# Fusing Multimodal Signals on Hyper-complex Space for Extreme Abstractive Text Summarization (TL;DR) of Scientific Contents

Yash Kumar Atri
yashk@iiitd.ac.in
IIIT Delhi

Vikram Goyal
vikram@iiitd.ac.in
IIIT Delhi

Tanmoy Chakraborty
tanchak@iitd.ac.in
IIT Delhi

## ABSTRACT

The realm of scientific text summarization has experienced remarkable progress due to the availability of annotated brief summaries and ample data. However, the utilization of multiple input modalities, such as videos and audio, has yet to be thoroughly explored. At present, scientific multimodal-input-based text summarization systems tend to employ longer target summaries like abstracts, leading to an underwhelming performance in the task of text summarization.

In this paper, we deal with a novel task of *extreme abstractive text summarization* (*aka TL;DR generation*) *by leveraging multiple input modalities*. To this end, we introduce mTLDR, a first-of-its-kind dataset for the aforementioned task, comprising videos, audio, and text, along with both author-composed summaries and expert-annotated summaries. The mTLDR dataset accompanies a total of 4, 182 instances collected from various academic conference proceedings, such as ICLR, ACL, and CVPR. Subsequently, we present mTLDRgen, an encoder-decoder-based model that employs a novel dual-fused hyper-complex Transformer combined with a Wasserstein Riemannian Encoder Transformer, to dexterously capture the intricacies between different modalities in a hyper-complex latent geometric space. The hyper-complex Transformer captures the intrinsic properties between the modalities, while the Wasserstein Riemannian Encoder Transformer captures the latent structure of the modalities in the latent space geometry, thereby enabling the model to produce diverse sentences. mTLDRgen outperforms 20 baselines on mTLDR as well as another non-scientific dataset (How2) across three Rouge-based evaluation measures. Furthermore, based on the qualitative metrics, BERTScore and FEQA, and human evaluations, we demonstrate that the summaries generated by mTLDRgen are fluent and congruent to the original source material.

## CCS CONCEPTS

• >**Computing methodologies → Natural language processing**.

## KEYWORDS

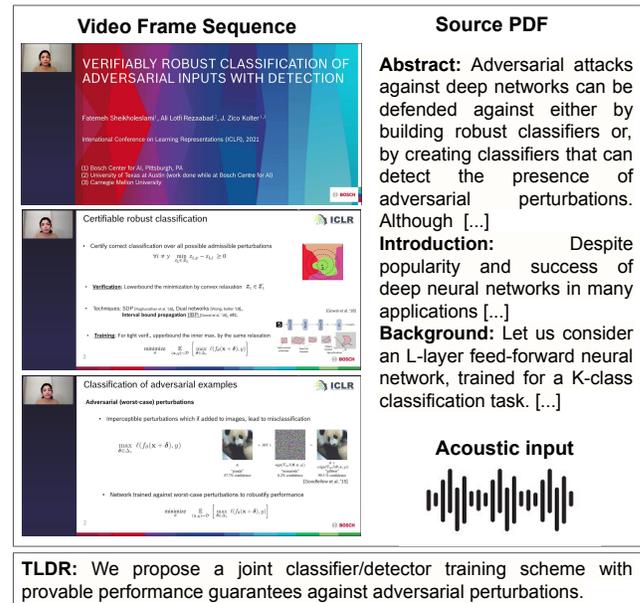Abstractive summarization, multi-modal summarization, neural networks

**Video Frame Sequence**

VERIFIABLY ROBUST CLASSIFICATION OF ADVERSARIAL INPUTS WITH DETECTION

**Source PDF**

**Abstract:** Adversarial attacks against deep networks can be defended against either by building robust classifiers or, by creating classifiers that can detect the presence of adversarial perturbations. Although [...]
**Introduction:** Despite popularity and success of deep neural networks in many applications [...]
**Background:** Let us consider an L-layer feed-forward neural network, trained for a K-class classification task. [...]

**Acoustic input**

**TLDR:** We propose a joint classifier/detector training scheme with provable performance guarantees against adversarial perturbations.

Figure 1: **A sample of mTLDR dataset with video, text and audio modalities along with the target TLDR. The feature representations for video frames are obtained by ResNext, audio features are extracted using Kaldi, and the text is extracted from the pdf of the article.**

## 1 INTRODUCTION

Abstractive text summarization enables one to promptly comprehend the essence of a written work, determining if it is worth perusing. In contrast to extractive summarization, which emphasizes the crucial passages within the original document as a summary, abstractive summarization recomposes the summary from scratch by synthesizing the core semantics and the entire substance of the document. Earlier studies dealt with abstractive summarization by solely utilizing textual input [9, 10, 15, 26, 46]; thereafter, multimodal inputs [29, 42, 47, 64] were integrated to enhance the quality of the generated summaries. Studies have revealed that multimodal data assists humans in comprehending the essence of a written work more effectively [17], thus leading us to the inference

that multimodal data can enrich the context and produce more comprehensive scientific summaries.

**Motivation:** With the emergence of deep learning architectures like LSTM, Attention, and Transformer, the literature in the scientific community has skyrocketed. It is extremely hard to keep up with the current literature by going through every piece of text in a research article. The abstract of a paper often serves as a bird's eye view of the paper, highlighting the problem statement, datasets, proposed methodology, analysis, etc. Recent studies [1] re-purpose abstracts to generate summaries of scientific articles. However, it is cumbersome to go through the abstract of each paper. The abstracts are nearly 300 tokens long, and reading the complete abstract of every paper to figure out the mutual alignment is tedious. The task of TL;DR (*aka*, tl;dr, too long; didn't read) [5, 55] was introduced to generate an extremely concise summary from the text-only article highlighting just the high-level contributions of the work. Later, Mao et al. [37] introduced the CiteSum dataset for generating text-only extreme summaries. However, the text alone can not comprehend the entire gist of the research article. The multimodal information, including the video of the presentation and audio, often provide crucial signals for extreme text summary generation.

**Problem statement:** In this work, we propose a new task of multimodal-input-based TL;DR generation for scientific contents which aims to generate an extremely-concise and informative text summary. We incorporate the visual modality to capture the visual elements, the audio modality to capture the tonal-specific details of the presenter, and the text modality to help the model align all three modalities. We also show the generalizability of the proposed model on another non-academic dataset (How2).

**State-of-the-art and limitations:** The pursuit of multimodal-input-based abstractive text summarization can be related to various other fields, such as image and video captioning [22, 34, 39, 48, 49], video story generation [16], video title generation [57], and multimodal sentence summarization [28]. However, these works generally produce summaries based on either images or short videos, and the target summaries are easier to predict due to the limited vocabulary diversity. On the other hand, scientific documents have a complex and structured vocabulary, which the existing methods [42] of generating short summaries are not equipped to handle. Recently, Atri et al. [1] proposed as a novel dataset for the multimodal text summarization of scientific presentations; however, it uses the abstract as the target summary, which falls short in producing coherent summaries for the extreme multimodal summarization (TL;DR) task.

