

Curriculum Data Augmentation for Low-Resource Slides Summarization

Anonymous ACL submission

Abstract

Data augmentation is commonly used in training in low-resource scenarios. However, there are sometimes large discrepancy between distributions of augmented data and target data. How to bridge the gap between the augmented and target data, especially when target data is harder-to-learn? In this paper, we study improved data augmentation strategies in the scenario of scientific slides text summarization, where we generate a textual summary based on texts of presentation slides. Since slides are messy and difficult to understand by current models, we introduce an easier form of data, i.e., articles in natural language. The basic idea is that we generate the transition data between slides and articles, and all three of them form a curriculum for neural models to learn the distribution transition from article data to slides data. We find that our approach achieves consistent improvements over different backbone summarization models. The curriculum-oriented data augmentation method can generate data that fill the gap between the easy-to-obtain data and the low-resource task data. We show that curriculum learning and data augmentation can be combined to help NLP models learn from otherwise hard-to-learn data. ¹

1 Introduction

Nowadays, presentation slides have become one of the main materials for disseminating ideas and thoughts in conferences, lectures or events. Similar to other formats of documents such as articles, presentation slides grow rapidly in volume. The difficulty for readers to exhaustively go through all contents for so many slides calls for the automatic summarization of such documents. With textual summaries for presentation slides, readers can easily find ones they are interested in and then dive into details.

Most existing works focus on summarizing articles for large-scale datasets, e.g., CNN/Daily Mail

¹We will release the code after paper’s acceptance.

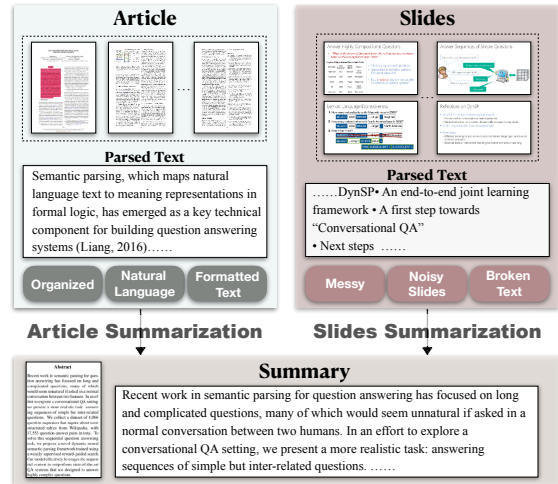


Figure 1: Difference between slides summarization and the corresponding article summarization.

(Hermann et al., 2015), XSum (Narayan et al., 2018) and PubMed/ArXiv (Cohan et al., 2018), leaving the summarization of slides an unsolved puzzle. Traditional extractive text summarization methods (Mihalcea and Tarau, 2004; Zhang et al., 2018; Liu, 2019). do not trivially apply to noisy and broken text as slides, as seen in Figure 1. On the other hand, abstractive methods are developed to generate summaries for given text (Rush et al., 2015; See et al., 2017; Gehrmann et al., 2018). In recent years, summarization with large-scale pre-trained language models (PLMs) (Devlin et al., 2019; Raffel et al., 2020; Lewis et al., 2020; Zhang et al., 2020) has become the dominant paradigm, due to their excellent ability in language modeling. However, slides differ greatly from regular text these models are pre-trained or fine-tuned on, leaving the potential of the PLMs not fully realized.

To fill in the blank, we propose the task of slides summarization in this work. There are two major challenges in this task: *limited resource* and *noisy input*. Summarizing slides is a low-resource task by nature, because we can not utilize any section

of the slides as summaries directly for building large-scale datasets. In contrast, article summarization often uses the first paragraph or abstract of an article as ground-truth summary. Moreover, the rich multi-modal information (i.e., textual, visual and layout information (Xu et al., 2020b, 2021)) in the source documents (i.e., presentation slides) are difficult to be utilized for current summarization techniques. When parsed into text, the source documents are much more noisy and ill-formatted than articles, posing great challenges for language models to understand in a low-resource setting.

Instead of focusing on designing summarization models, we tackle this task from the perspective of data augmentation. The intuition is that language models can learn better representations of unfamiliar slides if they start with something familiar, i.e., articles in natural language. We are motivated by the idea of *curriculum learning* (CL) (Bengio et al., 2009), where models benefit from starting small and gradually increasing the learning difficulty when training with limited data samples (Wang et al., 2021; Wu et al., 2020). In this sense, the accompanying articles facilitate the learning of the model by providing connections between what it knows well (natural language) and what it does not (noisy slides).

In this paper, we propose the LESSON framework for **Low-rEsource Slides SummarizatiON** with curriculum data augmentation. Different from traditional CL (Bengio et al., 2009; Liu et al., 2018; Platanios et al., 2019), we do not discover curriculum data within training data. Instead, we build extra curriculum data from the accompanying articles of slides. To bridge the gap between data distributions of articles and slides, we enrich the curriculum by generating in-between data samples, yielding a sequence of data samples to language models in an increasing order of difficulty. We adopt a simple generative model to realize the transition from articles to slides, which is controlled by a balancing weight in training objectives. Like CL, our method focuses on optimizing training and thus does not need such corresponding articles for slides in the test time. In this sense, LESSON is practical, because we do not usually have corresponding articles for test-time slides.

In summary, the contributions of this paper include: 1) We are the first to propose the slides summarization problem, and tackle low-resource scientific slides summarization as a concrete task;

2) We propose LESSON framework for this task, a training strategy with curriculum data augmentation, in order to learn a robust model for slides summarization; 3) We conduct extensive experiments to verify and understand the effectiveness of the LESSON framework.

