Peters, and Arman Cohan. Longformer: The long-document transformer, 2020.
- [2] Alexander Fabbri, Irene Li, Tianwei She, Suyi Li, and Dragomir Radev. Multi-news: A large-scale multi-document summarization dataset and abstractive hierarchical model. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1074–1084, Florence, Italy, July 2019. Association for Computational Linguistics.
- [3] Demian Gholipour Ghalandari, Chris Hokamp, Nghia The Pham, John Glover, and Georgiana Ifrim. A large-scale multi-document summarization dataset from the wikipedia current events portal, 2020.

- [4] Peter J. Liu, Mohammad Saleh, Etienne Pot, Ben Goodrich, Ryan Sepassi, Lukasz Kaiser, and Noam Shazeer. Generating wikipedia by summarizing long sequences, 2018.
- [5] Yao Lu, Yue Dong, and Laurent Charlin. Multi-XScience: A large-scale dataset for extreme multi-document summarization of scientific articles. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 8068–8074, Online, November 2020. Association for Computational Linguistics.
- [6] Ratish Puduppully and Mark Steedman. Multi-document summarization with centroid-based pretraining, 2022.
- [7] Abigail See, Peter J. Liu, and Christopher D. Manning. Get to the point: Summarization with pointer-generator networks, 2017.
- [8] Wen Xiao, Iz Beltagy, Giuseppe Carenini, and Arman Cohan. Primera: Pyramid-based masked sentence pre-training for multi-document summarization, 2022.
- [9] Jingqing Zhang, Yao Zhao, Mohammad Saleh, and Peter J. Liu. Pegasus: Pre-training with extracted gap-sentences for abstractive summarization, 2020.

Qualitative Analysis of Generated Summaries

Introduction

1. This document provides a qualitative analysis of the summaries written by 6 models tested in the W266 project, namely:
 - a) Baseline LED model ("allenai/led-base-16384"), with max input token length of 16384 & ~200M parameters
 - b) Off-the-shelf LED model ("allenai/led-large-16384-arxiv")¹, with max input token length of 16384 & ~512M parameters
 - c) Off-the-shelf LED Centrum model ("ratishtsp/Centrum"), with max input token length of 4096² & ~192M parameters
 - d) Finetuned LED model (our own model), with max input token length of 16384 but finetuned using inputs tokenized at 4096 max
 - e) Finetuned Centrum model (our own model), with max input token length of 4096
 - f) Two-step LED-to-Centrum model (our own model), with max input token length of 4096 at the second step³
2. A total of 25 randomly drawn samples⁴ (from the 5093 samples in the X-Science test dataset) are analyzed. These 25 samples include:
 - a) 10 short samples, i.e. with total token lengths below the lower quartile
 - b) 10 medium samples, i.e. with total token lengths between the lower and upper quartiles
 - c) 5 long samples, i.e. with token lengths above the upper quartile
3. In analyzing the results, a few dimensions are looked at:
 - a) general fluency, i.e. whether the summary looks coherent and contains blatant errors or strange symbols
 - b) indication of multi-document summary ("MDS"), e.g. multiple sentences each attempting to summarize the work done in one or more of the journals being referenced
 - c) validity of MDS, e.g. whether sentences appearing to be summarizing different journals are in fact extracting key information from that journal's abstract
 - d) in addition, highest scores are highlighted in green, while lowest ones are in red, while key ideas from different articles are highlighted in different colors

High-level summary

4. As set out in the main paper, this analysis provides us with further insights, and we note the following:
 - a) in general, the fine-tuned LED, fine-tuned Centrum and two-step model manages to summarize from multiple sources most of the time (see the different colors for excerpts for different parts of the source)
 - b) the writings from these three models are usually very fluent as well, though occasional strange formulations and factual inaccuracies are observed
 - c) furthermore, these models show clear signs of learning the desired writing style, especially in contrasting the main article with the other references
 - d) as a contrast, the baseline LED is prone to just copying from just the first part of the main article; while the off-the-shelf LED and Centrum either just copies or sometimes extracts from multiple sources

Results [only those referred to in the main paper included; full version of this document with all 25 samples available at Github.]

(a) Short samples

- 3157, 4820: Samples with used to illustrate the presence of long and compound nouns
- 4160: Short example where the off-the-shelf LED and off-the-shelf Centrum managed to extract information from multiple sources as well (i.e. in addition to the tuned and two-step models)

No.	Abstracts	Label	Base LED	Off-the-shelf LED	Off-the-shelf Centrum	Tuned LED	Tuned Centrum	Two-step
3157	<p>(1):</p> <p>In this paper, we study a novel approach for named entity recognition (NER) and mention detection in natural language processing. Instead of treating NER as a sequence labelling problem, we propose a new local detection approach, which rely on the recent fixed-size ordinally forgetting encoding (FOFE) method to fully encode each sentence fragment and its left right contexts into a fixed-size representation. Afterwards, a simple feedforward neural network is used to reject or predict entity label for each individual fragment. The proposed method has been evaluated in several popular NER and mention detection tasks, including the CoNLL 2003 NER task and TAC-KBP2015 and TAC-KBP2016 Tri-lingual Entity Discovery and Linking (EDL) tasks. Our methods have yielded pretty strong performance in all of these examined tasks. This local detection approach has shown many advantages over the traditional sequence labelling methods.</p> <p>(2):</p> <p>In this paper, we study a novel approach for named entity recognition (NER) and mention detection in natural language processing. Instead of treating NER as a sequence labelling problem, we propose a new local detection approach, which rely on the recent fixed-size ordinally forgetting encoding (FOFE) method to fully encode each sentence fragment and its left right contexts into a fixed-size representation. Afterwards, a simple feedforward neural network is used to reject or predict entity label for each individual fragment. The proposed method has been evaluated in several popular NER and mention detection tasks, including the CoNLL 2003 NER task and TAC-KBP2015 and TAC-KBP2016 Tri-lingual Entity Discovery and Linking (EDL) tasks. Our methods have yielded pretty strong performance in all of these examined tasks. This local detection approach has shown many advantages over the traditional sequence labelling methods.</p>	The success of word embedding @cite encourages researchers to focus on machine-learned representation instead of heavy feature engineering in NLP. Using word embedding as the typical feature representation for words, NNs become competitive to traditional approaches in NER. Many NLP tasks, such as NER, chunking and part-of-speech (POS) tagging can be formulated as sequence labeling tasks. In @cite, deep convolutional neural networks (CNN) and conditional random fields (CRF) are used to model the dependencies between words and their surrounding context. The proposed model achieves state-of-the-art performance on various NLP benchmarks.	<p>Rouge scores:</p> <ul style="list-style-type: none">- Rouge 2: 0.0316 (prec) 0.0463 (recall) 0.0376 (f-1)- Rouge L: 0.1132 (prec) 0.1651 (recall) 0.1343 (f-1) <p>Summary:</p> <p>In this paper, we study a novel approach for named entity recognition (NER) and mention detection in natural language processing. Instead of treating NER as a sequence labelling problem, we propose a new local detection approach, which rely on the recent fixed-size ordinally forgetting encoding (FOFE) method to fully encode each sentence fragment and its left right contexts into a fixed-size representation. Afterwards, a simple feedforward neural network is used to reject or predict entity label for each individual fragment. The proposed method has been evaluated in several popular NER and mention detection tasks, including the CoNLL 2003 NER task and TAC-KBP2015 and TAC-KBP2016 Tri-lingual Entity Discovery and Linking (EDL) tasks. Our methods have yielded pretty strong performance in all of these examined tasks. This local detection approach has shown many advantages over the traditional sequence labelling methods.</p>	<p>Rouge scores:</p> <ul style="list-style-type: none">- Rouge 2: 0.0303 (prec) 0.037 (recall) 0.0333 (f-1)- Rouge L: 0.0977 (prec) 0.1193 (recall) 0.1074 (f-1) <p>Summary:</p> <p>In this paper, we study a novel approach for named entity recognition (NER) and mention detection in natural language processing. Instead of treating NER as a sequence labelling problem, we propose a new local detection approach, which rely on the recent fixed-size ordinally forgetting encoding (FOFE) method to fully encode each sentence fragment and its left right contexts into a fixed-size representation. Afterwards, a simple feedforward neural network is used to reject or predict entity label for each individual fragment. The proposed method has been evaluated in several popular NER and mention detection tasks, including the CoNLL 2003 NER task and TAC-KBP2015 and TAC-KBP2016 Tri-lingual Entity Discovery and Linking (EDL) tasks. Our methods have yielded pretty strong performance in all of these examined tasks. This local detection approach has shown many advantages over the traditional sequence labelling methods.</p>	<p>Rouge scores:</p> <ul style="list-style-type: none">- Rouge 2: 0.0685 (prec) 0.0463 (recall) 0.0552 (f-1)- Rouge L: 0.1045 (prec) 0.1193 (recall) 0.1355 (f-1) <p>Summary:</p> <p>In this paper, we study a novel approach for named entity recognition (NER) and mention detection in natural language processing. Instead of treating NER as a sequence labelling problem, we propose a new local detection approach, which rely on the recent fixed-size ordinally forgetting encoding (FOFE) method to fully encode each sentence fragment and its left right contexts into a fixed-size representation. Afterwards, a simple feedforward neural network is used to reject or predict entity label for each individual fragment. The proposed method has been evaluated in several popular NER and mention detection tasks, including the CoNLL 2003 NER task and TAC-KBP2015 and TAC-KBP2016 Tri-lingual Entity Discovery and Linking (EDL) tasks. Our methods have yielded pretty strong performance in all of these examined tasks. This local detection approach has shown many advantages over the traditional sequence labelling methods.</p>	<p>Rouge scores:</p> <ul style="list-style-type: none">- Rouge 2: 0.0 (prec) 0.0 (recall) 0.0 (f-1)- Rouge L: 0.1957 (prec) 0.1835 (recall) 0.1887 (f-1) <p>Summary:</p> <p>The Skip-gram model @cite is an efficient method for learning high-quality distributed vector representations that can be applied to various natural language processing tasks including entity linking and question answering.</p>	<p>Rouge scores:</p> <ul style="list-style-type: none">- Rouge 2: 0.049 (prec) 0.0463 (recall) 0.0476 (f-1)- Rouge L: 0.1942 (prec) 0.1835 (recall) 0.1887 (f-1) <p>Summary:</p> <p>In @cite, the authors proposed a unified neural network architecture and learning algorithm that can be applied to various natural language processing tasks including entity linking and question answering.</p>	

¹ This model has been pre-trained by the authors for general summarization task using arXiv dataset.

² The publicly available Centrum checkpoint was built upon the LED model architecture using the 4096 token version. This means Centrum is not able to take in 16384 tokens.

³ For the first step, none of the individual source articles in the X-science dataset contain more than 4096 tokens, so there is no issue of the first step model (i.e. the finetuned LED) receiving truncated inputs.

⁴ For details, please refer to the "Qualitative_analysis.ipynb" notebook. We originally planned for 20 samples for each category but ended up doing 25 in total due to time limitations. The 10 are the first and last 5 of the randomly selected samples for the short and medium samples, and just the first 5 for the long samples.