In summary, the current paper offers the following contributions:

- **Novel problem:** We propose the task of extreme abstractive text summarization for scientific contents, by utilizing videos, audio and research articles as inputs.
- **Novel dataset:** The development and curation of the first large-scale dataset mTLDR for extreme multimodal-input-based text summarization of scientific contents. Figure 1 shows an excerpt from the mTLDR dataset. This dataset has been meticulously compiled from five distinct public websites and comprises articles and videos obtained from renowned international conferences in Computer Science. The target summaries are a fusion of

manually-annotated summaries and summaries written by the authors/presenters of the papers.
- **Novel model:** We propose mTLDRgen, a novel encoder-decoder-based model designed to effectively capture the dynamic interplay between various modalities. The model is implemented with a dual-fused hyper-complex Transformer and a Wasserstein Riemannian Encoder Transformer. The hyper-complex Transformer projects the modalities into a four-dimensional space consisting of one real component and three imaginary components, thereby capturing the intrinsic properties of individual modalities and their relationships with one another. Additionally, the Wasserstein Riemannian Encoder Transformer is employed to apprehend the latent structure of the modalities in the geometry of the latent space.
- **Evaluation:** We benchmark mTLDR over six extractive (text-only), eight abstractive (text-only), two video-based and four multimodal summarization baselines, demonstrating the effectiveness of incorporating multimodal signals in providing more context and generating more fluent and informative summaries. We evaluate the benchmark results over the quantitative (Rouge-1/2/L) and qualitative (BERTScore and FEQA) metrics. Our proposed modal, mTLDRgen, beats the best-performing baseline by +5.24 Rouge-1 and +3.35 Rouge-L points. We also show the generalizability of mTLDRgen on another non-scientific dataset (How2).
- **Deployment:** We further designed an in-house institute-wide web API based on the end-to-end pipeline of mTLDRgen. The web API is currently undergoing a beta testing phase and has gathered more than 100+ hits so far. The API will be made open across academic institutes and beyond upon successful completion of the beta testing.

**Reproducibility:** We discuss the detailed hyperparameters (Supplementary, Table 7) and experimentation setting in Section 6.1. We also provide a sample dataset of mTLDR and the source code of mTLDRgen at https://github.com/LCS2-IIITD/mTLDRgen.

## 2 RELATED WORK

The development and utilization of abstractive text summarization systems involve the formulation of textual summaries through the integration of two or more auxiliary signals. These signals, beyond the traditional text, may encompass video [12], images [51], and audio [19]. The integration of additional modalities, as compared to text-only systems, offers a plethora of opportunities to enhance the contextual richness and knowledge base of the generated summaries. Several recent studies [1, 45] demonstrated that the integration of multimodal signals such as video and audio can significantly improve the contextual accuracy and informativeness of the summaries generated by unimodal systems.

**Unimodal text summarization:** Text summarization is classified into two categories – extractive and abstractive. Extractive systems extract the most relevant sentences from the source document to form the summary, while abstractive systems paraphrase important sentences to generate a new summary. Conventional extractive summarization approaches either construct a graph representation of the source document [13, 32, 38] or pose the summarization task as a binary classification with ranking [8, 35, 40, 62]. On the other hand, abstractive summarization has significantly benefited from

the advent of deep neural networks. Early works [41, 46] utilized CNN/Dailymail dataset [21] to explore abstractive summarization on a large scale. Later, Pointer Generators (PG) [46] were extended to capture the latent structures of documents [14, 50]. The use of Transformers [54] and attention mechanisms [2] further improved the encoding of long sequential data. These improvements include leveraging Transformers [20] and repurposing attention heads as copy pointers [15] to enhance the quantitative performance. Large language models [11, 44, 59] have demonstrated impressive performance on multiple datasets [14, 21]. Models proposed in [11] and [59] are pre-trained using token and phrase masking techniques, respectively, while Raffel et al. [44] approached all downstream tasks as a text-to-text problem and pre-train using a single loss function.

**Extreme unimodal text summarization:** The objective of extreme text summarization is to drastically reduce the size of the source document while preserving its essence in the resulting summary. The concept of extreme summarization was first introduced by Völske et al. [55] with a novel dataset focused on social media summarization. Subsequently, Cachola et al. [5] and Mao et al. [37] presented new corpora, namely SciTLDR and CiteSum, respectively, for extreme summarization of scientific documents. However, it remains an open area to explore extreme abstractive text summarization using multimodal signals.

**Text summarization with multimodality:** The incorporation of multimodal information in text summarization enriches the comprehension of the data, elevating it to a representation that better reflects the source document [3, 25]. In the absence of multimodal information, the summarization model can only comprehend limited information; however, with the integration of multiple modalities, the models acquire a more comprehensive understanding, leading to the creation of highly fluent and semantically meaningful summaries. Multimodal summarization has been explored in various domains, including instructional YouTube videos [45], news videos [7, 30], and recipe videos [63]. The concept of generating multimodal outputs from multimodal inputs has also been studied [43, 60, 64], where a news event was summarized with a text summary along with a corresponding image. Zhang et al. [61] recently introduced cross-modal alignment to harmonize visual features with the text to produce a more coherent summary. Although existing datasets feature the use of either images or videos as a visual modality, the closest dataset to our task is the How2 dataset [45], housing all three modalities to generate a short summary. However, when compared to mTLDR, the How2 dataset falls short in terms of length, structuredness and complexity of vocabulary in the source and target documents. Evidently, the existing approaches over the How2 dataset fail to extend similar performance on the mTLDR dataset.

## 3   PROPOSED DATASET

To explore the efficacy of multimodal signals and enable enriched abstractive summaries aided by various modalities, we introduce mTLDR, the first large-scale **m**ultimodal-input based abstractive summarization (**TL;DR**) dataset with diverse lengths of videos. mTLDR is collected from various well-known academic conferences like ACL, ICLR, CVPR, etc. The only comparable dataset to mTLDR is

**Table 1: Statistics of the used datasets (mTLDR and How2) – the number of samples (#source), average token length of source documents (avg source len), average tokens in the target summaries (avg target len), and abstractness percentage (Abs) of datasets.**