2 Related Work

Data Augmentation In NLP, the goals of data augmentation include increasing training data size (Fadaee et al., 2017), achieving regularizing effects (Hernández-García and König, 2020; Fabbri et al., 2021), and diversifying data (Lu et al., 2020; Kumar et al., 2019). Common data augmentation approaches include rule-based methods, interpolation methods and model-based methods (Feng et al., 2021). Rule-based approaches utilize a set of pre-defined rules to transform existing data and create new samples (Wei and Zou, 2019). Interpolation approaches construct a continuous hidden space for text and interpolate new data from that hidden space (Chen et al., 2020; Cheng et al., 2020). Model-based approaches either leverage a pre-trained language model (Anaby-Tavor et al., 2020; Yang et al., 2020) or train an auxiliary model based on existing training data (Kumar et al., 2019), and then use the model to generate new samples.

Curriculum Learning Curriculum learning (Bengio et al., 2009) is a learning strategy that mimics the human learning process. Theoretically, curriculum can be regarded as an optimization strategy (Bengio et al., 2009) for non-convex training criteria. Starting from easier samples, the model can learn a smoother training objective, and approximate local minima that have better generalization ability towards the global minima (Wang et al., 2021). Previous empirical results have demonstrated the strengths of curriculum learning in a wide range of NLP tasks, such as question answering (Liu et al., 2018), natural language understanding (Xu et al., 2020a), machine translation (Platanios et al., 2019), and text classification (Wei et al., 2021). These works have testified to curriculum learning’s ability to learn from noisier data (Wu et al., 2020), reduce training time (Wu et al., 2020; Platanios et al., 2019), and gain performance improvements over randomized training (Liu et al., 2018; Platanios et al., 2019; Xu et al., 2020a).

Curriculum Data Augmentation Previously, Wei et al. (2021) explore curriculum data augmentation in a few-shot setting. They use rule-based approaches (Wei and Zou, 2019) to augment the original dataset through controlled noising, i.e., creating noisier and more difficult data for curriculum learning. In contrast, we employ curriculum learning to make the optimization process smoother and improve the model’s performance on the hard data. Both of our studies represent a new approach in curriculum learning: instead of discovering curriculum from existing data, we artificially create data of different difficulty levels.

Abstractive Summarization The goal of abstractive summarization is to generate concise and precise summary text for documents. Traditionally, such methods mostly applies sequence-to-sequence encoder-decoder architectures (Rush et al., 2015; See et al., 2017; Gehrmann et al., 2018; You et al., 2019). In recent years, pre-trained transformer-based (Vaswani et al., 2017) language models have achieves remarkable success on a wide range of NLP tasks (Radford et al., 2019; Devlin et al., 2019; Liu et al., 2019b; Raffel et al., 2020; Lewis et al., 2020). Many approaches using pre-trained language models have achieved state-of-the-art results on abstractive summarization tasks (Liu and Lapata, 2019; Rothe et al., 2020; Zhang et al., 2020; Beltagy et al., 2020). Most summarization studies work on text summarization, and some focus on summarizing from noisier input such as meeting transcripts and dialogues (Liu et al., 2019a; Zhao et al., 2019; Liu et al., 2019c; Zhu et al., 2020). In this work, we directly use parsed text as the input due to the difficulty in slides structure parsing. Thus, our input is even noisier with structured information mixed in plain text. As our proposed LESSON is model-agnostic, we leave the handling of slides structures and layouts to future work.

3 Proposed Approach

In this section, we detail the slides summarization task and the proposed LESSON method, a curriculum-based data augmentation training framework for slides summarization.

3.1 Task Formulation

Let \mathcal{X}^S denote slides text and \mathcal{Y} denote summary text, and we have $\mathbb{D}^S = (\mathcal{X}^S, \mathcal{Y})$ for a slides summarization dataset. We hypothesize that slides data are harder to learn compared with normal article

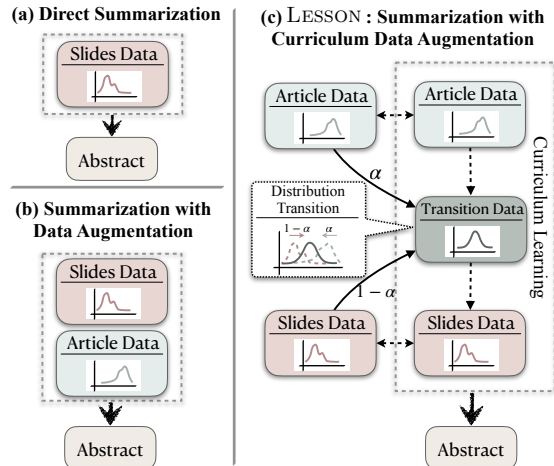


Figure 2: Comparison of (a) Direct summarization, input \mathcal{X}^S slides, output \mathcal{Y} summary and (b) Data augmentation with \mathcal{X}^A article and (c) LESSON. The acquisition of transition data $\mathcal{X}^{A \rightarrow S}$ is shown.

data, given limited parsing techniques and language models that are easier to process natural language.

Inspired by previous studies on curriculum learning, which demonstrate that learning from easier samples helps a model to learn harder samples, we introduce paralleled article data \mathcal{X}^A as augmentation, denoted as $\mathbb{D}^A = (\mathcal{X}^A, \mathcal{Y})$. Thus we have paralleled data triples $(\mathcal{X}^S, \mathcal{X}^A, \mathcal{Y})$. Since the distribution of the parallel articles is greatly different from target slides data, we further generate transition data to bridge the distributional discrepancy between them, which is denoted as $\mathcal{X}^{A \rightarrow S}$ and thus $\mathbb{D}^{A \rightarrow S} = (\mathcal{X}^{A \rightarrow S}, \mathcal{Y})$.

Finally, we set up a multi-stage $(\mathcal{X}^A, \mathcal{X}^{A \rightarrow S}, \mathcal{X}^S)$ curriculum learning strategy to learn from easier samples to harder ones during training. Hence, the model is able to deal with slides data directly during inference.