No.	Abstracts	Label	Base LED	Off-the-shelf LED	Off-the-shelf Centrum	Tuned LED	Tuned Centrum	Two-step
	<p>We propose a unified neural network architecture and learning algorithm that can be applied to various natural language processing tasks including: part-of-speech tagging, chunking, named entity recognition, and semantic role labeling. This versatility is achieved by trying to avoid task-specific engineering and therefore disregarding a lot of prior knowledge. Instead of exploiting man-made input features carefully optimized for each task, our system learns internal representations on the basis of vast amounts of mostly unlabeled training data. This work is then used as a basis for building a freely available tagging system with good performance and minimal computational requirements.</p> <p>(3):</p> <p>The recently introduced continuous Skip-gram model is an efficient method for learning high-quality distributed vector representations that capture a large number of precise syntactic and semantic word relationships. In this paper we present several extensions that improve both the quality of the vectors and the training speed. By subsampling of the frequent words we obtain significant speedup and also learn more regular word representations. We also describe a simple alternative to the hierarchical softmax called negative sampling. An inherent limitation of word representations is their indifference to word order and their inability to represent idiomatic phrases. For example, the meanings of "Canada" and "Air" cannot be easily combined to obtain "Air Canada". Motivated by this example, we present a simple method for finding phrases in text, and show that learning good vector representations for millions of phrases is possible.</p>	<p>F) are used to infer NER labels at a sentence level, where they still use many hand-crafted features to improve performance, such as capitalization features explicitly defined based on first-letter capital, non-initial capital and so on.</p>	<p>a sequence labelling problem, we propose a new local detection approach, which rely on the recent fixed-size ordinally forgetting encoding (FOFE) method to fully encode each sentence fragment and its left right contexts into a fixed-size representation. Afterwards, a simple feedforward neural network is used to reject or predict entity label for each individual fragment. The proposed method has been evaluated in several popular NER and mention detection tasks, including the CoNLL 2003 NER task and TAC-KBP2016 Tri-lingual Entity Discovery and Linking (EDL) tasks. Our methods have yielded pretty strong performance in all of these examined tasks. This local detection approach has shown many advantages over the traditional sequence labelling methods. We propose a unified neural network architecture and learning algorithm that can be applied to various natural language processing tasks including: part-of-speech tagging, chunking, named entity recognition, and semantic role labeling. This versatility is achieved by trying to avoid task-specific engineering and therefore disregarding a lot of prior knowledge. Instead of exploiting man-made input features carefully optimized for each task, our system learns internal representations on the basis of vast amounts of mostly unlabeled training data. This work is then used as a basis for building a freely available tagging system with good performance and minimal computational requirements.</p> <p>Comment:</p> <ul style="list-style-type: none"> - Copied the first tokens only - No indication of MDS 	<p>recognition as a sequence labelling problem, we propose a new local detection approach, which rely on the recent fixed-size ordinally forgetting encoding (FOFE) method to fully encode each sentence fragment and its left right contexts into a fixed-size representation. Afterwards, a simple feedforward neural network is used to reject or predict entity label for each individual fragment. The proposed method has been evaluated in several popular NER and mention detection tasks, including the CoNLL 2003 NER task and TAC-KBP2016 Tri-lingual Entity Discovery and Linking (EDL) tasks. Our methods have yielded pretty strong performance in all of these examined tasks. This local detection approach has shown many advantages over the traditional sequence labelling methods. We propose a unified neural network architecture and learning algorithm that can be applied to various natural language processing tasks including: part-of-speech tagging, chunking, named entity recognition, and semantic role labeling. This versatility is achieved by trying to avoid task-specific engineering and therefore disregarding a lot of prior knowledge. Instead of exploiting man-made input features carefully optimized for each task, our system learns internal representations on the basis of vast amounts of mostly unlabeled training data. This work is then used as a basis for building a freely available tagging system with good performance and minimal computational requirements.</p> <p>Comment:</p> <ul style="list-style-type: none"> - Copied the first tokens only - No indication of MDS 	<p>sequence labelling problem, we propose a new local detection approach, which rely on the recent fixed-size ordinally forgetting encoding (FOFE) method to fully encode each sentence fragment and its left right contexts into a fixed-size representation. Afterwards, a simple feedforward neural network is used to reject or predict entity label for each individual fragment. The proposed method has been evaluated in several popular NER and mention detection tasks, including the CoNLL 2003 NER task and TAC-KBP2016 Tri-lingual Entity Discovery and Linking (EDL) tasks. Our methods have yielded pretty strong performance in all of these examined tasks. This local detection approach has shown many advantages over the traditional sequence labelling methods. We propose a unified neural network architecture and learning algorithm that can be applied to various natural language processing tasks including: part-of-speech tagging, chunking, named entity recognition, and semantic role labeling. This versatility is achieved by trying to avoid task-specific engineering and therefore disregarding a lot of prior knowledge. Instead of exploiting man-made input features carefully optimized for each task, our system learns internal representations on the basis of vast amounts of mostly unlabeled training data. This work is then used as a basis for building a freely available tagging system with good performance and minimal computational requirements.</p> <p>Comment:</p> <ul style="list-style-type: none"> - Shows signs of MDS, which contrasts the 2nd (green) article with the 1st (yellow) one. - However, the 3rd article is not covered at all 	<p>uding: part-of-speech tagging, chunking, named entity recognition, and semantic role labeling. In this work, we propose a new local detection approach, which rely on the recent fixed-size ordinally forgetting encoding (FOFE) method to fully encode each sentence fragment and its left right contexts into a fixed-size representation. Afterwards, a simple feedforward neural network is used to reject or predict entity label for each individual fragment. The proposed method has been evaluated in several popular NER and mention detection tasks, including the CoNLL 2003 NER task and TAC-KBP2016 Tri-lingual Entity Discovery and Linking (EDL) tasks. Our methods have yielded pretty strong performance in all of these examined tasks. This local detection approach has shown many advantages over the traditional sequence labelling methods. We propose a unified neural network architecture and learning algorithm that can be applied to various natural language processing tasks including: part-of-speech tagging, chunking, named entity recognition, and semantic role labeling. This versatility is achieved by trying to avoid task-specific engineering and therefore disregarding a lot of prior knowledge. Instead of exploiting man-made input features carefully optimized for each task, our system learns internal representations on the basis of vast amounts of mostly unlabeled training data. This work is then used as a basis for building a freely available tagging system with good performance and minimal computational requirements.</p> <p>Comment:</p> <ul style="list-style-type: none"> - Shows signs of MDS, which contrasts the 2nd (green) article with the 1st (yellow) one. - However, the 3rd article is not covered at all 	<p>actic and semantic word relationships. By subsampling of the frequent words, it obtain significant speedup and also learn more regular word representations.</p> <p>Comment:</p> <ul style="list-style-type: none"> - Shows weak signs of MDS, providing a summary of the 3rd article only (blue) and presenting it as different from the main one. 	<p>cluding part-of-speech tagging, chunking, named entity recognition, and semantic role labeling. However, unlike our work, they do not use a deep neural network to train their system. Instead, they use a multi-task learning algorithm to train their network on a large amount of unlabeled training data. In contrast, we use a deep network to train our system on a much larger amount of training data, and we do not need to train our network on unlabeled data.</p> <p>Comment:</p> <ul style="list-style-type: none"> - Shows signs of MDS, providing a summary of the 2nd article (green) and contrasting it with the main one (red). - However, the red parts also contain hallucinations, not to mention the neglect of the 3rd article
4160	<p>(1):</p> <p>We present a method for extracting depth information from a rectified image pair. We train a convolutional neural network to predict how well two image patches match and use it to compute the stereo matching cost. The cost is refined by cross-based cost aggregation and semiglobal matching, followed by a left-right consistency check to eliminate errors in the occluded regions. Our stereo method achieves an error rate of 2.61 on the KITTI stereo dataset and is currently (August 2014) the top performing method on this dataset.</p> <p>(2):</p> <p>While machine learning has been instrumental to the ongoing progress in most areas of computer vision, it has not been applied to the problem of stereo matching with similar frequency or success. We present a supervised learning approach for predicting the correctness of stereo matches based on a random forest and a set of features that capture various forms of information about each pixel. We show highly competitive results in predicting the correctness of matches and in confidence estimation, which allows us to rank pixels according to the reliability of their assigned disparities. Moreover, we show how these confidence values can be used to improve the accuracy of disparity maps by integrating them with an MRF-based stereo algorithm. This is an important distinction from current literature that has mainly focused on sparsification by removing potentially erroneous disparities to generate quasi-dense disparity maps.</p> <p>(3):</p> <p>With the aim to improve accuracy of stereo confidence measures, we apply the random decision forest framework to a large set of diverse stereo confidence measures. Learning and testing sets were drawn from the recently introduced KITTI dataset, which currently poses higher challenges to stereo solvers than other benchmarks with ground truth for stereo evaluation. We experiment with semi global matching stereo (SGM) and a census data term, which is the best performing real-time capable stereo method known to date. On KITTI images, SGM still produces a significant amount of error. We obtain consistently improved area under curve values of sparsification measures in comparison to best performing single stereo confidence measures where numbers of stereo errors are large. More specifically, our method performs best in all but one out of 194 frames of the KITTI dataset.</p>	<p>Recent work @cite @cite focused on estimating the confidence of the computed matching cost. used a random forest classifier to combine several confidence measures. Similarly, trained a random forest classifier to predict the confidence of the matching cost and used the predictions as soft constraints in a Markov random field to decrease the error of the stereo method.</p>	<p>Rouge scores:</p> <ul style="list-style-type: none"> - Rouge 2: 0.0523 (prec) 0.1579 (recall) 0.0786 (f-1) - Rouge L: 0.0983 (prec) 0.2931 (recall) 0.1472 (f-1) <p>Summary:</p> <p>We present a method for extracting depth information from a rectified image pair. We train a convolutional neural network to predict how well two image patches match and use it to compute the stereo matching cost. The cost is refined by cross-based cost aggregation and semiglobal matching, followed by a left-right consistency check to eliminate errors in the occluded regions. Our stereo method achieves an error rate of 2.61 on the KITTI stereo dataset and is currently (August 2014) the top performing method on this dataset. While machine learning has been instrumental to the ongoing progress in most areas of computer vision, it has not been applied to the problem of stereo matching with similar frequency or success. We present a supervised learning approach for predicting the correctness of stereo matches based on a random forest and a set of features that capture various forms of information about each pixel. We show highly competitive results in predicting the correctness of matches and in confidence estimation, which allows us to rank pixels according to the reliability of their assigned disparities. Moreover, we show how these confidence values can be used to improve the accuracy of disparity maps by integrating them with an MRF-based stereo algorithm. This is an important distinction from current literature that has mainly focused on sparsification by removing potentially erroneous disparities to generate quasi-dense disparity maps.</p> <p>Comment:</p> <ul style="list-style-type: none"> - Clear sign of MDS, with the first green part showing an attempt to make comparison - The blue and 2nd green part shows possible mix up between the 2nd and 3rd models, caused probably by the repeated use of stereo and SGM - In a similar vein, the red part is from the 1st article (see numbers in yellow) but the sentence structure is unclear 	<p>Rouge scores:</p> <ul style="list-style-type: none"> - Rouge 2: 0.0826 (prec) 0.1579 (recall) 0.1084 (f-1) - Rouge L: 0.2 (prec) 0.3793 (recall) 0.2619 (f-1) <p>Summary:</p> <p>We present a supervised learning approach for predicting the correctness of stereo matches based on a random forest and a set of features that capture various forms of information about each pixel. We're using the random decision forest framework, which currently poses higher challenges to stereo solvers than other benchmarks with ground truth for stereo evaluation. We experiment with semi global matching stereo (SGM) and a census data term, which is the best performing real-time capable stereo method known to date. On KITTI images, SGM still produces a significant amount of error. We obtain consistently improved area under curve values of sparsification measures in comparison to best performing single stereo confidence measures where numbers of stereo errors are large. More specifically, our method performs best in all but one out of 194 frames of the KITTI dataset.</p> <p>Comment:</p> <ul style="list-style-type: none"> - Clear sign of MDS, with the first green part showing an attempt to make comparison - The blue and 2nd green part shows possible mix up between the 2nd and 3rd models, caused probably by the repeated use of stereo and SGM - In a similar vein, the red part is from the 1st article (see numbers in yellow) but the sentence structure is unclear 	<p>Rouge scores:</p> <ul style="list-style-type: none"> - Rouge 2: 0.0468 (prec) 0.1404 (recall) 0.0702 (f-1) - Rouge L: 0.1047 (prec) 0.3103 (recall) 0.1565 (f-1) <p>Summary:</p> <p>Our work is closely related to @cite @cite in that they use a random forest to predict the correctness of stereo matches based on a random forest and a set of features that capture various forms of information about each pixel. However, our work is different in that we train a convolutional neural network to predict how well two image patches match and use it to compute the stereo matching cost. the cost is refined by cross-based cost aggregation and semiglobal matching, followed by a left-right consistency check to eliminate errors in the occluded regions. Our stereo method achieves an error rate of 2.61 on the KITTI stereo dataset and is currently the top performing method on this dataset.</p> <p>Comment:</p> <ul style="list-style-type: none"> - Some signs of MDS, with the first green part showing an attempt to make comparison - The blue and 2nd green part shows possible mix up between the 2nd and 3rd models, caused probably by the repeated use of stereo and SGM - In a similar vein, the red part is from the 1st article (see numbers in yellow) but the sentence structure is unclear 	<p>Rouge scores:</p> <ul style="list-style-type: none"> - Rouge 2: 0.1515 (prec) 0.1754 (recall) 0.1626 (f-1) - Rouge L: 0.2687 (prec) 0.3103 (recall) 0.2393 (f-1) <p>Summary:</p> <p>Our work is also related to the semi-global matching stereo (SGM) method @cite @cite, which is based on a random forest and a set of features that capture various forms of information about each pixel. SGM achieves an error rate of 2.61 on the KITTI stereo dataset and is currently the top performing method on this dataset.</p> <p>Comment:</p> <ul style="list-style-type: none"> - Shows signs of MDS, with the first green part showing an attempt to make comparison - The blue and 2nd green part shows possible mix up between the 2nd and 3rd models, caused probably by the repeated use of stereo and SGM - In a similar vein, the red part is from the 1st article (see numbers in yellow) but the sentence structure is unclear 	<p>Rouge scores:</p> <ul style="list-style-type: none"> - Rouge 2: 0.0517 (prec) 0.0526 (recall) 0.0522 (f-1) - Rouge L: 0.2373 (prec) 0.2414 (recall) 0.2393 (f-1) <p>Summary:</p> <p>In @cite @cite, the authors used a random forest to predict the correctness of stereo matching based on a random forest and a set of features that capture various forms of information about each pixel. However, they do not consider sparsification. In addition, they did not use sparsification to improve the accuracy of the disparity maps, and they do not require sparsification in order to improve the performance of the matching. Moreover, their work does not rely on the sparsification of sparsification, but rather on sparsifying the disparity maps in order to achieve sparsification without sparsification and without sparsification.</p> <p>Comment:</p> <ul style="list-style-type: none"> - Clear sign of MDS, with the first green part showing an attempt to make comparison - The blue and 2nd green part shows possible mix up between the 2nd and 3rd models, caused probably by the repeated use of stereo and SGM - In a similar vein, the red part is from the 1st article (see numbers in yellow) but the sentence structure is unclear 	<p>Rouge scores:</p> <ul style="list-style-type: none"> - Rouge 2: 0.1122 (prec) 0.193 (recall) 0.1419 (f-1) - Rouge L: 0.2121 (prec) 0.3621 (recall) 0.2675 (f-1) <p>Summary:</p> <p>In @cite @cite, the authors used a random forest to predict the correctness of stereo matching based on a random forest and a set of features that capture various forms of information about each pixel. However, they do not consider sparsification. In addition, they did not use sparsification to improve the accuracy of the disparity maps, and they do not require sparsification in order to improve the performance of the matching. Moreover, their work does not rely on the sparsification of sparsification, but rather on sparsifying the disparity maps in order to achieve sparsification without sparsification and without sparsification.</p> <p>Comment:</p> <ul style="list-style-type: none"> - Clear sign of MDS, with the first green part showing an attempt to make comparison - The blue and 2nd green part shows possible mix up between the 2nd and 3rd models, caused probably by the repeated use of stereo and SGM - In a similar vein, the red part is from the 1st article (see numbers in yellow) but the sentence structure is unclear