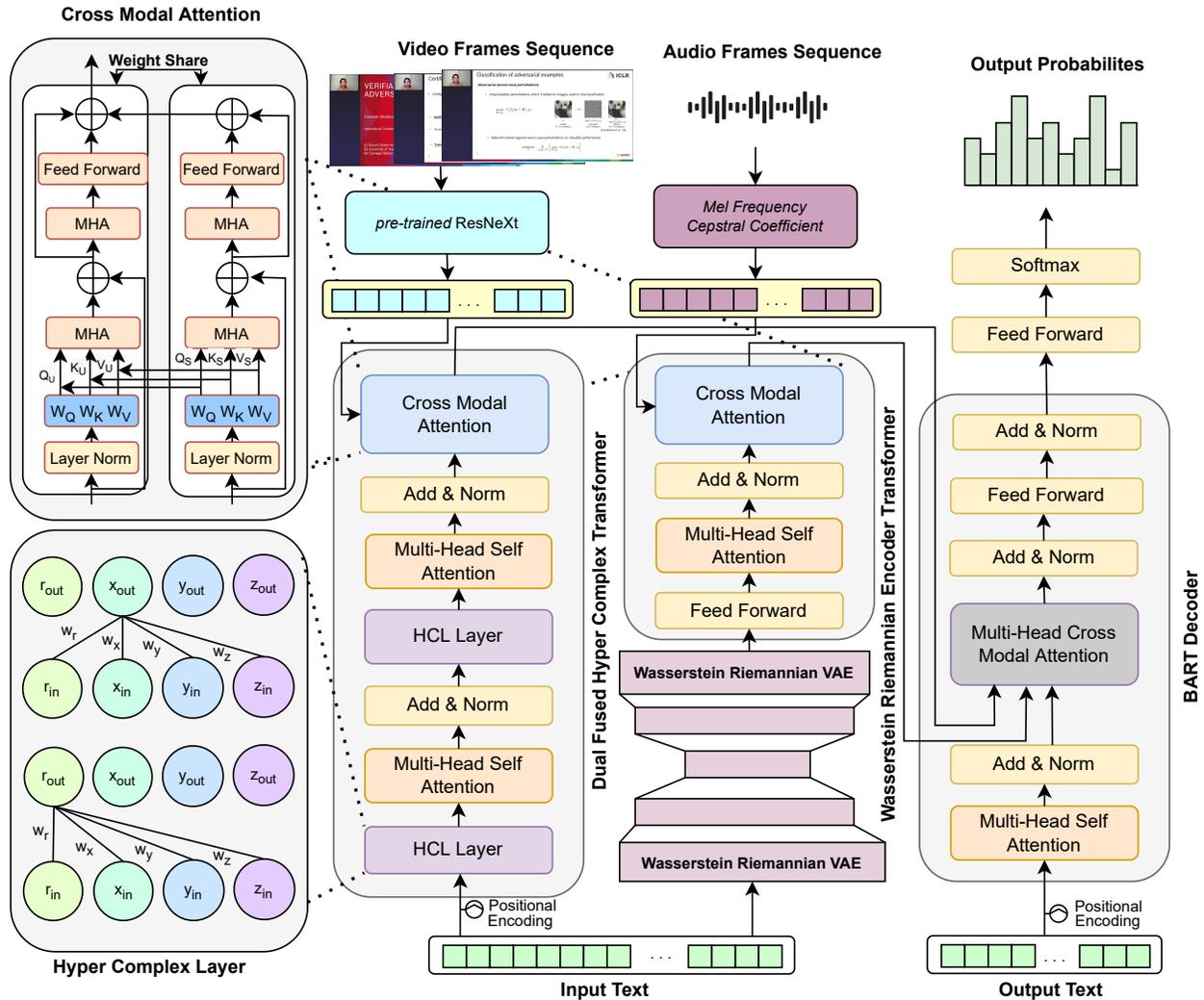| Dataset | #source | avg source len | avg target len | %Abs |
|---------|---------|----------------|----------------|------|
| How2    | 73993   | 291            | 33             | 14.2 |
| mTLDR   | 4182    | 5K             | 18             | 15.9 |

the How2 dataset, which comprises short instructional videos from various topics like gardening, yoga, sports, etc. Compared to How2, mTLDR contains structured and complex vocabulary, which requires attention to diverse information while generating summaries.

Our compilation encompasses video recordings from openreview.net and videolecture.net, in addition to the accompanying source pdf and metadata information, including the details of the authors, title, and keywords. The collected dataset comprises a total of 4, 182 video recordings, spanning a duration of over 1, 300 hours. Of these, we designated 2, 927 instances as the training set, 418 for validation, and 837 for testing. The average length of the videos is 14 minutes, and the TLDR summary has an average of 19 tokens. The target summaries for the data are a combination of human-annotated and author-generated summaries. In terms of abstractness, mTLDR contains 15.9% novel $n$-grams. Each data instance includes a video, audio extracted from the video, an article pdf, and a target summary. We opted not to annotate or retain multiple summaries for a single instance to ensure efficient training and testing processes. We assert that a single extreme summary is sufficient to convey the essence of the paper. The target summaries for papers obtained from the ACL anthology were annotated as they lacked any author-generated summaries. Of the 4, 182 videos, a total of 1, 128 summaries were manually annotated by 25 annotators. During the annotation process, the annotators were instructed to thoroughly read the abstract, introduction, and conclusion and to have a general understanding of the remaining content. Each summary was then verified by another to confirm that it accurately represents the paper's major contributions.

In contrast, the How2 dataset [45] consists of 73, 993 training, 2, 965 validation, and 2, 156 test instances. The average token length for the source documents is 291, while for the target summary, it is 33. Compared to the source document, the target summaries contains 14.2% novel $n$-grams. The transcripts for videos and the target summary are human-annotated. Table 1 shows brief statistics of the How2 and mTLDR datasets.

## 4   PROPOSED METHODOLOGY

This section presents our proposed system, **mTLDRgen**, a **m**ultimodal-input-based extreme abstractive text summary (**TL;DR**) **gen**erator. Figure 2 shows a schematic diagram. During an academic conference presentation, there are typically three major modalities present – visual, audio, and text, each of which complements the others, and when combined, contributes to a rich and expressive feature space, leading to the generation of coherent and fluent summaries. mTLDRgen initially extracts features from the independent modalities and then feeds them to the dual-fused hyper-complex

**Figure 2: An overview of the proposed model – mTLDRgen. It houses two parallel encoders, one with a hyper-complex layer fused with the video embeddings using cross-model attention and the other with Wasserstein Riemannian Encoder Transformer with audio embeddings fused with cross-model attention. The individual encoder representations are later fused with the multi-head attention of the pre-trained BART decoder to generate the final summary.**

Transformer (DFHC) and the Wasserstein Riemannian Encoder Transformer (WRET) blocks. Cross-modal attention is used to fuse the visual and audio features with the text representations. Finally, the fused representation is fed to a pre-trained BART [27] decoder block to produce the final summary. The rest of this section delves into the individual components of mTLDRgen.

### 4.1 Video Feature Extraction

The video modality in an academic presentation often comprises variations in frames and kinesthetic signals, highlighting key phrases or concepts during a presentation. To capture visual and kinesthetic aspects, we utilise the ResNeXt-152-3D [24] model as it is pre-trained on the Kinetics dataset for recognition of 400 human actions. Four frames per second are extracted from the video, cropped

to $112 \times 112$ pixels and normalized, and a 2048-dimensional feature vector is extracted from the ResNeXt-152-3D model for every 12 non-overlapping frames. The 2048-dimensional vector is then fed to the mean pooling layer to obtain a global representation of the video modality. Later, a feed-forward layer is applied to map the 2048-dimensional vector to a 512-dimensional vector.

### 4.2 Speech Feature Extraction

To capture the variations in the speaker's voice amplitudes, which are considered to signify the importance of specific topics or phrases [52], we extract audio features from the conference video. This is accomplished by extracting audio from the video using the FFM-PEG package[1], resampling it to a mono channel, processing it to

---

[1]https://ffmpeg.org/

a 16$K$ Hz audio sample, and dividing it into overlapping windows of 30 milliseconds. The extracted audio is then processed to obtain 512-dimensional Mel Frequency Cepstral Coefficients (MFCC) features. The final representation is obtained by applying a log Mel frequency filter bank and discrete cosine transformation, and the feature sequence is padded or clipped to a fixed length.