3.2 Transition Data Generation

We approximate the transition between slides data \mathcal{X}^S and article data \mathcal{X}^A , in the form of text generation, where we train a transition model $P_{A \rightarrow S}$ by fine-tuning a BART (Lewis et al., 2020), a pre-trained sequence-to-sequence model.

In order to control the transition process, we design a hyper-parameter schedule weight α in the training objective $\mathcal{L}_{A \rightarrow S}$, which is defined as

$$\mathcal{L}_{A \rightarrow S} = \alpha \mathcal{L}_S + (1 - \alpha) \mathcal{L}_A \quad (1)$$

where \mathcal{L}_S is the cross entropy loss for article-to-slides generation, and \mathcal{L}_A is the cross entropy loss

for article-to-article generation, which is basically input reconstruction. Note that a larger α empirically steers the generated text towards the distribution of slides, and vice versa.

After training the article-to-slides model, we use it for generating intermediate transition data $\mathcal{X}^{A \rightarrow S} = P_{A \rightarrow S}(\mathcal{X}^A, \alpha)$. The generated data is then compiled to form a transition data pair $\mathbb{D}^{A \rightarrow S} = (\mathcal{X}^{A \rightarrow S}, \mathcal{Y})$, reaching the full augmented dataset with quadruples: $\mathbb{D} = (\mathcal{X}^S, \mathcal{X}^A, \mathcal{X}^{A \rightarrow S}, \mathcal{Y})$.

3.3 Curriculum Learning

Finally, we employ curriculum learning strategy to train models on the augmented dataset \mathbb{D} .

Curriculum learning strategies usually include two components: the difficulty measurer and the curriculum scheduler. The difficulty measurer evaluates the difficulty level of the training samples. In this work, we employ a pre-defined approach to difficulty measurement that slides data are harder to learn compared with article data in natural text. Specifically, we treat \mathbb{D}^A , $\mathbb{D}^{A \rightarrow S}$, and \mathbb{D}^S as discrete training buckets that are trained during different curriculum stages. Note that the curriculum is fully extendable, as we can generate more stages of transition data $\mathbb{D}^{A \rightarrow S}$ by tuning the schedule weight α .

During each of the stages, the model takes as input \mathcal{X}^* and is optimized with a simple text summarization objective for summaries \mathcal{Y} . To automatically regulate curriculum stages, we observe the development set ROUGE score during training and move to the next curriculum stage when the score converges. In this way, the model starts with the easiest article data, then switch to the transition data, and finally learns the target slides data. During inference, the model is enhanced with the knowledge to understand noisy data and does not require the help of parallel article data, which may be hard to collect during test time.

4 Experimental Setup

4.1 Datasets

We use two slides summarization datasets for experiments, where one is an existing dataset with few adaptations and the other is collected ourselves. Statistics of these datasets are reported in Table 1, where they are randomly split into training, development, and test set at the ratio of 7:1:2. One major difference between them is the *noisiness* of

Dataset	Split	# Inst.	Src. Len.	Tgt. Len.
S ³	Train	191	697	138
	Dev	26	817	145
	Test	55	663	132
NOISYS ³	Train	306	2,878	298
	Dev	44	2,388	202
	Test	87	2,444	319

Table 1: Statistics of the datasets. “# Inst.” denotes the number of instances in the splits. “Src. Len.” and “Tgt. Len.” denote the *average* token number of parsed slides text (source) and summary text (target) respectively.

the parsed slides text, due to different slides parsing techniques. Empirically, noisier slides text as input makes it harder for models to generate summaries. We describe the details of these datasets as follows.

The Existing Dataset: S³ We use the public SCIDUET-ACL dataset released by Sun et al. (2021), which is a dataset for generating presentation slides from scientific papers.² Based on SCIDUET-ACL, we construct the dataset of slides summarization by keeping the abstract of a paper as the summary for the corresponding slides, referred to as S³ (Scientific Slides Summarization). We keep only the slides text parsed by Sun et al. (2021) and discard the figures and tables in the slides, which are identified by OpenCV (Bradski, 2000).

The Collected Dataset: NOISYS³ We also collect a relatively larger-scale dataset from scratch. We crawl from the ACL anthology and scholars’ home pages for pairs of slides and scientific papers and get 137 and 300 instances respectively. Since most of the slides are in PDF format, we transform all slides into PDF and use *pymupdf* to parse slides text.³ In other words, we do not exclude figures and tables, thus the parsed text is much noisier than one in S³ due to mixed numbers and text snippets.

4.2 Metrics & Baselines

Since the output of our task is a resemblance to the mainstream text summarization tasks, we evaluate the generated summaries with the commonly-used ROUGE-1/2/L (Lin, 2004), which calculates the n-gram overlap of generated and reference text.

There was no previous work for slides summarization, so we adopt several widely used

²The authors of SCIDUET only release the ACL portion of the dataset due to copyright issues.

³<https://github.com/pymupdf/PyMuPDF>

Model	The S ³ Dataset			The NOISYS ³ Dataset		
	R-1	R-2	R-L	R-1	R-2	R-L
Transformer	3.50	0.01	3.42	2.00	0.00	2.00
BERT2BERT	22.87	4.55	13.59	21.44	3.61	12.13
+ LESSON	24.97(+9.18%)	4.55(+0.00%)	13.79(+1.47%)	24.14(+12.59%)	4.07(+12.74%)	12.80(+5.52%)
SciB2SciB	24.89	4.10	12.94	21.22	3.50	11.56
+ LESSON	26.62(+6.95%)	4.97(+21.22%)	12.91(-0.23%)	23.78(+12.06%)	4.25(+21.43%)	12.14(+5.02%)
BART	36.65	9.46	20.25	30.05	6.72	16.68
+ LESSON	38.55(+5.18%)	10.30(+8.88%)	21.19(+4.64%)	34.21(+13.84%)	8.48(+26.19%)	18.27(+9.53%)
BART-ArXiv	37.46	9.64	20.77	32.18	7.46	17.19
+ LESSON	38.92(+3.90%)	10.30(+6.85%)	21.56(+3.80%)	35.88(+11.50%)	9.61(+28.82%)	19.02(+10.65%)

Table 2: Main summarization results of baselines with LESSON in two datasets. LESSON shows consistent performance improvements over baseline models.

abstractive text summarization models as baselines, where we prioritize on the Transformer-based generative language models pre-trained on large corpus.⁴ These models follow the dominant encoder-decoder architecture for sequence-to-sequence (Seq2Seq) generation. Note that LESSON is open to the choices of the backbone summarization model.