No.	Abstracts	Label	Base LED	Off-the-shelf LED	Off-the-shelf Centrum	Tuned LED	Tuned Centrum	Two-step
			<p>dicting the correctness of stereo matches based on a random forest and a set of features that capture various forms of information about each pixel. We show highly competitive results in predicting the accuracy of matches and in confidence estimation, which allows us to rank</p> <p>Comment:</p> <ul style="list-style-type: none"> - Copied the first tokens only - No indication of MDS 	<p>- Copied from the 1st (yellow) and 2nd (green) articles</p> <p>- No indication of MDS</p>	<p>re errors are large. More specifically, our method performs best in all but one out of 194 frames of the KITTI dataset. This is an important distinction from current literature that has mainly focused on sparsification by removing potentially erroneous disparities to generate quasi-dense disparity maps.</p> <p>Comment:</p> <ul style="list-style-type: none"> - Some signs of MDS by contrasting the 2nd (green) and 3rd (blue articles) - However, the phrase in red is also copied from the 2nd article so it is uncertain if the model is actually doing MDS - It also mixed up the facts in the 2nd and 3rd articles 		<p>as to whether the model is referring to the main article or not.</p>	<p>contrasting it with the main one.</p> <ul style="list-style-type: none"> - However, the red parts are hallucinations
4820	<p>(1): This paper presents a new method and a constraint-based objective function to solve two problems related to the design of optical telecommunication networks, namely the Synchronous Optical Network Ring Assignment Problem (SRAP) and the Intra-ring Synchronous Optical Network Design Problem (IDP). These network topology problems can be represented as a graph partitioning with capacity constraints as shown in previous works. We present here a new objective function and a new local search algorithm to solve these problems. Experiments conducted in C</p> <p>(2): This paper considers two problems that arise in the design of optical telecommunication networks when a ring-based topology is adopted, namely the SONET Ring Assignment Problem and the Intra-ring Synchronous Optical Network Design Problem. We show that these two network topology problems correspond to graph partitioning problems with capacity constraints: the first is a vertex partitioning problem, while the latter is an edge partitioning problem. We consider solution methods for both problems, based on metaheuristic algorithms. We first describe variable objective functions that depend on the transition from one solution to a neighboring one, then we apply several diversification and intensification techniques including Path Relinking, exploring Tabu Search and Scatter Search. Finally we propose a diversification method based on the use of multiple neighborhoods. A set of extensive computational results is used to compare the behaviour of the proposed methods and objective functions.</p>	<p>The two other methods described in @cite are based on the idea that to save ADMs a good solution should have very dense rings. They are both greedy and rely on a clique algorithm. In graph theory, a clique in an undirected graph \mathcal{G} is a subset of the vertex set V, such that for every two vertices in V, there exists an edge connecting the two. Finding a clique is not that easy, a way to do it is to use an "Union-Find" strategy, two clique C_1 and C_2 such that each node in C_1 is adjacent to each node in C_2 then merge the two cliques $C_1 \cup C_2$. The associated heuristic starts by considering each node to be a clique of size one, and to merge two cliques into a larger clique until there are no more possible merges.</p>	<p>Rouge scores:</p> <ul style="list-style-type: none"> - Rouge 2: 0.0132 (prec) 0.0144 (recall) 0.0138 (f-1) - Rouge L: 0.1316 (prec) 0.1429 (recall) 0.137 (f-1) <p>Summary: This paper presents a new method and a constraint-based objective function to solve two problems related to the design of optical telecommunication networks, namely the Synchronous Optical Network Ring Assignment Problem (SRAP) and the Intra-ring Synchronous Optical Network Design Problem (IDP). These network topology problems can be represented as a graph partitioning with capacity constraints as shown in previous works. We present here a new objective function and a new local search algorithm to solve these problems. Experiments conducted in CERN and the University of California, Berkeley. This paper considers two problems that arise in the design of Optical telecommunication networks when a ring-based topology is adopted, namely the SONET Ring Assignment Problem and the Intra-ring Synchronous Optical Network Design Problem. We show that these two network topology problems correspond to graph partitioning problems with capacity constraints: the first is a vertex partitioning problem, while the latter is an edge partitioning problem.</p> <p>Comment:</p> <ul style="list-style-type: none"> - Copied the first tokens only - No indication of MDS 	<p>Rouge scores:</p> <ul style="list-style-type: none"> - Rouge 2: 0.0 (prec) 0.0 (recall) 0.0 (f-1) - Rouge L: 0.1579 (prec) 0.0857 (recall) 0.1111 (f-1) <p>Summary: this paper presents a new method and a constraint-based objective function to solve two problems related to the design of optical telecommunication networks, namely the Synchronous Optical Network Ring Assignment Problem (SRAP) and the Intra-ring Synchronous Optical Network Design Problem (IDP). These network topology problems can be represented as a graph partitioning with capacity constraints as shown in previous works. We present here a new objective function and a new local search algorithm to solve these problems.</p> <p>Comment:</p> <ul style="list-style-type: none"> - Copied the first tokens only - No indication of MDS 	<p>Rouge scores:</p> <ul style="list-style-type: none"> - Rouge 2: 0.0205 (prec) 0.0288 (recall) 0.024 (f-1) - Rouge L: 0.1276 (prec) 0.1786 (recall) 0.1488 (f-1) <p>Summary: This paper presents a new method and a constraint-based objective function to solve two problems related to the design of optical telecommunication networks, namely the Synchronous Optical Network Ring Assignment Problem (SRAP) and the Intra-ring Synchronous Optical Network Design Problem (IDP). These network topology problems can be represented as a graph partitioning with capacity constraints as shown in previous works. We present here a new objective function and a new local search algorithm to solve these problems.</p> <p>Comment:</p> <ul style="list-style-type: none"> - Copied the first tokens only - No indication of MDS 	<p>Rouge scores:</p> <ul style="list-style-type: none"> - Rouge 2: 0.0189 (prec) 0.0072 (recall) 0.0104 (f-1) - Rouge L: 0.2407 (prec) 0.0929 (recall) 0.134 (f-1) <p>Summary: In the context of optical telecommunication networks, the SONET Ring Assignment Problem (SRAP) and the Intra-ring Synchronous Optical Network Design Problem (IDP) are two problems that arise when a ring-based topology is adopted. The SRAP problem can be represented as a graph partitioning problem, while the IDP problem is an edge partitioning problem.</p> <p>Comment:</p> <ul style="list-style-type: none"> - Shows weak signs of MDS, with the model just extracting the common theme of the 2 articles 	<p>Rouge scores:</p> <ul style="list-style-type: none"> - Rouge 2: 0.0357 (prec) 0.0072 (recall) 0.012 (f-1) - Rouge L: 0.2414 (prec) 0.05 (recall) 0.0828 (f-1) <p>Summary: In @cite, an objective function and a diversification method are proposed to solve the SONET Ring Assignment Problem (SRAP) and the Intra-ring Synchronous Optical Network Design Problem (IDP).</p> <p>Comment:</p> <ul style="list-style-type: none"> - Shows weak signs of MDS, which presents the 2nd (green) article as different from the main one 	<p>Rouge scores:</p> <ul style="list-style-type: none"> - Rouge 2: 0.1233 (prec) 0.0647 (recall) 0.0849 (f-1) - Rouge L: 0.2703 (prec) 0.1429 (recall) 0.1869 (f-1) <p>Summary: In @cite, the objective functions of the SONET ring assignment problem and the IDP problem are discussed in terms of a metaheuristic algorithm. In particular, the objective function for SONET Ring Assignment Problem is described in @cite. The objective functions of SONET Ring Assignment Problem are defined as: where n is the number of neighbors, m the number of nodes in the network, k and l are the number of vertices in the network.</p> <p>Comment:</p> <ul style="list-style-type: none"> - Shows weak signs of MDS, which presents the 2nd (green) article as different from the main one - Contains more relevant information than the other tuned models