## 4.3 Textual Feature Extraction

In order to extract the feature representations for the article text, the pdf content is obtained through the Semantic Scholar Open Research pipeline (SSORP) [36]. SSORP uses the SCIENCEPARSE[2] and GROBID[3] APIs for text extraction from pdf. For the How2 dataset, the video transcriptions are manually annotated and transformed into a text feature set for training. In contrast, the acoustic features for mTLDRgen are not transformed into the text as they are characterized by a variety of non-native English accents and a high error rate for speech-to-text models. Both the textual representations are tokenized using the vanilla BART tokenizer and transformed to word vectors using standard Transformer positional encoding embeddings.

## 4.4 Dual-fused Hyper-complex Transformer

We propose a dual-fused hyper-complex Transformer (DFHC) for the task of multimodal text summarization. Compared to multi-head attention, the hyper-complex layer allows mTLDRgen to efficiently capture the intricacies between different modalities and learn better representations [58] in the hyper-complex space. For the block DFHC, we represent the text as $X$ and pass it through the hyper-complex layer to extract the (Q) Query, (K) Key and (V) Value transformations as follows: $Q, K, V = \Phi(\text{HCL}(X))$, where $HCL(X) = Hx + b$. Here $H \in \mathbb{R}^{m \times n}$ is constructed by a sum of Kronecker products and is given by $H = \sum_{i=1}^{n} P_i \otimes Q_i$. The $P_i$ and $Q_i$ are the parameter matrices, and $\otimes$ represents the Kronecker product.

The final attention score $A$ is obtained as:

$$A = \text{softmax}(\frac{QK^T}{\sqrt{d_k}})$$

Here $Q$ represents the query value, $K$ represents the key value, and $d_k$ represents the key dimension.

The HCL layers share attention weights among the multi-head attention heads. The multi-head attention weights are concatenated and represented as

$$X = \text{HCL}([H_1 + ... + H_{Num_h}])$$

Here $Num_h$ represents the attention head. The final output obtained from the HCL layer is represented as:

$$Y = \text{HCL}(\text{ReLU}(\text{HCL}(X))),$$

The transformation $Y$ is passed through a multi-head attention block and is fused with the visual embeddings using the cross-model attention as discussed in Section 4.6.

[2]https://github.com/allenai/science-parse
[3]https://github.com/kermitt2/grobid

## 4.5 Wasserstein Riemannian Encoder Transformer

We base our idea from Wang and Wang [56] to repurpose the Wasserstein Riemannian Autoencoder to Wasserstein Riemannian Encoder Transformer (WRET) in the summarization setting.

For an input $X$ and a manifold $M$, the Riemannian manifold is represented as $(M, G)$, where $G$ represents the Riemannian tensor unit. For two vectors $u$ and $v$ in the tangent space $T_z M$, the inner product is computed using $\langle u, v \rangle_G = u^T G(z)v$. The Wasserstein block acts as a Variational Autoencoder. However, we extract the feature dimension from the last layer and feed it to the attention block.

The Wasserstein Autoencoder optimizes the cost between the target data distribution $A_x(x)$ and the predicted data distribution $B_x(x)$ using:

$$Dist(A_X, B_G) = \inf_{Q(Z|X) \in Q} E_{P_X} E_{Q(Z|X)} [c(X, G(Z))] \\ + \lambda MMD(Q_Z, P_Z)$$

where $G$ is the generator function, $\lambda$ is a learnable metric, $c$ is the optimal cost, and $D_z$ is approximated between $B_G$ and $Q_Z(z) = \int q(z|x)p(x)dx$ using the Maximum Mean Discrepancy (MMD) [53]. The MMD is computed using

$$MMD_k(P_Z, Q_Z) = || \int_{\mathcal{Z}} k(z, \cdot) dP_Z - \int_{\mathcal{Z}} k(z, \cdot) dQ_Z ||$$

We formulate a RNF function $F = f_K \ldots f_1$, and optimize the following RNF-Wasserstein function,

$$Dist(A_X, B_G) = \inf_{Q(Z|X) \in Q} P_X Q(Z|X) [c(X, G(Z'))] \\ + \lambda MMD(Q_{Z'}, P_{Z'}) \\ + \alpha(KLD(q(z|x)||p(z)) - \sum \log |det \frac{\partial f'}{\partial z}|)$$

where $Z' = F(Z)$, KLD is KL [23] divergence, and $p(z)$ represents the posterior probability distribution. The MMD term is approximated using the Gaussian kernel $k(z, z') = e^{-||z-z'||^2}$. The term $G(Z')$ represents the reconstructed feature set, which is then passed to a feed-forward layer. The attention weights are computed for $G(Z')$ and fused with the audio feature using cross-modal attention as discussed in Section 4.6.

## 4.6 Cross-model Attention

We fuse the text-video and text-audio features using cross-modal attention to align the attention distribution obtained from the last layer. The text feature set projects the Query ($Q$) value, while the video and audio features project the key ($K$) and value ($V$), respectively. The obtained $Q, K$, and $V$ representations are passed through cross-modal attention, and the final encoder representation $E_s$ is obtained.

$$Q = Z_t W_q; \ K = Z_v W_k; \ V = Z_v W_v$$

$$E_s = \text{softmax} \left( \frac{X_\alpha W_{Q_\alpha} W_{K_\beta}^\top X_\beta^\top}{\sqrt{d_k}} \right) X_\beta W_{V_\beta}$$

The amalgamated representation, represented as $E_s$, is subsequently integrated with the multi-head attention mechanism of the BART decoder to yield the final summary.

To the best of our knowledge, the application of hyper-complex Transformer and Wasserstein Riemannian flow for abstractive text summarization has not been explored yet.

## 5 BASELINES

We benchmark our proposed model against 20 baselines – six extractive text summarization, eight abstractive text summarization, two video-based abstractive text summarization, and four multimodal-input-based abstractive text summarization baselines. We briefly elaborate on them below.

**Text-only extractive summarization:** (i) **Lead-2:** The top 2 sentences of the source documents are marked as the generated summary and evaluated against the target summary. (ii) **LexRank:** It [13] represents the source document sentence as nodes of a graph and edges as a similarity measure. The edge weights are computed using the eigenvector centrality and the token frequency. (iii) **TextRank:** Similar to LexRank, it [38] also represents the source document as a fully-connected graph. All edge weights are given a unit weight, and later a derived version of PageRank re-ranks the edge weights. (iv) **MMR (Maximal Marginal Loss):** The redundancy between the sentences is computed by mapping the query to the sentence [6]. The relevant sentences are kept in a cluster and filtered based on the similarity ranking. (v) **ICSISumm:** The coverage of the sentence in the final summary is optimized using the linear optimization framework [18]. Given a summary length bound, integer linear programming (ILP) solvers try to maximize the global topic coverage. (vi) **BERTExtrative:** The task of summarization is transformed into a binary classification problem [35]. The sentences are classified into two classes representing whether the sentence is a part of the final summary or not.