Transformer (Vaswani et al., 2017) is the classic attention-based Seq2Seq model, which serves as a non-pre-trained baseline for this task.

BERT2BERT (Rothe et al., 2020) leverages pre-trained checkpoints, such as BERT (Devlin et al., 2019), to initialize both encoder and decoder, where the only variables initialized randomly is the encoder-decoder attention. In addition, we also use SciBERT (Beltagy et al., 2019) as the pre-trained checkpoint for in-domain scientific text.

BART (base version) (Lewis et al., 2020) is a transformer-based PLM for Seq2Seq generation, pre-trained with a set of self-supervised sequence denoising tasks. BART has shown its effectiveness on a variety of text generation tasks, which makes it our primary backbone model for LESSON in the experiments.

BART-ArXiv We also study the effect of *transfer learning* from article summarization to slides summarization, since these two tasks are in a similar domain but the former is much more high-resourced than the latter. To this end, we finetune a BART (base version) with article summarization

⁴Note that we do not incorporate some summarization models such as T5 (Raffel et al., 2020) and PEGASUS (Zhang et al., 2020) as baselines because part of the summary data in S³ is leaked in the C4 dataset upon which these models are pre-trained.

objective on ArXiv dataset (Cohan et al., 2018), which consists of 215,913 scientific papers.⁵

4.3 Implementation Details

We use pre-trained checkpoints provided by HuggingFace (Wolf et al., 2019) in our experiments.

For transition data $\mathcal{X}^{A \rightarrow S}$ generation, we finetune a BART model for 3 epochs on the training set. We set the maximum generation length to 1024, the maximum length for BART. We use Adam (Kingma and Ba, 2015) as the optimizer for all of our models and set the learning rate to 5×10^{-5} .

For scheduling LESSON training, we train our models for 20 epochs in each curriculum stage and select the checkpoint with the highest development set ROUGE score to enter the next stage.

Each model has a different maximum input token length configuration. Since our input data is usually longer than most of the models' maximum possible input length, we set each of them to the max possible length. The max input token lengths for BERT- and BART- models are set to 512, 1024. For summary generation, we set the minimum length to 50 tokens and maximum to 400 tokens. We also set the number of beams to 4 and length penalty to 2.0.

5 Results & Analysis

5.1 Main Results

We report the main results of slides summarization in Table 2. In general, we observe a consistent performance boost with LESSON across different base models, including out-of-the-box pre-trained language models (BERT2BERT, BART), models pre-trained on in-domain texts (SciB2SciB), and

⁵We ensure there is no data leakage in the ArXiv dataset.

Dataset	Ablation	R-1	R-2	R-L
S^3	BART	36.65	9.46	20.25
	+ \mathcal{X}^A	36.84	9.37	19.98
	+ $\mathcal{X}^A, \mathcal{X}^{A \rightarrow S}$	35.88	8.24	19.13
	+ \mathcal{X}^A + CL	37.25	9.65	20.73
	+ LESSON	38.55	10.30	21.19
NOISYS ³	BART	30.05	6.72	16.68
	+ \mathcal{X}^A	32.83	7.25	17.35
	+ $\mathcal{X}^A, \mathcal{X}^{A \rightarrow S}$	30.28	6.65	16.41
	+ \mathcal{X}^A + CL	33.93	7.42	17.53
	+ LESSON	34.21	8.48	18.27

Table 3: Ablation study on different components of LESSON using BART as the base model, including article data \mathcal{X}^A , transition data $\mathcal{X}^{A \rightarrow S}$, and curriculum learning strategy CL. If “+ CL” is not indicated, the data is trained in random order. “+ LESSON” is equivalent to “+ $\mathcal{X}^A, \mathcal{X}^{A \rightarrow S}$ + CL”. A leap in performance emerges when data augmentation and curriculum learning are combined.

pre-trained models finetuned on downstream in-domain summarization task (BART-ArXiv). The ROUGE-2 boosts are the greatest in most cases, up to 28.82% for BART-ArXiv.

When we compare LESSON-BART with BART-ArXiv, we find that a base model with LESSON outperforms the base model with transfer learning on much larger in-domain data. It shows that in cases where abundant in-domain texts are not available, LESSON can still achieve similar performance under the low-resource constraint.

We also observe bigger performance improvements on the NOISYS³ dataset. This corroborates LESSON’s ability to learn noisy and difficult data better through curriculum learning.

Ablation Study We perform an ablation study on different components of LESSON with BART, the results are presented in Table 3. We observe that there is little effect when we use data augmentation without curriculum learning (“+ \mathcal{X}^A ” and “+ $\mathcal{X}^A, \mathcal{X}^{A \rightarrow S}$ ”). However, if we use the curriculum to strategically order the training of the augmented data, we can take full advantage of \mathcal{X}^A and $\mathcal{X}^{A \rightarrow S}$. This provides strong evidence for our assumption that combining data augmentation and curriculum learning leads to better performance on hard-to-learn data.