No.	Abstracts	Label	Base LED	Off-the-shelf LED	Off-the-shelf Centrum	Tuned LED	Tuned Centrum	Two-step
					- No indication of MDS			

(b) Medium samples

- 831, 5068: examples where the tuned and two-step models show the ability to contrast and compare the main article with other references.
- 845: sample where the tuned models avoided quoting from the irrelevant input
- 4371: sample where even the off-the-shelf Centrum resorted to just copying from one of the articles
- 4858: sample where highly technical language use is present causing the models to show weaker MDS performance

No.	Abstracts	Label	Base LED	Off-the-shelf LED	Off-the-shelf Centrum	Tuned LED	Tuned Centrum	Two-step
831	<p>(1): Learning to drive faithfully in highly stochastic urban settings remains an open problem. To that end, we propose a Multi-task Learning from Demonstration (MT-LfD) framework which uses supervised auxiliary task prediction to guide the main task of predicting the driving commands. Our framework involves an end-to-end trainable network for imitating the expert demonstrator's driving commands. The network intermediately predicts visual affordances and action primitives through direct supervision which provide the aforementioned auxiliary supervised guidance. We demonstrate that such joint learning and supervised guidance facilitates hierarchical task decomposition, assisting the agent to learn faster, achieve better driving performance and increases transparency of the otherwise black-box end-to-end network. We run our experiments to validate the MT-LfD framework in CARLA, an open-source urban driving simulator. We introduce multiple non-player agents in CARLA and induce temporal noise in them for realistic stochasticity.</p> <p>(2): An artificial agent is developed that learns to play a diverse range of classic Atari 2600 computer games directly from sensory experience, achieving a performance comparable to that of an expert human player; this work paves the way to building general-purpose learning algorithms that bridge the divide between perception and action.</p> <p>(3): Reinforcement learning is the learning of a mapping from situations to actions so as to maximize a scalar reward or reinforcement signal. The learner is not told which action to take, as in most forms of machine learning, but instead must discover which actions yield the highest reward by trying them. In the most interesting and challenging cases, actions may affect not only the immediate's reward, but also the next situation, and through that all subsequent rewards. These two characteristics—trial-and-error search and delayed reward—are the two most important distinguishing features of reinforcement learning.</p> <p>(4): Several deep learning approaches have been applied to the autonomous driving task, many employing end-to-end deep neural networks. Autonomous driving is complex, utilizing multiple behavioral modalities ranging from lane changing to turning and stopping. However, most existing approaches do not factor in the different behavioral modalities of the driving task into the training strategy. This paper describes a technique for using Multi-Modal Multi-Task Learning that considers multiple behavioral modalities as distinct modes of operation for an end-to-end autonomous deep neural network, utilizing the insertion of modal information as secondary input data. Using labeled data from hours of driving our fleet of 1 10th scale model cars, we trained multiple neural networks to imitate the steering angle and driving speed of human control of a car. We show that in each case, our models trained with MTL can match or outperform multiple networks trained on individual tasks, while using a fraction of the parameters and having more distinct modes of operation than a network trained without MTL on the same multi-modal data. These results should encourage Multi-Modal MTL-style training with the insertion of Modal Information for tasks with related behaviors.</p> <p>(5): In recent years different lines of evidence have led to the idea that motor actions and movements in both vertebrates and invertebrates are composed of elementary building blocks. The entire motor repertoire can be spanned by applying a well-defined set of operations and transformations to these primitives and by combining them in many different ways according to well-defined syntactic rules. Motor and movement primitives and modules might exist at the neural, dynamic and kinematic levels with complicated mapping among the elementary building blocks subserving these different levels of representation. Hence, while considerable progress has been made in recent years in unravelling the nature of these primitives, new experimental, computational and conceptual approaches are needed to further advance our understanding of motor compositionality.</p>	<p>Multi-task learning (MTL) research shows the joint training of auxiliary related side-tasks along with the main task enhances the training performance. MTL in neural networks has been successfully demonstrated in many tasks previously including text-to-speech conversion, natural language processing, speech processing and computer vision. In the field of sequential decision making, @cite demonstrate MTL for 3D game playing, @cite and @cite demonstrate MTL in 3D maze navigation task whereas @cite utilize the MTL framework for autonomous driving. Instead of employing future control outputs as auxiliary tasks as shown by @cite, in this work we employ action and visual abstractions to guide the driving behavior.</p>	<p>Rouge scores:</p> <ul style="list-style-type: none"> - Rouge 2: 0.0721 (prec) 0.0748 (recall) 0.0734 (f-1) - Rouge L: 0.1786 (prec) 0.1852 (recall) 0.1818 (f-1) <p>Summary: Learning to drive faithfully in highly stochastic urban settings remains an open problem. To that end, we propose a Multi-task Learning from Demonstration (MT-LfD) framework which uses supervised auxiliary task prediction to guide the main task of predicting the driving commands. Our framework involves an end-to-end trainable network for imitating the expert demonstrator's driving commands. The network intermediately predicts visual affordances and action primitives through direct supervision which provide the aforementioned auxiliary supervised guidance. We demonstrate that such joint learning and supervised guidance facilitates hierarchical task decomposition, assisting the agent to learn faster, achieve better driving performance and increases transparency of the otherwise black-box end-to-end network. We run our experiments to validate the MT-LfD framework in CARLA, an open-source urban driving simulator. We introduce multiple non-player agents in CARLA and induce temporal noise in them for realistic stochasticity.</p>	<p>Rouge scores:</p> <ul style="list-style-type: none"> - Rouge 2: 0.0276 (prec) 0.0374 (recall) 0.0317 (f-1) - Rouge L: 0.1301 (prec) 0.1759 (recall) 0.1496 (f-1) <p>Summary: learning to drive faithfully in highly stochastic urban settings remains an open problem. this paper describes a technique for using Multi-Modal Multi-Task learning that considers multiple behavioral modalities as distinct modes of operation for an end-to-end autonomous deep neural network utilizing the insertion of modal information as secondary input data for tasks with related behaviors. using labeled data from hours of driving our fleet of 1 10th scale model cars, we trained multiple neural networks to imitate the steering angle and driving speed of human control of a car. @cite proposed a multi-modal multi-task learning from demonstration (MT-LfD) framework to train an agent to play Atari games.</p>	<p>Rouge scores:</p> <ul style="list-style-type: none"> - Rouge 2: 0.0192 (prec) 0.0254 (f-1) - Rouge L: 0.0861 (prec) 0.1019 (f-1) <p>Summary: "Machine learning is the learning of a mapping from situations to actions so as to maximize a scalar reward or reinforcement signal. The learner is not told which action to take, as in most forms of machine learning, but instead must discover which actions yield the highest reward by trying them. In the most interesting and challenging cases, actions may affect not only the immediate's reward, but also the next situation, and through that all subsequent rewards. These two characteristics—trial-and-error search and delayed reward—are the two most important distinguishing features of reinforcement learning." "In recent years different lines of evidence have led to the idea that motor actions and movements in both vertebrates and invertebrates are composed of elementary building blocks. The entire motor repertoire can be spanned by applying a well-defined set of operations and transformations to these primitives and by combining them in many different ways according to well-defined syntactic rules. Motor and movement primitives and modules might exist at the neural, dynamic and kinematic levels with complicated mapping among the elementary building blocks subserving these different levels of representation. Hence, while considerable progress has been made in recent years in unravelling the nature of these primitives, new experimental, computational and conceptual approaches are needed to further advance our understanding of motor compositionality.</p>	<p>Rouge scores:</p> <ul style="list-style-type: none"> - Rouge 2: 0.0492 (prec) 0.028 (recall) 0.0357 (f-1) - Rouge L: 0.1774 (prec) 0.1667 (recall) 0.1294 (f-1) <p>Summary: In @cite, an end-to-end deep neural network is used to learn to play Atari 2600 Atari games directly from sensory experience. The network is trained to imitate the steering angle and driving speed of human control of a car. @cite proposed a multi-modal multi-task learning from demonstration (MT-LfD) framework to train an agent to play Atari games.</p>	<p>Rouge scores:</p> <ul style="list-style-type: none"> - Rouge 2: 0.0494 (prec) 0.0374 (recall) 0.0426 (f-1) - Rouge L: 0.1585 (prec) 0.1204 (recall) 0.1368 (f-1) <p>Summary: Multi-task learning (MTL) @cite @cite aims to learn a mapping from situations to actions so as to maximize a scalar reward or reinforcement signal. The learner is not told which action to take, as in most forms of machine learning, but instead must discover which actions yield the highest reward by trying them. In the most interesting and challenging cases, actions may affect not only the immediate's reward, but also the next situation, and through that all subsequent rewards. @cite. comment: - Shows some sign of MDS, summarizing the 2nd (green) and 4th (purple) articles. - However, there is some mix up about the networks for the two articles - The model then mentions the 1st article (yellow) but erroneously treated it as a non-main article (red @cite) and that it is used to play Atari games instead of driving (2nd red part)</p>	<p>Rouge scores:</p> <ul style="list-style-type: none"> - Rouge 2: 0.0541 (prec) 0.0561 (recall) 0.055 (f-1) - Rouge L: 0.1786 (prec) 0.1852 (recall) 0.1818 (f-1) <p>Summary: In @cite, the authors proposed a multi-modal MTL method to train a neural network to imitate the steering angle and driving speed of a car. In @cite, the authors proposed an end-to-end trainable network for imitating the expert demonstrator's driving commands. However, they do not provide a direct supervision mechanism for the agent to learn the driving commands, which is not the case in our case. Moreover, the authors do not propose a method to induce temporal noise in the network. In contrast, in our work, we propose a method that intermediately predicts the driving commands to guide the agent to achieve better driving performance and increase transparency.</p> <p>Comment: - Shows some sign of MDS, contrasting the main article (yellow) with the other articles. - However, the only other article explicitly mentioned is the 4th one (purple).</p>