**Text-only abstractive summarization:** (i) **Seq2Seq:** It [41] uses the standard RNN network for both encoder and decoder with a global attention mechanism. (ii) **Pointer Generator (PG):** It [46] extends the Seq2Seq network with the addition of the Pointing mechanism. The Pointing mechanism allows the network to either generate tokens from the vocabulary or directly copy from the source document. (iii) **CopyTransformer:** It [15] uses the standard Transformer network. Similar to PG, a random attention head acts as a pointing mechanism. It also leverages a content selection module to enrich the generated summary. (iv) **Longformer:** Unlike the standard Transformer [54], Longformer [4] uses linear attention to cater to the long source document. The computed attention weights are a combination of global and windowed attention. (v) **BERT:** It [11] is an encoder-only language model trained on the token masking technique. We fine-tune BERT over the text-only setting till convergence. (vi) **BART:** It [27] is an encoder-decoder-based language model pre-trained on the phrase masking technique. Similar to BERT, we fine-tune BART till convergence. (vii) **T5:** It [44] considers all NLP downstream tasks as text-to-text problem. As text-to-text uses the same model architecture and loss throughout all NLP problems. We fine-tune t5-base on the mTLDR training

data. (viii) **Pegasus:** Similar to T5, Pegasus [59] is also an encoder-decoder based model. However, the self-supervision objective is to mask sentences rather than token masking, helping the model generate contextual sentences. This pre-training objective vastly helps text generation tasks like summarization.

**Video-only abstractive summarization:** (i) **Action features only:** The video feature representations are trained over a convolution layer and passed through attention and an RNN-based decoder. (ii) **RNN (Action features):** The video features are trained over the convolution layer and passed through attention and an RNN-based encoder. The latent representation is finally fused using the hierarchical attention and passed onto the RNN-based decoder to generate the final summary.

**Multimodal abstractive summarization:** (i) **HA**: The work of [31] is repurposed for the multimodal summarization task. The visual and textual features are fused with hierarchical attention [42] to align features and capture more context while generating the summary. (ii) **MFFG:** It [33] introduces multistage fusion with forget gate. The encoder part uses a cross-attention-based fusion with forget gates. The decoder is assembled using a hierarchical fusion network to capture only the important concepts and forget redundant information. (iii) **FLORAL:** It [1] proposes a Factorized Multimodal Transformer based language model consisting of guided attention and multimodal attention layer to align attention scores of each modality and use speech and OCR text to guide the generated summary. (iv) **ICAF:** It [61] utilises recurrent and contrastive alignment to capture the relationship between the video and text. It makes use of contrastive loss to align modalities in the embedding space resulting in enriched aligned summaries.

## 6 EXPERIMENTS

We perform extensive ablations to evaluate the efficacy of the proposed modules of mTLDRgen and individual modalities. We explore how text, visual and acoustic features perform separately and jointly over mTLDR and How2.

**Evaluation measures:** We evaluate the performance of mTLDRgen using both quantitative metrics - Rouge-1, Rouge-2, and Rouge-L and qualitative metrics – BERTScore and FEQA. Rouge measures the recall of unigrams (Rouge-1), bigrams (Rouge-2), and $n$-grams (Rouge-L) between the generated and target summaries. BERTScore assesses the similarity between the generated and target summaries in the embedding space through the cosine similarity of the BERT embeddings. FEQA, a question-answer generation metric, evaluates the quality of the generated summaries by determining the number of answers mapped to questions generated from the target summaries.

We further perform human evaluations[4] In the evaluations, we benchmark the summaries over the following parameters — Informativeness, Fluency, Coherence and Relevance (c.f. Supplementary, Section B). We randomly sample 40 instances from the test set and evaluate them against the target summaries. We perform

---

[4]The human evaluations were performed by 15 individuals with sufficient background in NLP, machine learning and deep learning. The participants were aged between 22-28 years.

**Table 2: Performance benchmark over six text-only Extractive (Extr) baselines (Lead-2, Lexrank, TextRank, MMR, ICSISumm, and BERTExtractive), eight text-only Abstrative (Abst) baselines (Seq2Seq, Pointer Generator (PG), CopyTransformer, Long-former, BERT, BART, T5, and Pegasus), two video-only baselines (Action feature, and Action feature with RNN), and four Multimodal baselines (HA, FLORAL, MFFG, and ICAF) over the datasets – `mTLDRgen` and How2. The benchmarks are evaluated over the Quantitative metric – Rouge (Rouge-1 (R1), Rouge-2 (R2), and Rouge-L (RL)), and Qualitative metric – BERTScore (BERTSc.) and FEQA.**

| Type | System | mTLDR | | | | | How2 | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | R1 | R2 | RL | BERTSc. | FEQA | R1 | R2 | RL | BERTSc. | FEQA |
| Extr-text | Lead-2 | 22.82 | 4.61 | 15.47 | 61.27 | 32.45 | 43.96 | 13.31 | 39.28 | 71.56 | 32.28 |
| | LexRank | 27.18 | 6.82 | 17.22 | 63.23 | 34.21 | 27.93 | 12.88 | 16.93 | 64.52 | 31.89 |
| | TextRank | 27.43 | 6.86 | 17.41 | 63.34 | 34.29 | 27.49 | 12.61 | 16.71 | 64.55 | 31.92 |
| | MMR | 29.54 | 8.19 | 18.84 | 64.59 | 35.67 | 28.24 | 13.12 | 17.86 | 64.87 | 31.98 |
| | ICSISumm | 31.57 | 9.52 | 19.42 | 65.84 | 36.14 | 28.53 | 13.44 | 17.93 | 65.14 | 32.16 |
| | BERTExtractive | 31.52 | 9.49 | 19.31 | 65.83 | 36.13 | 27.18 | 12.47 | 15.38 | 63.47 | 31.67 |
| Abst-text | Seq2Seq | 23.54 | 5.61 | 15.48 | 62.47 | 31.57 | 55.37 | 23.08 | 53.86 | 76.15 | 36.48 |
| | PG | 23.59 | 5.78 | 16.21 | 62.71 | 31.84 | 51.68 | 22.63 | 50.29 | 73.47 | 35.37 |
| | CopyTransformer | 25.63 | 7.82 | 18.54 | 63.11 | 37.86 | 52.94 | 23.25 | 50.26 | 73.58 | 35.43 |
| | Longformer | 21.37 | 6.47 | 15.12 | 61.05 | 32.14 | 49.24 | 21.39 | 47.41 | 72.39 | 35.28 |
| | BERT | 24.87 | 8.85 | 18.33 | 62.91 | 31.89 | 53.74 | 23.86 | 48.06 | 73.45 | 35.62 |
| | BART | 26.13 | 9.69 | 19.62 | 64.24 | 38.64 | 53.81 | 23.89 | 48.15 | 73.51 | 35.68 |
| | T5 | 25.87 | 9.24 | 18.63 | 64.13 | 38.45 | 53.21 | 22.51 | 47.48 | 73.42 | 35.65 |
| | Pegasus | 26.66 | 9.83 | 19.26 | 64.85 | 36.98 | 53.87 | 23.91 | 48.17 | 73.61 | 35.70 |
| Video only | Action features only | 26.38 | 6.47 | 15.37 | 62.48 | 30.41 | 45.24 | 24.42 | 38.47 | 69.74 | 31.28 |
| | RNN (Action features) | 26.73 | 6.51 | 15.75 | 63.14 | 31.35 | 48.27 | 27.74 | 46.37 | 72.32 | 35.11 |
| Multimodal | HA | 29.32 | 11.84 | 26.18 | 67.24 | 39.37 | 55.82 | 38.31 | 54.96 | 77.15 | 38.55 |
| | FLORAL | 31.69 | 13.54 | 31.55 | 69.56 | 41.19 | 56.84 | 39.86 | 56.93 | 79.84 | 39.14 |
| | MFFG | 33.19 | 18.88 | 33.28 | 71.54 | 43.13 | 61.49 | 44.61 | 57.21 | 80.16 | 41.59 |
| | ICAF | *36.38* | *20.54* | *34.52* | *73.94* | *45.63* | *63.84* | *44.78* | *58.24* | *82.39* | *42.86* |
| | **mTLDRgen** | **41.62** | **22.69** | **37.87** | **78.39** | **48.46** | **67.33** | **48.71** | **61.83** | **84.11** | **44.82** |
| Δ_mTLDRgen | BEST | ↑ 5.24 | ↑ 2.15 | ↑ 3.35 | ↑ 4.45 | ↑ 2.83 | ↑ 3.49 | ↑ 3.93 | ↑ 3.59 | ↑ 1.72 | ↑ 1.96 |