5.2 Analysis

In this section, we analyze the curriculum data augmentation in LESSON to figure out two research questions: Does the augmented data form a cur-

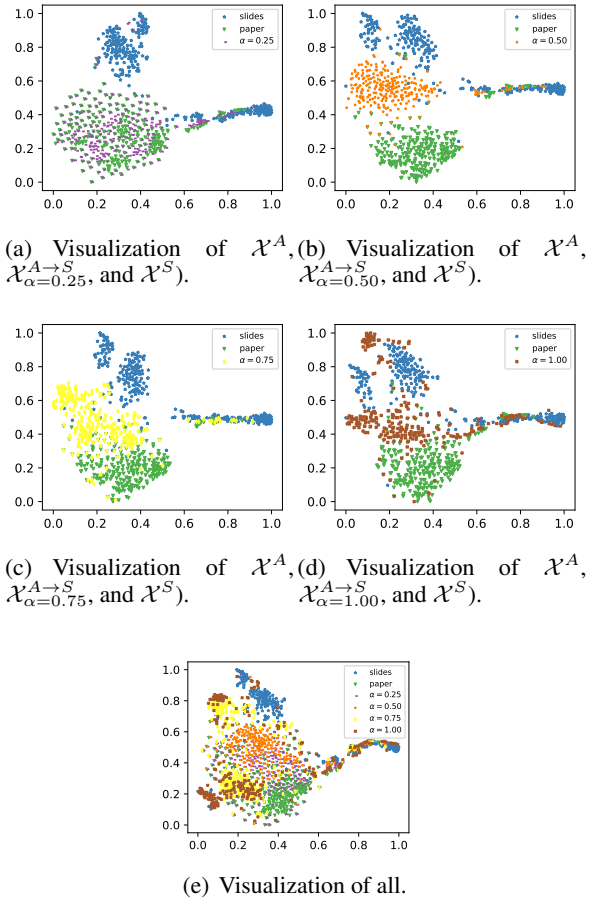


Figure 3: t-SNE visualization of the SciBERT embeddings of articles, transition data and slides text in order to show the distribution shift within curriculum. t-SNE is computed individually for each sub-figure. Therefore, the same data distribution in different sub-figures are slightly different. The transition data is controlled by the schedule weight α . Best viewed in color.

riculum? Does curriculum learning work for slides summarization? ⁶

Data Visualization Theoretically, CL helps the model learn from the easy distribution (articles) to the hard distribution (slides). For an empirical understanding of the distribution shift between them, we visualize the SciBERT embedding (Beltagy et al., 2019) of articles, slides, and in-between transition data in the curriculum as seen in Figure 3. The ideal transition data bridges the distribution gap between the slides data and the article data. Our qualitative evaluation concludes that Figure 3(b) ($\mathcal{X}_{\alpha=0.5}^{A \rightarrow S}$) shows the most desirable case, where the generated transition data lies in the middle of the article and slides data. For other configurations

⁶For the rest of the experiments, we run LESSON with BART (base) by default.

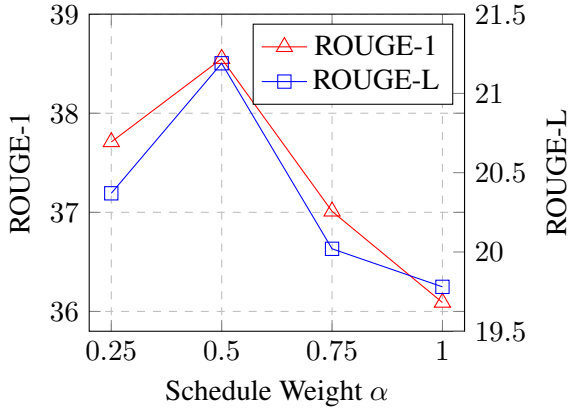


Figure 4: Results of LESSON-BART using transition data generated from different schedule weight α . The result indicates that $\mathcal{X}_{\alpha=0.5}^{A \rightarrow S}$ yields the best performance, which is consistent with our qualitative evaluation that $\mathcal{X}_{\alpha=0.5}^{A \rightarrow S}$ makes up the smoothest transition between the article data and the slides data.

of α , the transition data does not make as good a distributional shift between the article and the slides data.

Then, we verify how $\mathcal{X}^{A \rightarrow S}$ generated with different α affects the curriculum. The results are presented in Figure 4, which we find is consistent with the pattern we find in Figure 3. The transition data $\mathcal{X}_{\alpha=0.5}^{A \rightarrow S}$ achieve the best result because it forms the smoothest distributional shift between the article data \mathcal{X}^A and the slides data \mathcal{Y} . Other transition data, especially $\mathcal{X}_{\alpha=1.0}^{A \rightarrow S}$, performs poorly because the transitions are not smooth, which would lead to a harder optimization process.

Multi-Staged Curriculum We extend the curriculum to multiple stages. In particular, we include $\{\mathcal{X}_{\alpha=i}^{A \rightarrow S}\}_i$ as $\mathcal{X}^{A \rightarrow S}$. The results are shown in Table 4. We organize stage order based on Equation 1. Hence, $\mathcal{X}_{\alpha=0.25}^{A \rightarrow S}$ is closest to the article data and $\mathcal{X}_{\alpha=1.00}^{A \rightarrow S}$ is closest to the slides data.

On the S^3 dataset, the single-stage curriculum learning strategy proves to be the best, and we observe performance decrease when we add more stages to training. This is because, as shown previously in Figure 3, the transition data for $\alpha \in \{0.25, 0.75, 1.00\}$ do not make up a smooth distributional shift between the article data and the target slides data. The learning process would be complicated by these noisy transition data.