No.	Abstracts	Label	Base LED	Off-the-shelf LED	Off-the-shelf Centrum	Tuned LED	Tuned Centrum	Two-step	
845	<p>(1): We present a neural framework for opinion summarization from online product reviews which is knowledge-lean and only requires light supervision (e.g., in the form of product domain labels and user-provided ratings). Our method combines two weakly supervised components to identify salient opinions and form extractive summaries from multiple reviews: an aspect extractor trained under a multi-task objective, and a sentiment predictor based on multiple instance learning. We introduce an opinion summarization dataset that includes a training set of product reviews from six diverse domains and human-annotated development and test sets with gold standard aspect annotations, salience labels, and opinion summaries. Automatic evaluation shows significant improvements over baselines, and a large-scale study indicates that our opinion summaries are preferred by human judges according to multiple criteria.</p> <p>(2): The specification discloses a luggage carrier made up of a generally U-shaped frame. The frame has two spaced legs with a hook on the front which hooks over the bumper of an automobile. Two braces are attached to the cross member of the U-shaped member and the front portion of the braces is received on fastening means welded to the under side of the car frame. The cross members provide a supporting surface for carrying articles, boats and the like. A platform may be supported on the frame.</p> <p>(3): We develop a Ranking framework upon Recursive Neural Networks (R2N2) to rank sentences for multi-document summarization. It formulates the sentence ranking task as a hierarchical regression process, which simultaneously measures the salience of a sentence and its constituents (e.g., phrases) in the parsing tree. This enables us to draw on word-level to sentence-level supervisions derived from reference summaries. In addition, recursive neural networks are used to automatically learn ranking features over the tree, with hand-crafted feature vectors of words as inputs. Hierarchical regressions are then conducted with learned features concatenating raw features. Ranking scores of sentences and words are utilized to effectively select informative and non-redundant sentences to generate summaries. Experiments on the DUC 2001, 2002 and 2004 multi-document summarization datasets show that R2N2 outperforms state-of-the-art extractive summarization approaches.</p> <p>(4): We propose a neural multi-document summarization (MDS) system that incorporates sentence relation graphs. We employ a Graph Convolutional Network (GCN) on the relation graphs, with sentence embeddings obtained from Recurrent Neural Networks as input node features. Through multiple layer-wise propagation, the GCN generates high-level hidden sentence features for salience estimation. We then use a greedy heuristic to extract salient sentences while avoiding redundancy. In our experiments on DUC 2004, we consider three types of sentence relation graphs and demonstrate the advantage of combining sentence relations in graphs with the representation power of deep neural networks. Our model improves upon traditional graph-based extractive approaches and the vanilla GRU sequence model with no graph, and it achieves competitive results against other state-of-the-art multi-document summarization systems.</p>	<p>A few extractive neural models have been recently applied to generic multi-document summarization. train a recursive neural network using a ranking objective to identify salient sentences, while follow-up work @cite employs a multi-task objective to improve sentence extraction, an idea we adapted to our task. propose a graph convolution network to represent sentence relations and estimate sentence salience. Our summarization method is tailored to the opinion extraction task, it identifies aspect-specific and salient units, while minimizing the redundancy of the final summary with a greedy selection algorithm @cite. Redundancy is also addressed in who propose a graph-based framework for abstractive summarization. introduce an encoder-decoder neural method for extractive opinion summarization. Their approach requires direct supervision via gold-standard extractive summaries for training, in contrast to our weakly supervised formulation.</p>	<p>Rouge scores:</p> <ul style="list-style-type: none"> - Rouge 2: 0.0732 (prec) 0.0889 (recall) 0.0803 (f-1) <p>- Rouge L:</p> <ul style="list-style-type: none"> 0.1394 (prec) 0.1691 (recall) 0.1528 (f-1) <p>Summary: We present a neural framework for opinion summarization from online product reviews which is knowledge-lean and only requires light supervision (e.g., in the form of product domain labels and user-provided ratings). Our method combines two weakly supervised components to identify salient opinions and form extractive summaries from multiple reviews: an aspect extractor trained under a multi-task objective, and a sentiment predictor based on multiple instance learning. We introduce an opinion summarization dataset that includes a training set of product reviews from six diverse domains and human-annotated development and test sets with gold standard aspect annotations, salience labels, and opinion summaries. Automatic evaluation shows significant improvements over baselines, and a large-scale study indicates that our opinion summaries are preferred by human judges according to multiple criteria.</p> <p>Comment:</p> <ul style="list-style-type: none"> - Copied the first tokens only, except that the phrase "Automatic evaluation shows significant improvements over baselines, and" is skipped before the red part - No indication of MDS 	<p>Rouge scores:</p> <ul style="list-style-type: none"> - Rouge 2: 0.0902 (prec) 0.0815 (recall) 0.0856 (f-1) <p>- Rouge L:</p> <ul style="list-style-type: none"> 0.1707 (prec) 0.1544 (recall) 0.1622 (f-1) <p>Summary: we present a neural framework for opinion summarization from online product reviews which is knowledge-lean and only requires light supervision (e.g., in the form of product domain labels and user-provided ratings). Our method combines two weakly supervised components to identify salient opinions and form extractive summaries from multiple reviews: an aspect extractor trained under a multi-task objective, and a sentiment predictor based on multiple instance learning. We introduce an opinion summarization dataset that includes a training set of product reviews from six diverse domains and human-annotated development and test sets with gold standard aspect annotations, salience labels, and opinion summaries. Automatic evaluation shows significant improvements over baselines, and a large-scale study indicates that our opinion summaries are preferred by human judges according to multiple criteria.</p> <p>Comment:</p> <ul style="list-style-type: none"> - Copied the first tokens only - No indication of MDS 	<p>Rouge scores:</p> <ul style="list-style-type: none"> - Rouge 2: 0.0651 (prec) 0.1037 (recall) 0.08 (f-1) <p>- Rouge L:</p> <ul style="list-style-type: none"> 0.1157 (prec) 0.1838 (recall) 0.142 (f-1) <p>Summary: we present a neural framework for opinion summarization from online product reviews which is knowledge-lean and only requires light supervision (e.g., in the form of product domain labels and user-provided ratings). Our method combines two weakly supervised components to identify salient opinions and form extractive summaries from multiple reviews: an aspect extractor trained under a multi-task objective, and a sentiment predictor based on multiple instance learning. We introduce an opinion summarization dataset that includes a training set of product reviews from six diverse domains and human-annotated development and test sets with gold standard aspect annotations, salience labels, and opinion summaries. Automatic evaluation shows significant improvements over baselines, and a large-scale study indicates that our opinion summaries are preferred by human judges according to multiple criteria.</p> <p>Comment:</p> <ul style="list-style-type: none"> - Copied the first tokens only - No indication of MDS 	<p>Rouge scores:</p> <ul style="list-style-type: none"> - Rouge 2: 0.1695 (prec) 0.0741 (recall) 0.1031 (f-1) <p>- Rouge L:</p> <ul style="list-style-type: none"> 0.3333 (prec) 0.2167 (recall) 0.2041 (f-1) <p>Summary: In the context of multi-document summarization, R2N2 @cite is a recurrent neural network (RNN) for extracting salience labels. The RNN is trained on sentence relation graphs, and the salience labels are extracted using a greedy heuristic to extract salient sentences while avoiding redundancy @cite. In this work, we use the RNN to extract salience labels from product reviews.</p> <p>Comment:</p> <ul style="list-style-type: none"> - Shows clear sign of MDS, summarizing the 3rd (blue) and 4th (purple) articles, while clearly contrasting the main article (yellow). - The RNN notation proved slightly confusing for the model though, the 2nd article used "Recursive Neural Network", not "Recurrent Neural Network" - The 2nd article appears completely unrelated to the task, and the model probably correctly omits it. 	<p>Rouge scores:</p> <ul style="list-style-type: none"> - Rouge 2: 0.1186 (prec) 0.0519 (recall) 0.0722 (f-1) <p>- Rouge L:</p> <ul style="list-style-type: none"> 0.2429 (prec) 0.125 (recall) 0.165 (f-1) <p>Summary: There is a large body of work on multi-document summarization, @cite @cite and @cite. However, most of these works focus on extracting features from a single review. In this work, we focus on extractive summaries from multiple reviews: an aspect extractor trained under a multi-task objective, and a sentiment predictor trained based on multiple instance learning.</p> <p>Comment:</p> <ul style="list-style-type: none"> - Shows some sign of MDS, with the model summarizing what the main article does (yellow), and trying to contrast that with other literature (red) - However, the red part is actually not factually based, and the phrase "a large body of work" is probably picked up from the summarization training. 	<p>Rouge scores:</p> <ul style="list-style-type: none"> - Rouge 2: 0.1594 (prec) 0.0815 (recall) 0.1078 (f-1) <p>- Rouge L:</p> <ul style="list-style-type: none"> 0.2429 (prec) 0.125 (recall) 0.165 (f-1) <p>Summary: In the context of summarization, @cite employ a Graph Convolutional Network (GCN) on the relation graphs, with sentence embeddings obtained from Recurrent Neural Networks as input node features for salience estimation. However, they do not consider sentence relation graphs and do not consider the representation power of deep neural networks. In contrast, we use salience labels and salience labels to identify salient opinions and form extractive summaries from multiple reviews.</p> <p>Comment:</p> <ul style="list-style-type: none"> - Shows some sign of MDS, the model summarizing what the main article does (yellow), and trying to contrast that with other literature (purple) - However, the red-and-purple part contains erroneous facts. The model wrongly assumes that the 4th article did not look at sentence relation graphs when the authors explicitly say they did. 	
4371	<p>(1): In contrast to XML query languages as e.g. XPath which require knowledge on the query language as well as on the document structure, keyword search is open to anybody. As the size of XML sources grows rapidly, the need for efficient search indices on XML data that support keyword search increases. In this paper, we present an approach of XML keyword search which is based on the DAG of the XML data, where repeated substructures are considered only once, and therefore, have to be searched only once. As our performance evaluation shows, this DAG-based extension of the set intersection search algorithm [1], [2], can lead to search times that are on large documents more than twice as fast as the search times of the XML-based approach. Additionally, we utilize a smaller index, i.e., we consume less main memory to compute the results.</p> <p>(2): Recent works have shown the benefits of keyword proximity search in querying XML documents in addition to text documents. For example, given query k</p>	<p>Similar to XRank @cite is the stack-based approach presented in @cite. In contrast to the previous stack-based approaches, the authors do not use the DeweyID to identify a node and to calculate the ancestor-descendant or even parent-child relationships, but they propose to use a combination of preorder position, postorder position, and depth of the node.</p>	<p>Rouge scores:</p> <ul style="list-style-type: none"> - Rouge 2: 0.0244 (prec) 0.069 (recall) 0.036 (f-1) <p>- Rouge L:</p> <ul style="list-style-type: none"> 0.0788 (prec) 0.2203 (recall) 0.1161 (f-1) <p>Summary: In contrast to XML query languages as e.g. XPath which require knowledge on the query language as well as on the document structure, keyword search is open to anybody. In this paper, we consider the problem of efficiently producing ranked results for keyword search queries over hyperlinked XML documents.</p>	<p>Rouge scores:</p> <ul style="list-style-type: none"> - Rouge 2: 0.0 (prec) 0.0 (recall) 0.0 (f-1) <p>- Rouge L:</p> <ul style="list-style-type: none"> 0.1282 (prec) 0.1695 (recall) 0.146 (f-1) <p>Summary: "We consider the problem of efficiently producing ranked results for keyword search queries over hyperlinked XML documents, as opposed to (conceptually) flat</p>	<p>Rouge scores:</p> <ul style="list-style-type: none"> - Rouge 2: 0.0061 (prec) 0.0172 (recall) 0.009 (f-1) <p>- Rouge L:</p> <ul style="list-style-type: none"> 0.0793 (prec) 0.2203 (recall) 0.1166 (f-1) <p>Summary: XRank @cite is an extension of the set intersection search algorithm @cite, which is based on the DAG of the XML data. XRank can be used to query a mix of HTML and XML documents.</p>	<p>Rouge scores:</p> <ul style="list-style-type: none"> - Rouge 2: 0.125 (prec) 0.0517 (recall) 0.0732 (f-1) <p>- Rouge L:</p> <ul style="list-style-type: none"> 0.4 (prec) 0.1695 (recall) 0.2381 (f-1) <p>Summary: XRank @cite is an extension of the set intersection search algorithm @cite. XRank can be used to perform keyword search over hyperlinked XML documents. However, it does not consider</p>	<p>Rouge scores:</p> <ul style="list-style-type: none"> - Rouge 2: 0.04 (prec) 0.0862 (recall) 0.0546 (f-1) <p>- Rouge L:</p> <ul style="list-style-type: none"> 0.1349 (prec) 0.2881 (recall) 0.1838 (f-1) <p>Summary: XRank @cite is a hyperlink based HTML search engine. It is designed to perform keyword search over hyperlinked XML documents. However, it does not consider</p>		