**Table 3: Ablation study to show the efficacy of each module of `mTLDRgen`.**

| System | mTLDR | | | | How2 | | | |
|---|---|---|---|---|---|---|---|---|
| | Rouge-1 | Rouge-2 | Rouge-L | BERTScore | Rouge-1 | Rouge-2 | Rouge-L | BERTScore |
| Transformer | 25.63 | 7.82 | 18.54 | 63.11 | 52.94 | 23.25 | 50.26 | 73.58 |
| + DFHC | 29.37 | 11.78 | 23.19 | 67.81 | 57.34 | 28.71 | 56.02 | 77.31 |
| + WRET | 34.52 | 14.82 | 26.54 | 72.06 | 61.12 | 36.89 | 58.1 | 81.44 |
| + DFHC & WRET | 37.34 | 18.32 | 32.49 | 74.58 | 64.23 | 42.61 | 59.02 | 82.45 |
| mTLDRgen | 41.62 | 22.69 | 37.87 | 78.39 | 67.33 | 48.71 | 61.83 | 84.11 |

human evaluations for BART (text-only), T5 (text-only), MFFG (multimodal), FLORAL (multimodal) and `mTLDRgen` (multimodal).

## 6.1 Training

The implementation of the `mTLDRgen` model was carried out by utilizing the Pytorch 1.8.1 framework on an NVIDIA A6000 GPU equipped with 46 GB of dedicated memory and CUDA-11 and cuDNN-7 libraries. The model was initialized with pre-trained BART language model weights for the encoder and decoder and fine-tuned on the summarization dataset. In the implementation, the loss computation was only performed over the target sequence in adherence to the encoder-decoder paradigm. The hyper-parameters used in both the pre-training and fine-tuning phases are detailed in Section A.1, and Table 7 (Supplementary).

## 6.2 Quantitative Analysis

Table 2 compares the performance of `mTLDRgen` with that of its baselines across the How2 and mTLDR datasets. Our results demonstrate that `mTLDRgen` outperforms the best baseline, ICAF, with a Rouge-1 score of 41.62 and a Rouge-L score of 37.87, an improvement of +5.24 and +3.35, respectively. When benchmarked against the How2 dataset, `mTLDRgen` exhibits superior results, obtaining Rouge-1 of 67.33 and Rouge-L of 61.83, outperforming the best baseline (ICAF) by +3.49 and +3.59 points, respectively. With respect

KDD '23, August 6–10, 2023, Long Beach, CA, USA

Yash Kumar Atri, Vikram Goyal, & Tanmoy Chakraborty

**Table 4: Comparison of target summary with six models – Extractive (ICSISumm), Abstractive (Pointer Generator (PG), BART, Pegasus) and multimodal (ICAF and `mTLDRgen`) models.**

| Model | Output |
|---|---|
| Target | In this paper, we propose an adversarial multi-task learning framework, where the shared and private latent feature spaces donot interfere with each other. This task is achieved by introducing orthogonality constraints. |
| ICSISumm | To prevent the shared and private latent feature spaces from interfering with each other, we introduce two strategies: adversarial training and orthogonality constraints. |
| PG | propose multi-task learning for the generative , propose latent feature for multi task learning where the shared knowledge regarded as off the self knowledge and trasferred to new task. |
| BART | In this paper, we conduct experiment on 16 tasks demonstrate the benefits and propose multi-task learning framework, The dataset are shared and latent feature spaces. the dataset is prone. |
| Pegasus | In this paper, we propose an multitask learning framework, where we conduct experiments on 16 text classification tasks. our model is off the shelf and donot interfere with each other. |
| ICAF | we propose an adversarial multi-task framework, where we conduct experiments demostrating private feature space do not interefere with eachother. The model is regarded as off the shelf and transferred to new task. |
| `mTLDRgen` | In this paper, we propose an multi-task framework, the shared and private latent feature spaces not interferd with each other. We conduct experiments on 16 text classification tasks. |

to the best text-only abstractive baseline, Pegasus (Rouge-1 and Rouge-2) and BART (Rouge-L), `mTLDRgen` shows an improvement of +14.96 Rouge-1, +12.86 Rouge-2, and +18.25 Rouge-L. Similarly, `mTLDRgen` surpasses ICSISumm, the best extractive baseline, with an improvement of +10.05, +13.17, and +18.45 on Rouge-1, Rouge-2 and Rouge-L, respectively.

We also perform ablations to study the efficacy of individual modalities and modules of `mTLDRgen`. Table 5 demonstrates the performance improvements obtained when all three modalities are fused, while Table 3 showcases the contribution of individual modules of `mTLDRgen`. These results serve to affirm our hypothesis that models specifically designed for longer summarization sequences are inadequate in extreme summarization tasks and that the integration of multiple modalities with text modality enhances the quality of the generated summary.