Curriculum Schedules To study the effects of curriculum learning, we schedule the curriculum learning in three settings: 1) *regular curricu-*

$\{\mathcal{X}_{\alpha=i}^{A \rightarrow S}\}_i$ Setting	R-1	R-2	R-L
$i = 0.50$ (LESSON baseline)	38.55	10.30	21.19
$i \in \{0.25, 0.50\}$	35.76	6.87	18.01
$i \in \{0.25, 0.50, 0.75\}$	31.73	5.00	14.95
$i \in \{0.25, 0.50, 0.75, 1.00\}$	29.94	3.99	15.06

Table 4: Results of LESSON-BART with multiple curriculum stages for different values of α on S^3 . Increasing curriculum stages leads to performance drop, as the transition data for $\alpha \in \{0.25, 0.75, 1.00\}$ do not make up a smooth distributional shift between the article data and the target slides data.

Dataset	Schedule	R-1	R-2	R-L
S^3	CL	38.55	10.30	21.19
	Anti-CL	36.06	10.03	19.58
	Rand-CL	35.88	8.24	19.13
NOISYS ³	CL	34.21	8.48	18.27
	Anti-CL	32.03	7.16	17.23
	Rand-CL	30.28	6.65	16.41

Table 5: Results of LESSON-BART trained with different curriculum schedules. Anti-CL refers to the inverse training order (first \mathcal{X}^S , then $\mathcal{X}^{A \rightarrow S}$, and last \mathcal{X}^A). Curriculum learning shows better results than Anti-CL and random training schedule.

1) *regular curriculum* (CL) as described in Section 3.3; 2) *anti-curriculum* (Anti-CL) where we reverse the training order of the curriculum; and 3) *random curriculum* (Rand-CL) where we shuffle all data during training. As seen in Table 5, we observe LESSON trained with regular CL achieves the best result whereas LESSON with random CL has the lowest scores, which is consistent with the findings of Wu et al. (2020). The results also demonstrate the difficulty order of slides, transition data, and articles, which is consistent with our intuition and findings mentioned above.

5.3 Case Study

We examine the generation results from a baseline BART model and BART equipped with LESSON on the NOISYS³ test set and compare them with the reference abstract and the input slides. We can identify some key observations in the results, as demonstrated in Table 6. Due to the difficulty of the input slides, the BART baseline model is unable to generate meaningful summarization. On the other hand, as we expected, the improvement of LESSON over the baseline is particularly visible. LESSON makes it possible to digest the noisy input and provide more consistent and factually correct

Reference

Conformity is a type of social influence involving a change in opinion or behavior to fit in with a group. Employing several social networks as the source for our experimental data, we study how the effect of conformity plays a role in changing users' online behavior. We formally define several major types of conformity in individual, peer, and group levels. We propose Confluence model to formalize the effects of social conformity into a probabilistic model. Confluence can distinguish and quantify the effects of the different types of conformities. To scale up to large social networks, we propose a distributed learning method that can construct the Confluence model efficiently with near-linear speedup. Our experimental results on four different types of large social networks, i.e., Flickr, Gowalla, Weibo and Co-Author, verify the existence of the conformity phenomena. Leveraging the conformity information, Confluence can accurately predict actions of users. Our experiments show that Confluence significantly improves the prediction accuracy by up to 5-10% compared with several alternative methods.

BART Baseline

Aspect Term Extraction (ASE) is a novel way to measure the sentiment of an aspect term. In this paper, we propose a new dataset, the Aspect-based Sentiment Estimation (AGE) dataset. Aspect term extraction (aspect-term-extraction) is used to quantify the frequency of distinct aspect terms. We propose a novel approach, which is based on a dataset that is used for the task of estimating the frequency and frequency of the most frequent aspect terms in a domain. Our dataset, which consists of two datasets, consists of three datasets, and consists of four datasets. The first two datasets (the second dataset) are based on the same dataset, and the second dataset on the third dataset (the third dataset). The second dataset, our dataset consists of a dataset consisting of two different datasets, which are both based on our dataset. The third dataset is a dataset based on two datasets. We evaluate our dataset on a variety of domain-specific domains, and we show that our dataset outperforms the previous datasets on a wide range of domains. We show that we outperform the previous dataset on both domains. We also show that aspect term Extraction is more accurate than the first dataset. In our dataset, we find that our approach outperforms our previous dataset by more than 90% of the time. We use a new approach to evaluate the accuracy of our datasets.

LESSON-BART

Aspect-based sentiment analysis estimates the sentiment expressed for each particular aspect (e.g., battery, screen). Different words or phrases, however, may be used to refer to the same aspect, and similar aspects may need to be aggregated at coarser or finer granularity to fit the available space or satisfy user preferences. We introduce the problem of aspect similarity at multiple levels of learning. We decompose it in two processing phases, to allow previous work on term similarity and hierarchical clustering to be reused. We show that the second phase, where aspects are clustered, is almost a solved problem. We also introduce a novel sense pruning mechanism for WordNet-based similarity measures, which improves their performance in the first phase. Finally, we provide publicly available benchmark datasets.

Consistent with the slides. Inconsistent with the slides or flawed generation.

Table 6: Case study for generated slides summaries for different methods on the NOISYS³ test set.

492 summarization. There are, however, several prob- 509
493 lems in LESSON's generation that we observed: 510
494 1) LESSON sometimes generates contents that are 511
495 consistent with the input slides, but inconsistent 512
496 within the passage. The example in Table 6 shows 513
497 LESSON can point out that three factors are influ- 514
498 encing social conformity, but it is unable to name 515
499 them, making the marked sentence abrupt in the 516
500 summary; 2) LESSON sometimes generates fac- 517
501 tually inconsistent information. For example, the 518
502 datasets mentioned in Table 6's LESSON generation 519
503 is incorrect. These errors show that as a generation 520
504 model in its essence, LESSON still faces common 521
505 obstacles in natural language generation. 522

6 Conclusion and Future Work

506
507 In this paper, we tackle the slides summarization 523
508 problem, which is under-studied but of much prac- 524
525

tical use. We formulate this problem as a text sum-
marization task, and propose LESSON with cur-
riculum data augmentation to overcome the *limited*
resource and *noisy input* challenges in this task.
Experiments on both the public dataset S³ and our
collected dataset NOISYS³ show that LESSON con-
sistently improves summarization results over base-
line models. Further analyses show the data aug-
mentation process successfully creates transition
data that bridges the gap between the article data
and the slides data. The transition data enables cur-
riculum training, which proves to boost the model's
ability to learn from the noisy slides data. In the
future, we will emphasize on multi-modal slides
summarization to utilize the layout information of
the slides, and explore few-shot adaptation to slides
of unseen domains.