No.	Abstracts	Label	Base LED	Off-the-shelf LED	Off-the-shelf Centrum	Tuned LED	Tuned Centrum	Two-step
	<p>Keywords over Shakespeare's plays in XML, the user might be interested in knowing how the keywords cooccur. In this paper, we focus on XML trees and define XML keyword proximity queries to return the (possibly heterogeneous) set of minimum connecting trees (MCTs) of the matches to the individual keywords in the query. We consider efficiently executing keyword proximity queries on labeled trees (XML) in various settings: 1) when the XML database has been preprocessed and 2) when no indices are available on the XML database. We perform a detailed experimental evaluation to study the benefits of our approach and show that our algorithms considerably outperform prior algorithms and other applicable approaches.</p> <p>(3):</p> <p>We consider the problem of efficiently producing ranked results for keyword search queries over hyperlinked XML documents. Evaluating keyword search queries over hierarchical XML documents, as opposed to (conceptually) flat HTML documents, introduces many new challenges. First, XML keyword search queries do not always return entire documents, but can return deeply nested XML elements that contain the desired keywords. Second, the nested structure of XML implies that the notion of ranking is no longer at the granularity of a document, but at the granularity of an XML element. Finally, the notion of keyword proximity is more complex in the hierarchical XML data model. In this paper, we present the XRank system that is designed to handle these novel features of XML keyword search. Our experimental results show that XRank offers both space and performance benefits when compared with existing approaches. An interesting feature of XRank is that it naturally generalizes a hyperlink based HTML search engine such as Google. XRank can thus be used to query a mix of HTML and XML documents.</p>		<p>n to anybody. As the size of XML sources grows rapidly, the need for efficient search indices on XML data that support keyword search increases. In this paper, we present an approach of XML keyword search which is based on the DAG of the XML data, where repeated substructures are considered only once, and therefore, have to be searched only once.</p> <p>We perform a detailed experimental evaluation to study the benefits of our approach and show that our algorithms considerably outperform prior algorithms and other applicable approaches.</p> <p>Comment:</p> <ul style="list-style-type: none"> - Copied the first tokens only - No indication of MDS 	<p>we present an approach of XML keyword search which is based on the set intersection search algorithm, where repeated substructures are considered only once, and therefore, have to be searched only once.</p> <p>We perform a detailed experimental evaluation to study the benefits of our approach and show that our algorithms considerably outperform prior algorithms and other applicable approaches.</p> <p>Comment:</p> <ul style="list-style-type: none"> - Very weak signs of MDS, with all 3 articles covered. - However, they are all presented as a single document. - The model also performs like an extractive model 	<p>t HTML documents, introduces many new challenges," the authors write. "First, Xie XRank does not always return entire documents, but can return deeply nesting markup that contains the desired keywords. Second, the nesting structure of markup implies that the notion of ranking is no longer at the granularity of a document, but at the granularity of an XML element. Finally, the notion of keyword proximity is more complex in the hierarchical XML data model. In this paper, we present the XRank system that is designed to handle these novel features of XML keyword search. Our experimental results show that XRank offers both space and performance benefits when compared with existing approaches. An interesting feature of XRank is that it naturally generalizes a hyperlink based HTML search engine such as Google. XRank can thus be used to query a mix of HTML and XML documents."</p> <p>Comment:</p> <ul style="list-style-type: none"> - Copied the 3rd article in full - No indication of MDS 	<p>ch engine that can be used to query a mix of hyperlink based XML documents. X-RANK is designed to handle the problem of efficiently producing ranked results for keyword search queries over hyperlinked XML documents, as opposed to (conceptually) flat HTML documents.</p> <p>Comment:</p> <ul style="list-style-type: none"> - Shows very weak signs of MDS, with an attempt to reference multiple @cite. - Even worse than tuned LED - However, none of the sentences generated are coherent, showing difficulties of the model when dealing with articles filled with acronyms (XML, HTML) which it does not understand 	<p>Comment:</p> <ul style="list-style-type: none"> - Shows very weak signs of MDS, with an attempt to reference multiple @cite. - However, none of the sentences generated are coherent, showing difficulties of the model when dealing with articles filled with acronyms (XML, HTML) which it does not understand 	<p>der the nested structure of the XML documents. XRank can be used to query a mix of HTML and XML documents, but it does not require a hyperlink search engine, and it does not need a hyperlink, and it is not required to search for hyperlinks in a hyperlink document. Moreover, XRank is designed to handle the nested structures of the XML document. In contrast, our approach is based on a DAG-based extension of the set intersection search algorithm, which can lead to search times that are more than twice as fast as the search time of the XML-based approach.</p> <p>Comment:</p> <ul style="list-style-type: none"> - Shows signs of MDS, with the yellow part contrasting the approach of the main article with that of the 3rd (blue) - However, the red parts still show hallucination - The first step LED model in the two-step model probably helped simplify the meaning of the individual passages?
4858	<p>(1):</p> <p>It is shown that one can count @math-edge paths in an @math-vertex graph and @math-set @math-packings on an @math-element universe, respectively, in time @math and @math, up to a factor polynomial in @math, @math, and @math; in polynomial space, the bounds hold if multiplied by @math or @math, respectively. These are implications of a more general result: given two set families on an @math-element universe, one can count the disjoint pairs of sets in the Cartesian product of the two families with @math basic operations, where @math is the number of members in the two families and their subsets.</p> <p>(2):</p> <p>We present a fast algorithm for the subset convolution problem: given functions f and g defined on the lattice of subsets of ann-element set n, compute the subset convolution f^*g, defined for $S \subseteq N$ by $[f^*g](S) = [T \subseteq S] f(T)g(S \setminus T)$, where addition and multiplication is carried out in an arbitrary ring. Via Möbius transform and inversion, our algorithm evaluates the subset convolution in $O(n^2 2n)$ additions and multiplications, substantially improving upon the straightforward $O(3n)$ algorithm. Specifically, if the input functions have an integer range $[-M, M+1, \dots, M]$, their subset convolution over the ordinary sum-product ring can be computed in $O(2n \log M)$ time; the notation O suppresses polylogarithmic factors. Furthermore, using a standard embedding technique we can compute the subset convolution over the max-sum or min-sum semiring in $O(2nM)$ time. To demonstrate the applicability of fast subset convolution, we present the first $O(2kn^2 + nm)$ algorithm for the Steiner tree problem in graphs with n vertices, k terminals, and m edges with bounded integer weight, improving upon the $O(3kn + 2kn^2 + nm)$ time bound of the classical Dreyfus-Wagner algorithm. We also discuss extensions to recent $O(2n)$-time algorithms for covering and partitioning problems (Björklund and Husfeldt, FOCS 2006, Koivisto, FOCS 2006).</p> <p>(3):</p> <p>Given a set @math with @math elements and a family @math of subsets, we show how to partition @math into @math such subsets in @math time. We also consider variations of this problem where the subsets may overlap or are weighted, and we solve the decision, counting, summation, and optimization versions of these problems. Our algorithms are based on the principle of inclusion-exclusion and the zeta transform. In effect we get exact algorithms in @math time for several well-studied partition problems including domatic number, chromatic number, maximum @math-cut, bin packing, list coloring, and the chromatic polynomial. We also have applications to Bayesian learning with decision graphs and to model-based data clustering. If only polynomial space is available, our algorithms run in time @math if membership in @math can be decided in polynomial time. We solve chromatic number in @math time and domatic number in @math time. Finally, we present a family of polynomials</p>	<p>Concerning set packings the situation is analogous, albeit the research has been somewhat less extensive. Deciding whether a given family of @math subsets of a @math-element universe contains a @math-packing is known to be W[1]-hard @cite, and thus it is unlikely that the problem is fixed parameter tractable, that is, solvable in time @math for some function @math and constant α. If @math is fairly large, say exponential in @math, the fastest known algorithms actually count the packings by employing the inclusion-exclusion machinery @cite and run in time @math. This bound holds also for the presented algorithm (cf. Theorem).</p>	<p>Rouge scores:</p> <ul style="list-style-type: none"> - Rouge 2: 0.1079 (prec) 0.1442 (recall) 0.1235 (f-1) - Rouge L: 0.1714 (prec) 0.2286 (recall) 0.1959 (f-1) <p>Summary:</p> <p>It is shown that one can count @math-edge paths in an @math-vertex graph and @math-set @math-packings on an @math-element universe, respectively, in time @math and @math, up to a factor polynomial in @math, @math, and @math; in polynomial space, the bounds hold if multiplied by @math or @math, respectively. These are implications of a more general result: given two set families on an @math-element universe, one can count the disjoint pairs of sets in the Cartesian product of the two families with @math basic operations, where @math is the number of members in the two families and their subsets.</p>	<p>Rouge scores:</p> <ul style="list-style-type: none"> - Rouge 2: 0.0612 (prec) 0.0577 (recall) 0.0594 (f-1) - Rouge L: 0.2222 (prec) 0.2095 (recall) 0.2157 (f-1) <p>Summary:</p> <p>given a set @math with @math elements and a family @math of subsets, we show how to partition @math into @math such subsets in @math time.</p>	<p>Rouge scores:</p> <ul style="list-style-type: none"> - Rouge 2: 0.0791 (prec) 0.1346 (recall) 0.0996 (f-1) - Rouge L: 0.1348 (prec) 0.2286 (recall) 0.1696 (f-1) <p>Summary:</p> <p>It is shown that one can count @math-edge paths in an @math-vertex graph and @math-set @math-packings on an @math-element universe, respectively, in time @math and @math, up to a factor polynomial in @math, @math, and @math; in polynomial space, the bounds hold if multiplied by @math or @math, respectively. These are implications of a more general result: given two set families on an @math-element universe, one can count the disjoint pairs of sets in the Cartesian product of the two families with @math basic operations, where @math is the number of members in the two families and their subsets.</p>	<p>Rouge scores:</p> <ul style="list-style-type: none"> - Rouge 2: 0.125 (prec) 0.0769 (recall) 0.0952 (f-1) - Rouge L: 0.2923 (prec) 0.181 (recall) 0.2235 (f-1) <p>Summary:</p> <p>For the Steiner tree problem @cite, one can count the disjoint pairs of sets in the Cartesian product of the two families with @math basic operations, where @math is the number of members in the two families and their subsets. They also showed how to partition the set into @math such subsets in @math time, where @math and @math are the subsets of @math. They also gave a family of polynomial space approximation algorithms that find a number between @math and the @math in @math. @cite gave a fast algorithm for the subset convolution problem in graphs with @math vertices, @math terminals, and @math edges with bounded integer weights, improving upon the classical Dreyfus-Wagner algorithm. They also discussed extensions to recent @math-time algorithms for covering and partitioning problems.</p> <p>Comment:</p> <ul style="list-style-type: none"> - No indication of MDS behavior, basically just mixing in some words of the 2nd article (green) with the 1st (yellow) - Shows once again how the model fails to understand words used in unusual settings (e.g. mathematics) 	<p>Rouge scores:</p> <ul style="list-style-type: none"> - Rouge 2: 0.082 (prec) 0.0962 (recall) 0.0885 (f-1) - Rouge L: 0.1951 (prec) 0.2286 (recall) 0.2105 (f-1) <p>Summary:</p> <p>In @cite, Björklund and Husfeldt showed that one can partition a set @math into @math subsets in @math time, where @math is the number of members in the two families and their subsets. They also showed how to partition the set into @math such subsets in @math time, where @math and @math are the subsets of @math. They also gave a family of polynomial space approximation algorithms that find a number between @math and the @math in @math. @cite gave a fast algorithm for the subset convolution problem in graphs with @math vertices, @math terminals, and @math edges with bounded integer weights, improving upon the classical Dreyfus-Wagner algorithm. They also discussed extensions to recent @math-time algorithms for covering and partitioning problems.</p> <p>Comment:</p> <ul style="list-style-type: none"> - Best performing of all 3 tuned models, showing strong signs of MDS with the 2nd (green) and 3rd (blue) articles reflected. - The red part appears to contain some factual errors, with the original "O(3kn + 2kn^2 + nm)" changed to "O(3n - @math)" - The first step LED model in the two-step model probably 	