***Congruency of multi-modalities***. The performance of various unimodal and multimodal text summarization systems is shown in Table 2. The results demonstrate that for unimodal variants, the lead2, which was reported to be a strong baseline for datasets like CNN/Dailymail [21] and MultiNews [14], fails to perform effectively, indicating that the latent structure of the scientific text is distinct, and the information has a heterogeneous distribution throughout the document. In a similar vein, the text-only abstractive baselines perform inadequately over both the How2 and mTLDR datasets. On the other hand, the extractive baselines, which are able to identify the prominent sentences that start with "we propose" or "we introduce", perform better than the text-only abstractive baselines yet still provide only a limited context of the whole article. Meanwhile, the two video-only baselines display performance that is comparable to the best abstractive baselines, signifying that multimodal features do indeed contribute to generating more informative and coherent summaries. No baselines using audio-only features were run as audio captures only the amplitude shift and intonations,

which do not constitute a sufficient feature set in the vector space. As indicated in Table 2, the multimodal baselines outperform the text-only and video-only baselines. The HA model outperforms the best abstractive baseline by +2.66 Rouge-1 and +6.92 Rouge-L, demonstrating the significance of combining multimodal signals with text-only modalities. The fusion of video with text helps the model attain better context in the vector space, even the audio features aid in the mutual alignment of modalities leading to more diverse and coherent summaries. Evidently, all the remaining multimodal baselines show a remarkable improvement in performance over all the text-only extractive, text-only abstractive, and video-based baselines.

**Table 5: Performance benchmark for each modality of `mTLDRgen`.**

| Modality | Rouge-1 | Rouge-2 | Rouge-L |
|---|---|---|---|
| Text +Audio | 27.46 | 7.47 | 19.62 |
| Audio +Video | 27.62 | 7.53 | 20.11 |
| Text +Video | 28.05 | 7.83 | 24.49 |
| Text +Audio +Video | 41.62 | 22.69 | 37.87 |

## 6.3 Qualitative Analysis

We also conduct a qualitative evaluation of the generated summaries utilizing BERTSCore and FEQA (c.f. Table 2). Both metrics use the text modality from the source and target to assess the quality. On the mTLDR data, `mTLDRgen` achieves 78.39 BERTSCore and 48.46 FEQA, surpassing the best baseline (ICAF) by +4.45 and +2.83 points, respectively. For the How2 dataset, `mTLDRgen` obtains 84.11 BERTSCore and 44.82 FEQA, outperforming the best baseline (ICAF) by +1.72 and +1.96 points, respectively. Similar to the quantitative benchmarks, the multimodal baselines outperform the text-only

**Table 6: Human evaluation scores over the metrics – Informativeness (Infor.), Fluency, Coherence, and Relevance for the text-based baselines (BART and T5), multimodal baselines (MFFG, FLORAL, and mTLDRgen) on the mTLDRgen and How datasets.**

| Modality | System | mTLDR | | | | How2 | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | Infor. | Fluency | Coherence | Relevance | Infor. | Fluency | Coherence | Relevance |
| Abstractive-text | BART | 2.81 | 2.51 | 2.94 | 2.85 | 2.34 | 2.37 | 2.46 | 2.54 |
| Abstractive-text | T5 | 2.78 | 2.49 | 2.81 | 2.74 | 2.33 | 2.28 | 2.43 | 2.54 |
| Multimodal | FLORAL | 3.2 | 3.03 | 3.02 | 3.11 | 3.13 | 3.14 | 3.08 | 3.13 |
| Multimodal | MFFG | 3.21 | 3.05 | 3.09 | 3.11 | 3.17 | 3.21 | 3.04 | 3.11 |
| Multimodal | mTLDRgen | **3.46** | **3.32** | **3.27** | **3.29** | **3.34** | **3.27** | **3.21** | **3.18** |

extractive, abstractive, and video-only baselines by a substantial margin.

## 6.4 Human Evaluation

The quantitative enhancements are further reinforced by human assessments. As shown in Table 6, mTLDRgenscores highest over the datasets – mTLDR and How2, demonstrating that the generated summaries are highly faithful, pertinent, and coherent in comparison to the other baselines. Although mTLDRgen demonstrates some deficiencies in the coherence criterion in the human evaluations, it still performs significantly better than the other baselines. A manual examination of the generated summaries and an analysis of the findings are presented in Section 6.5.

## 6.5 Error Analysis

The limitations of extractive and abstractive baselines in generating extreme summaries are evident. Extractive systems rely on direct copying of phrases from the source document, often resulting in a single-line summary containing limited information diversity. This is reflected in their performance compared to abstractive text-only and a few multimodal (HA and FLORAL) baselines, as seen in Table 4. The text-only abstractive baselines like Seq2Seq and PG fail to extract the paper's main contributions, while Transformer based methods like Longformer, BERT, etc., struggle to summarize the contributions in very few lines.

However, mTLDRgen stands out as it is able to condense the three key contributions of the source article into a single sentence, demonstrating superiority over the other baselines. A manual inspection of instances where mTLDRgenfailed to generate a good summary reveals that the cause was often due to noisy text modality extracted from the article pdf, leading to non-coherent connections between phrases. Further, the data noise in the video and audio modalities arising due to different aspect ratios of presentation and speaker and non-native English accent speakers adds to the perplexity of modality alignment.

## 7 DEPLOYMENT - CONTINUOUS HUMAN FEEDBACK

The performance improvements across the quantitative and qualitative metrics over the mTLDR dataset motivated us to assess mTLDRgen more rigorously. After controlled alpha testing of mTLDRgen using human evaluations, we deployed mTLDRgen as a web-based API (the technical details of API are discussed in Section A.2 (Supplementary)) in an in-house tool. The API is currently hosted on a

local server with an A6000 (48 GB) GPU, and the endpoints are accessible across the institute. For input, the API takes either a web URL consisting of direct links to the video and the article pdf or a separate file pointer to upload files directly from the local workstation. The current response time for the API is $2 - 3$ minutes, which is considerably high as a wait time. However, to cut down on the revert-back time, we cache all the responses to provide immediate output to the already processed queries. During the production stage, our aim will be to reduce the inference time to less than 1 minute. However, to achieve this, we will aim to optimize the model by reducing model parameters and distilling it. Catering to the data regulations, we remove all the videos and article pdf's from our system after processing and only store the generated summaries and metadata for mapping queries to the cached data. The metadata includes the article title, author description, keywords, and month/year of publishing. Further, we do not track users' identities nor store any user-specific information on our servers.

## 8 CONCLUSION

We introduced a novel task of extreme abstractive text summarization using multimodal inputs. we curated mTLDR, a unique large-scale dataset for extreme abstractive text summarization that encompasses videos, audio, and text, as well as both author-written and expert-annotated summaries. Subsequently, we introduced mTLDRgen, a novel model that employs a dual fused hyper-complex Transformer and a Wasserstein Riemannian Encoder Transformer to efficiently capture the relationships between different modalities in a hyper-complex and latent geometric space. The hyper-complex Transformer captures the intrinsic properties between the modalities, while the Wasserstein Riemannian Encoder Transformer captures the latent structures of the modalities in the latent space geometry, enabling the model to generate diverse sentences. To assess the mTLDRgen, we conducted thorough experiments on mTLDR and How2 datasets and compared their performance with 20 baselines. Overall, mTLDRgen demonstrated superior performance both qualitatively and quantitatively.