References

- 527 Ateret Anaby-Tavor, Boaz Carmeli, Esther Goldbraich,
528 Amir Kantor, George Kour, Segev Shlomov, Naama
529 Tepper, and Naama Zwerdling. 2020. [Do not have
530 enough data? deep learning to the rescue!](#) In *The
531 Thirty-Fourth AAAI Conference on Artificial Intelli-
532 gence, AAAI 2020, The Thirty-Second Innovative Ap-
533 plications of Artificial Intelligence Conference, IAAI
534 2020, The Tenth AAAI Symposium on Educational
535 Advances in Artificial Intelligence, EAAI 2020, New
536 York, NY, USA, February 7-12, 2020*, pages 7383–
537 7390. AAAI Press.
- 538 Iz Beltagy, Kyle Lo, and Arman Cohan. 2019. [SciB-
539 ERT: A pretrained language model for scientific text.](#)
540 In *Proceedings of the 2019 Conference on Empirical
541 Methods in Natural Language Processing and the
542 9th International Joint Conference on Natural Lan-
543 guage Processing (EMNLP-IJCNLP)*, pages 3615–
544 3620, Hong Kong, China. Association for Computa-
545 tional Linguistics.
- 546 Iz Beltagy, Matthew E Peters, and Arman Cohan.
547 2020. Longformer: The long-document transformer.
548 *arXiv preprint arXiv:2004.05150*.
- 549 Yoshua Bengio, Jérôme Louradour, Ronan Collobert,
550 and Jason Weston. 2009. [Curriculum learning.](#) In
551 *Proceedings of the 26th Annual International Con-
552 ference on Machine Learning, ICML 2009, Mon-
553 treal, Quebec, Canada, June 14-18, 2009*, volume
554 382 of *ACM International Conference Proceeding
555 Series*, pages 41–48. ACM.
- 556 G. Bradski. 2000. The OpenCV Library. *Dr. Dobb's
557 Journal of Software Tools*.
- 558 Jiaao Chen, Zichao Yang, and Diyi Yang. 2020. [Mix-
559 Text: Linguistically-informed interpolation of hid-
560 den space for semi-supervised text classification.](#) In
561 *Proceedings of the 58th Annual Meeting of the Asso-
562 ciation for Computational Linguistics*, pages 2147–
563 2157, Online. Association for Computational Lin-
564 guistics.
- 565 Yong Cheng, Lu Jiang, Wolfgang Macherey, and Ja-
566 cob Eisenstein. 2020. [AdvAug: Robust adversarial
567 augmentation for neural machine translation.](#) In
568 *Proceedings of the 58th Annual Meeting of the Asso-
569 ciation for Computational Linguistics*, pages 5961–
570 5970, Online. Association for Computational Lin-
571 guistics.
- 572 Arman Cohan, Franck Dernoncourt, Doo Soon Kim,
573 Trung Bui, Seokhwan Kim, Walter Chang, and Na-
574 zli Goharian. 2018. [A discourse-aware attention
575 model for abstractive summarization of long docu-
576 ments.](#) In *Proceedings of the 2018 Conference of
577 the North American Chapter of the Association for
578 Computational Linguistics: Human Language Tech-
579 nologies, Volume 2 (Short Papers)*, pages 615–621,
580 New Orleans, Louisiana. Association for Computa-
581 tional Linguistics.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and
582 Kristina Toutanova. 2019. [BERT: Pre-training of
583 deep bidirectional transformers for language under-
584 standing.](#) In *Proceedings of the 2019 Conference
585 of the North American Chapter of the Association
586 for Computational Linguistics: Human Language
587 Technologies, Volume 1 (Long and Short Papers)*,
588 pages 4171–4186, Minneapolis, Minnesota. Associ-
589 ation for Computational Linguistics. 590
- Alexander Fabbri, Simeng Han, Haoyuan Li, Haoran
591 Li, Marjan Ghazvininejad, Shafiq Joty, Dragomir
592 Radev, and Yashar Mehdad. 2021. [Improving zero
593 and few-shot abstractive summarization with inter-
594 mediate fine-tuning and data augmentation.](#) In *Pro-
595 ceedings of the 2021 Conference of the North Amer-
596 ican Chapter of the Association for Computational
597 Linguistics: Human Language Technologies*, pages
598 704–717, Online. Association for Computational
599 Linguistics. 600
- Marzieh Fadaee, Arianna Bisazza, and Christof Monz.
601 2017. [Data augmentation for low-resource neural
602 machine translation.](#) In *Proceedings of the 55th An-
603 nual Meeting of the Association for Computational
604 Linguistics (Volume 2: Short Papers)*, pages 567–
605 573, Vancouver, Canada. Association for Computa-
606 tional Linguistics. 607
- Steven Y. Feng, Varun Gangal, Jason Wei, Sarath Chan-
608 dar, Soroush Vosoughi, T. Mitamura, and E. Hovy.
609 2021. A survey of data augmentation approaches
610 for nlp. *ArXiv*, abs/2105.03075. 611
- Sebastian Gehrmann, Yuntian Deng, and Alexander
612 Rush. 2018. [Bottom-up abstractive summarization.](#)
613 In *Proceedings of the 2018 Conference on Empir-
614 ical Methods in Natural Language Processing*,
615 pages 4098–4109, Brussels, Belgium. Association
616 for Computational Linguistics. 617
- Karl Moritz Hermann, Tomas Kocisky, Edward Grefen-
618 stette, Lasse Espeholt, Will Kay, Mustafa Suleyman,
619 and Phil Blunsom. 2015. [Teaching machines to
620 read and comprehend.](#) In *Advances in Neural Infor-
621 mation Processing Systems 28: Annual Conference
622 on Neural Information Processing Systems 2015,
623 December 7-12, 2015, Montreal, Quebec, Canada*,
624 pages 1693–1701. 625
- Alex Hernandez-Garcıa and Peter Konig. 2020. [Data
626 augmentation instead of explicit regularization.](#) 627
- Diederik P. Kingma and Jimmy Ba. 2015. [Adam: A
628 method for stochastic optimization.](#) In *3rd Inter-
629 national Conference on Learning Representations,
630 ICLR 2015, San Diego, CA, USA, May 7-9, 2015,
631 Conference Track Proceedings*. 632
- Ashutosh Kumar, Satwik Bhattamishra, Manik Bhan-
633 dari, and Partha Talukdar. 2019. [Submodular
634 optimization-based diverse paraphrasing and its ef-
635 fectiveness in data augmentation.](#) In *Proceedings of
636 the 2019 Conference of the North American Chap-
637 ter of the Association for Computational Linguistics:*
638