No.	Abstracts	Label	Base LED	Off-the-shelf LED	Off-the-shelf Centrum	Tuned LED	Tuned Centrum	Two-step
	I space approximation algorithms that find a number between @math and @math in time @math .				Comment: - Copied the first tokens only - No indication of MDS - Random word change noted (red)		words used in unusual settings (e.g. mathematics)	helped simplify the meaning of the individual passages?
5068	<p>(1): Convolutional neural networks have gained a remarkable success in computer vision. However, most usable network architectures are hand-crafted and usually require expertise and elaborate design. In this paper, we provide a block-wise network generation pipeline called BlockQNN which automatically builds high-performance networks using the Q-Learning paradigm with epsilon-greedy exploration strategy. The optimal network block is constructed by the learning agent which is trained sequentially to choose component layers. We stack the block to construct the whole auto-generated network. To accelerate the generation process, we also propose a distributed asynchronous framework and an early stop strategy. The block-wise generation brings unique advantages: (1) it performs competitive results in comparison to the hand-crafted state-of-the-art networks on image classification, additionally, the best network generated by BlockQNN achieves 3.54 top-1 error rate on CIFAR-10 which beats all existing auto-generate networks. (2) in the meanwhile, it offers tremendous reduction of the search space in designing networks which only spends 3 days with 32 GPUs, and (3) moreover, it has strong generalizability that the network built on CIFAR also performs well on a larger-scale ImageNet dataset.</p> <p>(2): This paper introduces the application of gradient descent methods to meta-learning. The concept of "meta-learning", i.e. of a system that improves or discards a learning algorithm, has been of interest in machine learning for decades because of its appealing applications. Previous meta-learning approaches have been based on evolutionary methods and, therefore, have been restricted to small models with few free parameters. We make meta-learning in large systems feasible by using recurrent neural networks with attendant learning routines as meta-learning systems. Our system derived complex well performing learning algorithms from scratch. In this paper we also show that our approach performs non-stationary time series prediction.</p> <p>(3): Different researchers hold different views of what the term meta-learning exactly means. The first part of this paper provides our own perspective view in which the goal is to build self-adaptive learners (i.e. learning algorithms that improve their bias dynamically through experience by accumulating meta-knowledge). The second part provides a survey of meta-learning as reported by the machine-learning literature. We find that, despite different views and research lines, a question remains constant: how can we exploit knowledge about learning (i.e. meta-knowledge) to improve the performance of learning algorithms? Clearly the answer to this question is key to the advancement of the field and continues being the subject of intensive research.</p> <p>(4): The move from hand-designed features to learned features in machine learning has been wildly successful. In spite of this, optimization algorithms are still designed by hand. In this paper we show how the design of an optimization algorithm can be cast as a learning problem, allowing the algorithm to learn to exploit structure in the problems of interest in an automatic way. Our learned algorithms, implemented by LSTMs, outperform generic, hand-designed competitors on the tasks for which they are trained, and also generalize well to new tasks with similar structure. We demonstrate this on a number of tasks, including simple convex problems, training neural networks, and styling images with neural art.</p>	Another bunch of related works include hyper-parameter optimization @cite , meta-learning @cite and learning to learn methods @cite @cite . However, the goal of these works is to use meta-data to improve the performance of the existing algorithms, such as finding the optimal learning rate of optimization methods or the optimal number of hidden layers to construct the network. In this paper, we focus on learning the entire topological architecture of network blocks to improve the performance.	Rouge scores: - Rouge 2: 0.0596 (prec) 0.1169 (recall) 0.0789 (f-1) - Rouge L: 0.1053 (prec) 0.2051 (recall) 0.1391 (f-1)	Rouge scores: - Rouge 2: 0.0723 (prec) 0.0779 (recall) 0.075 (f-1) - Rouge L: 0.1548 (prec) 0.1667 (recall) 0.1605 (f-1)	Rouge scores: - Rouge 2: 0.0513 (prec) 0.1299 (recall) 0.0735 (f-1) - Rouge L: 0.0969 (prec) 0.2436 (recall) 0.1387 (f-1)	Rouge scores: - Rouge 2: 0.0811 (prec) 0.039 (recall) 0.0526 (f-1) - Rouge L: 0.1698 (prec) 0.0897 (recall) 0.1207 (f-1)	Rouge scores: - Rouge 2: 0.0952 (prec) 0.1299 (recall) 0.1099 (f-1) - Rouge L: 0.1698 (prec) 0.2308 (recall) 0.1957 (f-1)	Rouge scores: - Rouge 2: 0.0899 (prec) 0.1039 (recall) 0.0964 (f-1) - Rouge L: 0.2111 (prec) 0.2436 (recall) 0.2262 (f-1)

(c) Long samples

- 485: sample where the strong performance of the tuned and two-step models can be seen

No.	Abstracts	Label	Base LED	Off-the-shelf LED	Off-the-shelf Centrum	Tuned LED	Tuned Centrum	Two-step
485	<p>(1): Learning from web data has attracted lots of research interest in recent years. However, crawled web images usually have two types of noises, label noise and background noise, which induce extra difficulties in utilizing them effectively. Most existing methods either rely on human supervision or ignore the background noise. In this paper, we propose a novel method, which is capable of h</p>	In learning classifier with web data, previous works focus on handling the label noise in three directions, removing label noise @cite @cite @cite @cite @cite , building noise-robust m	Rouge scores: - Rouge 2: 0.0175 (prec) 0.0732 (recall) 0.0283 (f-1) - Rouge L: 0.0283	Rouge scores: - Rouge 2: 0.0258 (prec) 0.0976 (recall) 0.0408 (f-1) - Rouge L: 0.0244	Rouge scores: - Rouge 2: 0.0048 (prec) 0.0244 (recall) 0.0081 (f-1) - Rouge L: 0.0081	Rouge scores: - Rouge 2: 0.0485 (prec) 0.122 (recall) 0.0694 (f-1) - Rouge L: 0.0437	Rouge scores: - Rouge 2: 0.0266 (prec) 0.122 (recall) 0.0447 (f-1) - Rouge L: 0.0447	Rouge scores: - Rouge 2: 0.0287 (prec) 0.1463 (recall) 0.048 (f-1) - Rouge L: 0.048