754	Xin Wang, Yudong Chen, and Wenwu Zhu. 2021. A survey on curriculum learning. <i>IEEE Transactions on Pattern Analysis and Machine Intelligence</i> .	812
755		813
756		814
757	Jason Wei, Chengyu Huang, Soroush Vosoughi, Yu Cheng, and Shiqi Xu. 2021. Few-shot text classification with triplet networks, data augmentation, and curriculum learning. In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 5493–5500, Online. Association for Computational Linguistics.	815
758		816
759		817
760		818
761		819
762		820
763		821
764		822
765	Jason Wei and Kai Zou. 2019. EDA: Easy data augmentation techniques for boosting performance on text classification tasks. In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 6382–6388, Hong Kong, China. Association for Computational Linguistics.	823
766		824
767		825
768		826
769		827
770		828
771		829
772		
773	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2019. Huggingface’s transformers: State-of-the-art natural language processing. <i>ArXiv</i> , abs/1910.03771.	830
774		831
775		832
776		833
777		834
778		835
779		
780		
781		
782	Xiaoxia Wu, Ethan Dyer, and Behnam Neyshabur. 2020. When do curricula work? <i>arXiv preprint arXiv:2012.03107</i> .	836
783		837
784		838
785		839
786	Benfeng Xu, Licheng Zhang, Zhendong Mao, Quan Wang, Hongtao Xie, and Yongdong Zhang. 2020a. Curriculum learning for natural language understanding. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 6095–6104, Online. Association for Computational Linguistics.	840
787		841
788		842
789		843
790		844
791		845
792		846
793		847
794		848
795		
796	Yang Xu, Yiheng Xu, Tengchao Lv, Lei Cui, Furu Wei, Guoxin Wang, Yijuan Lu, Dinei Florencio, Cha Zhang, Wanxiang Che, Min Zhang, and Li-dong Zhou. 2021. LayoutLMv2: Multi-modal pre-training for visually-rich document understanding. In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 2579–2591, Online. Association for Computational Linguistics.	
797		
798		
799		
800		
801		
802		
803	Yiheng Xu, Minghao Li, Lei Cui, Shaohan Huang, Furu Wei, and Ming Zhou. 2020b. Layoutlm: Pre-training of text and layout for document image understanding. In <i>KDD '20: The 26th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, Virtual Event, CA, USA, August 23-27, 2020</i> , pages 1192–1200. ACM.	
804		
805		
806		
807		
808		
809		
810	Yiben Yang, Chaitanya Malaviya, Jared Fernandez, Swabha Swayamdipta, Ronan Le Bras, Ji-Ping	
811		
	Wang, Chandra Bhagavatula, Yejin Choi, and Doug Downey. 2020. Generative data augmentation for commonsense reasoning. In <i>Findings of the Association for Computational Linguistics: EMNLP 2020</i> , pages 1008–1025, Online. Association for Computational Linguistics.	
	Yongjian You, Weijia Jia, Tianyi Liu, and Wenmian Yang. 2019. Improving abstractive document summarization with salient information modeling. In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 2132–2141, Florence, Italy. Association for Computational Linguistics.	
	Jingqing Zhang, Yao Zhao, Mohammad Saleh, and Peter Liu. 2020. Pegasus: Pre-training with extracted gap-sentences for abstractive summarization. In <i>International Conference on Machine Learning</i> , pages 11328–11339. PMLR.	
	Xingxing Zhang, Mirella Lapata, Furu Wei, and Ming Zhou. 2018. Neural latent extractive document summarization. In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 779–784, Brussels, Belgium. Association for Computational Linguistics.	
	Zhou Zhao, Haojie Pan, Changjie Fan, Yan Liu, Linlin Li, and Min Yang. 2019. Abstractive meeting summarization via hierarchical adaptive segmental network learning. In <i>The World Wide Web Conference, WWW 2019, San Francisco, CA, USA, May 13-17, 2019</i> , pages 3455–3461. ACM.	
	Chenguang Zhu, Ruochen Xu, Michael Zeng, and Xuedong Huang. 2020. A hierarchical network for abstractive meeting summarization with cross-domain pretraining. In <i>Findings of the Association for Computational Linguistics: EMNLP 2020</i> , pages 194–203, Online. Association for Computational Linguistics.	