No.	Abstracts	Label	Base LED	Off-the-shelf LED	Off-the-shelf Centrum	Tuned LED	Tuned Centrum	Two-step
	<p>Handling these two types of noises together, without the supervision of clean images in the training stage. Particularly, we formulate our method under the framework of multi-instance learning by grouping ROIs (i.e., images and their region proposals) from the same category into bags. ROIs in each bag are assigned with different weights based on the representative discriminative scores of their nearest clusters, in which the clusters and their scores are obtained via our designed memory module. Our memory module could be naturally integrated with the classification module, leading to an end-to-end trainable system. Extensive experiments on four benchmark datasets demonstrate the effectiveness of our method.</p> <p>(2):</p> <p>We study the problem of automatically removing outliers from noisy data, with application for removing outlier images from an image collection. We address this problem by utilizing the reconstruction errors of an autoencoder. We observe that when data are reconstructed from low-dimensional representations, the inliers and the outliers can be well separated according to their reconstruction errors. Based on this basic observation, we gradually inject discriminative information in the learning process of an autoencoder to make the inliers and the outliers more separable. Experiments on a variety of image datasets validate our approach.</p> <p>(3):</p> <p>We present a theoretically grounded approach to train deep neural networks, including recurrent networks, subject to class-dependent label noise. We propose two procedures for loss correction that are agnostic to both application domain and network architecture. They simply amount to at most a matrix inversion and multiplication, provided that we know the probability of each class being corrupted into another. We further show how one can estimate these probabilities, adapting a recent technique for noise estimation to the multi-class setting, and thus providing an end-to-end framework. Extensive experiments on MNIST, IMDB, CIFAR-10, CIFAR-100 and a large scale dataset of clothing images employing a diversity of architectures &#x2014; stacking dense, convolutional, pooling, dropout, batch normalization, word embedding, LSTM and residual layers &#x2014; demonstrate the noise robustness of our proposal s. Incidentally, we also prove that, when ReLU is the only non-linearity, the loss curvature is immune to class-dependent label noise.</p> <p>(4):</p> <p>In this paper, we study the problem of learning image classification models with label noise. Existing approaches depending on human supervision are generally not scalable as manually identifying correct or incorrect labels is time-consuming, whereas approaches not relying on human supervision are scalable but less effective. To reduce the amount of human supervision for label noise cleaning, we introduce CleanNet, a joint neural embedding network, which only requires a fraction of the classes being manually verified to provide the knowledge of label noise that can be transferred to other classes. We further integrate CleanNet and conventional convolutional neural network classifier into one framework for image classification learning. We demonstrate the effectiveness of the proposed algorithm on both of the label noise detection task and the image classification on noisy data task on several large-scale datasets. Experimental results show that CleanNet can reduce label noise detection error rate on held-out classes where no human supervision available by 41.5 compared to current weakly supervised methods. It also achieves 47 of the performance gain of verifying all images with only 3.2 images verified on an image classification task. Source code and dataset will be available at kuanghuei.github.io/CleanNetProject.</p> <p>(5):</p> <p>We present a simple yet efficient approach capable of training deep neural networks on large-scale weakly-supervised web images, which are crawled raw from the Internet by using text queries, without any human annotation. We develop a principled learning strategy by leveraging curriculum learning, with the goal of handling a massive amount of noisy labels and data imbalance effectively. We design a new learning curriculum by measuring the complexity of data using its distribution density in a feature space, and rank the complexity in an unsupervised manner. This allows for an efficient implementation of curriculum learning on large-scale web images, resulting in a high-performance CNN model, where the negative impact of noisy labels is reduced substantially. Importantly, we show by experiments that those images with highly noisy labels can surprisingly improve the generalization capability of model, by serving as a manner of regularization. Our approaches obtain state-of-the-art performance on four benchmarks: WebVision, ImageNet, Clothing-1M and Food-101. With an ensemble of multiple models, we achieved a top-5 error rate of 5.2 on the WebVision challenge [18] for 1000-category classification. This result was the top performance by a wide margin, outperforming second place by a nearly 50 relative error rate. Code and models are available at: https://github.com/MalongTech/CurriculumNet.</p> <p>(6):</p> <p>Learning from web data is increasingly popular due to abundant free web resources. However, the performance gap between webly supervised learning a</p>	<p>odel @cite @cite @cite @cite @cite @cite , and curriculum learning @cite .</p>	<p>0.064 (prec) 0.2619 (recall) 0.1028 (f-1)</p> <p>Summary: Learning from web data has attracted lots of research interest in recent years. However, it is still very difficult to learn from web data without the supervision of a human. In this paper, we study the problem of learning from web data. The problem is that web images usually have two types of noises, label noise and background noise, which induce extra difficulties in utilizing them effectively. Most existing methods either rely on human supervision or ignore the background noise. In this article, we propose a novel method, which is capable of handling these two types of noise together, without the supervision or supervision of clean images in the training stage. Particularly, we formulate our method under the framework of multi-instance learning by grouping ROIs (i.e., images and their region proposals) from the same category into bags. ROIs in each bag are assigned with different weights based on the representative discriminative scores of their nearest clusters, in which the clusters and their scores are obtained via our designed memory module.</p> <p>Comment: - Copied the first tokens only - No indication of MDS</p>	<p>0.0705 (prec) 0.2619 (recall) 0.1111 (f-1)</p> <p>Summary: learning from web data is increasingly popular due to abundant free web resources. However, the performance gap between webly supervised learning and traditional supervised learning is still very large, due to the label noise of web data as well as the domain shift between web data and test data. To fill this gap, most existing methods propose to purify or augment web data using instance-level supervision, which generally requires heavy annotation. Instead, we propose to address the label noise and domain shift by using more accessible category-level supervision. In particular, we build our deep probabilistic framework upon variational autoencoder (VAE), in which classification network and VAE can jointly leverage category-level hybrid information. Our memory module could be naturally integrated with the classification module, leading to an end-to-end trainable system. we design a new learning curriculum by measuring the complexity of data using its distribution density, in an unsupervised manner in an</p> <p>Comment: - Copied the 6th article, except replacing the last 2 sentences with another sentence from the 5th article (orange) - No indication of MDS</p>	<p>0.0433 (prec) 0.2143 (recall) 0.072 (f-1)</p> <p>Summary: We present a simple yet efficient approach capable of training deep neural networks on large-scale weakly-supervised web images, which are crawled raw from the Internet by using text queries, without any human annotation. We develop a principled learning strategy by leveraging curriculum learning, with the goal of handling a massive amount of noisy labels and data imbalance effectively. We design a new learning curriculum by measuring the complexity of data using its distribution density in a feature space, and rank the complexity in an unsupervised manner. This allows for an efficient implementation of curriculum learning on large-scale web images, resulting in a high-performance CNN model, where the negative impact of noisy labels is reduced substantially. Importantly, we show by experiments that those images with highly noisy labels can surprisingly improve the generalization capability of model, by serving as a manner of regularization. Our approaches obtain state-of-the-art performance on four benchmarks: WebVision, ImageNet, Clothing-1M and Food-101. With an ensemble of multiple models, we achieved a top-5 error rate of 5.2 on the WebVision challenge [18] for 1000-category classification. This result was the top performance by a wide margin, outperforming second place by a nearly 50 relative error rate. Code and models are available at: https://github.com/MalongTech/CurriculumNet.</p> <p>Comment: - Copied the 5th article only - No indication of MDS</p>	<p>0.1058 (prec) 0.2619 (recall) 0.1507 (f-1)</p> <p>Summary: Label noise is an important issue in image classification, with many potential negative consequences. For example, the label noise can be class-conditional, which is not suitable for large-scale weakly-supervised learning. To address this problem, several methods have been proposed. @cite proposed a method for learning from unlabeled web data, which is based on the assumption that the unlabeled data can be used to improve the classification performance. However, these methods do not consider the label noise of web data. In contrast, our method is able to deal with the label noise without the supervision of the classification module.</p> <p>Comment: - Some signs of MDS, using the context information from the 1st article (yellow) as background, and also mentioned correctly some elements of the current study - Elements of 4th article (purple) is clear while the 1st red part is true for many samples.</p>	<p>0.0688 (prec) 0.3095 (recall) 0.1126 (f-1)</p> <p>Summary: The label noise problem has been widely studied in the literature. In @cite, the label noise is treated as a set of instances that are mislabeled, and the label cleansing algorithm is proposed to deal with the label noise. In addition, @cite proposed an active learning approach for removing outliers from an image collection @cite, which does not require a large amount of training data to train the model, and they do not need any annotated labels for training. Moreover, the label noise cleansing algorithms are defined as follows: first, the labels of the instances are labeled, and second, the labels on the labels are labeled respectively. The labels are labeled by the label noise correction algorithm and the labels of mislabeled instances are labeled by label noise cleansing. @cite used curriculum learning to train a deep neural network on large-scale web images, which are crawled raw from the Internet by using text queries, without any human annotation. However, curriculum learning can significantly improve the generalization capability of the model, by reducing the number of noisy labels, and by using curriculum learning to reduce the amount of data imbalance in the training process. Moreover, curriculum learning is effective in handling noisy labels and the data imbalance effectively.</p> <p>Comment: - Clear signs of MDS, with the 1st article (yellow) used for both context (at the beginning) and contrast with other studies in the latter half - 3 other articles are cited, though the red parts shows a poorly written sentence.</p>	<p>0.0667 (prec) 0.3333 (recall) 0.1111 (f-1)</p> <p>Summary: Recently, there has been a lot of research on label noise cleaning. In @cite, the label noise is treated as a set of instances that are mislabeled, and the label cleansing algorithm is proposed to deal with the label noise. In addition, @cite proposed an active learning approach for removing outliers from an image collection @cite, which does not require a large amount of training data to train the model, and they do not need any annotated labels for training. Moreover, the label noise cleansing algorithms are defined as follows: first, the labels of the instances are labeled, and second, the labels on the labels are labeled respectively. The labels are labeled by the label noise correction algorithm and the labels of mislabeled instances are labeled by label noise cleansing. @cite used curriculum learning to train a deep neural network on large-scale web images, which are crawled raw from the Internet by using text queries, without any human annotation. However, curriculum learning can significantly improve the generalization capability of the model, by reducing the number of noisy labels, and by using curriculum learning to reduce the amount of data imbalance in the training process. Moreover, curriculum learning is effective in handling noisy labels and the data imbalance effectively.</p> <p>Comment: - Some signs of MDS, with the 1st article (yellow) used for context, and the 2nd (green), 7th (pink) and 8th (deep blue) articles referenced. - Contains much hallucinations though as noted in the red parts.</p>

No.	Abstracts	Label	Base LED	Off-the-shelf LED	Off-the-shelf Centrum	Tuned LED	Tuned Centrum	Two-step
	<p>nd traditional supervised learning is still very large, due to the label noise of w eb data as well as the domain shift between web data and test data. To fill thi s gap, most existing methods propose to purify or augment web data using in stance-level supervision, which generally requires heavy annotation. Instead, we propose to address the label noise and domain shift by using more accessi ble category-level supervision. In particular, we build our deep probabilistic fr amework upon variational autoencoder (VAE), in which classification network and VAE can jointly leverage category-level hybrid information. Then, we ext end our method for domain adaptation followed by our low-rank refinement strategy. Extensive experiments on three benchmark datasets demonstrate th e effectiveness of our proposed method.</p> <p>(7):</p> <p>Label noise is an important issue in classification, with many potential negativ e consequences. For example, the accuracy of predictions may decrease, whe reas the complexity of inferred models and the number of necessary training samples may increase. Many works in the literature have been devoted to the study of label noise and the development of techniques to deal with label noi se. However, the field lacks a comprehensive survey on the different types of label noise, their consequences and the algorithms that consider label noise. This paper proposes to fill this gap. First, the definitions and sources of label n oise are considered and a taxonomy of the types of label noise is proposed. S econd, the potential consequences of label noise are discussed. Third, label n oise-robust, label noise cleansing, and label noise-tolerant algorithms are revi ewed. For each category of approaches, a short discussion is proposed to help the practitioner to choose the most suitable technique in its own particular fi eld of application. Eventually, the design of experiments is also discussed, wh at may interest the researchers who would like to test their own algorithms.</p> <p>In this paper, label noise consists of mislabeled instances: no additional inform ation is assumed to be available like e.g., confidences on labels.</p> <p>(8):</p> <p>In this paper, we study a classification problem in which sample labels are ran domly corrupted. In this scenario, there is an unobservable sample with noise -free labels. However, before being observed, the true labels are independent ly flipped with a probability α, and the random label noise can be class- conditional. Here, we address two fundamental problems raised by this scena rio. The first is how to best use the abundant surrogate loss functions designe d for the traditional classification problem when there is label noise. We prov e that any surrogate loss function can be used for classification with noisy lab els by using importance reweighting, with consistency assurance that the labe l noise does not ultimately hinder the search for the optimal classifier of the n oise-free sample. The other is the open problem of how to obtain the noise ra te α. We show that the rate is upper bounded by the conditional proba bility α of the noisy sample. Consequently, the rate can be estimated, b ecause the upper bound can be easily reached in classification problems. Exp erimental results on synthetic and real datasets confirm the efficiency of our methods.</p> <p>(9):</p> <p>Current approaches for fine-grained recognition do the following: First, recrui t experts to annotate a dataset of images, optionally also collecting more stru ctured data in the form of part annotations and bounding boxes. Second, trai n a model utilizing this data. Toward the goal of solving fine-grained recogniti on, we introduce an alternative approach, leveraging free, noisy data from th e web and simple, generic methods of recognition. This approach has benefits in both performance and scalability. We demonstrate its efficacy on four fine -grained datasets, greatly exceeding existing state of the art without the man ual collection of even a single label, and furthermore show first results at scali ng to more than 10,000 fine-grained categories. Quantitatively, we achieve to p-1 accuracies of (92.3,) on CUB-200-2011, (85.4,) on Birdsnap, (93.4,) on F GVC-Aircraft, and (80.8,) on Stanford Dogs without using their annotated tra ining sets. We compare our approach to an active learning approach for expan ding fine-grained datasets.</p>							