# REFERENCES

[1] Yash Kumar Atri, Shraman Pramanick, Vikram Goyal, and Tanmoy Chakraborty. 2021. See, hear, read: Leveraging multimodality with guided attention for abstractive text summarization. *Knowledge-Based Systems* 227 (2021), 107152. https://doi.org/10.1016/j.knosys.2021.107152

[2] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2015. Neural Machine Translation by Jointly Learning to Align and Translate. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, Yoshua Bengio and Yann LeCun (Eds.). http://arxiv.org/abs/1409.0473

[3] Marco Baroni. 2016. Grounding distributional semantics in the visual world. *Language and Linguistics Compass* 10, 1 (2016), 3–13.

[4] Iz Beltagy, Matthew E. Peters, and Arman Cohan. 2020. Longformer: The Long-Document Transformer. *arXiv:2004.05150* (2020).

[5] Isabel Cachola, Kyle Lo, Arman Cohan, and Daniel Weld. 2020. TLDR: Extreme Summarization of Scientific Documents. In *Findings of the Association for Computational Linguistics: EMNLP 2020*. Association for Computational Linguistics, Online, 4766–4777. https://doi.org/10.18653/v1/2020.findings-emnlp.428

[6] Jaime Carbinell and Jade Goldstein. 2017. The Use of MMR, Diversity-Based Reranking for Reordering Documents and Producing Summaries. *SIGIR Forum* 51, 2 (Aug. 2017), 209–210. https://doi.org/10.1145/3130348.3130369

[7] Jingqiang Chen and Hai Zhuge. 2018. Abstractive text-image summarization using multi-modal attentional hierarchical rnn. In *Proceedings of the 2018 Conference on EMNLP*. 4046–4056.

[8] Jianpeng Cheng and Mirella Lapata. 2016. Neural Summarization by Extracting Sentences and Words. In *Proceedings of the 54th Annual Meeting of the ACL (Volume 1: Long Papers)*. ACL, Berlin, Germany, 484–494. https://doi.org/10.18653/v1/P16-1046

[9] Tanya Chowdhury, Sachin Kumar, and Tanmoy Chakraborty. 2020. Neural Abstractive Summarization with Structural Attention. In *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence, IJCAI-20*, Christian Bessiere (Ed.). IJCAI, 3716–3722. https://doi.org/10.24963/ijcai.2020/514 Main track.

[10] Tonglee Chung, Yongbin Liu, and Bin Xu. 2020. Monotonic alignments for summarization. *Knowledge-Based Systems* 192 (2020), 105363. https://doi.org/10.1016/j.knosys.2019.105363

[11] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the ACL: Human Language Technologies, Volume 1 (Long and Short Papers)*. Association for Computational Linguistics, Minneapolis, Minnesota, 4171–4186. https://doi.org/10.18653/v1/N19-1423

[12] Aniqa Dilawari and Muhammad Usman Ghani Khan. 2019. ASoVS: Abstractive Summarization of Video Sequences. *IEEE Access* 7 (2019), 29253–29263. https://doi.org/10.1109/ACCESS.2019.2902507

[13] G. Erkan and D. R. Radev. 2004. LexRank: Graph-based Lexical Centrality as Salience in Text Summarization. *Journal of Artificial Intelligence Research* 22 (Dec 2004), 457–479. https://doi.org/10.1613/jair.1523

[14] Alexander Fabbri, Irene Li, Tianwei She, Suyi Li, and Dragomir Radev. 2019. Multi-News: A Large-Scale Multi-Document Summarization Dataset and Abstractive Hierarchical Model. In *Proceedings of the 57th Annual Meeting of the ACL*. ACL, Florence, Italy, 1074–1084. https://doi.org/10.18653/v1/P19-1102

[15] Sebastian Gehrmann, Yuntian Deng, and Alexander Rush. 2018. Bottom-Up Abstractive Summarization. In *Proceedings of the 2018 Conference on EMNLP*. ACL, Brussels, Belgium, 4098–4109. https://doi.org/10.18653/v1/D18-1443

[16] Spandana Gella, Mike Lewis, and Marcus Rohrbach. 2018. A dataset for telling the stories of social media videos. In *Proceedings of the 2018 Conference on EMNLP*. 968–974.

[17] Michail N. Giannakos, Kshitij Sharma, Ilias O. Pappas, Vassilis Kostakos, and Eduardo Velloso. 2019. Multimodal data as a means to understand the learning experience. *International Journal of Information Management* 48 (2019), 108–119. https://doi.org/10.1016/j.ijinfomgt.2019.02.003

[18] Daniel Gillick, Benoit Favre, and Dilek Hakkani-Tür. 2008. The ICSI Summarization System at TAC 2008.. In *Tac*.

[19] Carlos-Emiliano Gonz'alez-Gallardo, Romain Deveaud, Eric SanJuan, and Juan-Manuel Torres. 2020. Audio Summarization with Audio Features and Probability Distribution Divergence. *ArXiv* abs/2001.07098 (2020).

[20] Anushka Gupta, Diksha Chugh, Anjum, and Rahul Katarya. 2022. Automated News Summarization Using Transformers. In *Sustainable Advanced Computing*, Sagaya Aurelia, Somashekhar S. Hiremath, Karthikeyan Subramanian, and Saroj Kr. Biswas (Eds.). Springer Singapore, Singapore, 249–259.

[21] Karl Moritz Hermann, Tomas Kocisky, Edward Grefenstette, Lasse Espeholt, Will Kay, Mustafa Suleyman, and Phil Blunsom. 2015. Teaching machines to read and comprehend. In *Advances in neural information processing systems*. 1693–1701.

[22] Vladimir Iashin and Esa Rahtu. 2020. Multi-modal Dense Video Captioning. In *Proceedings of the IEEE/CVF Conference on CVPR Workshops*. 958–959.

[23] James M. Joyce. 2011. *Kullback-Leibler Divergence*. Springer Berlin Heidelberg, Berlin, Heidelberg, 720–722. https://doi.org/10.1007/978-3-642-04898-2_327

[24] Will Kay, Joao Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijayanarasimhan, Fabio Viola, Tim Green, Trevor Back, Paul Natsev, et al. 2017. The kinetics human action video dataset. *CoRR* (2017).

[25] Douwe Kiela. 2017. *Deep embodiment: grounding semantics in perceptual modalities*. Technical Report. University of Cambridge, Computer Laboratory. 11–129 pages.

[26] Logan Lebanoff, Kaiqiang Song, and Fei Liu. 2018. Adapting the Neural Encoder-Decoder Framework from Single to Multi-Document Summarization. In *Proceedings of the 2018 Conference on EMNLP*. ACL, Brussels, Belgium, 4131–4141. https://doi.org/10.18653/v1/D18-1446

[27] Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Ves Stoyanov, and Luke Zettlemoyer. 2019. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. *arXiv preprint arXiv:1910.13461* (2019).

[28] Haoran Li, Junnan Zhu, Tianshang Liu, Jiajun Zhang, and Chengqing Zong. 2018. Multi-modal Sentence Summarization with Modality Attention and Image Filtering. In *Proceedings of the Twenty-Seventh IJCAI-18*. IJCAI, 4152–4158. https://doi.org/10.24963/ijcai.2018/577

[29] Haoran Li, Junnan Zhu, Cong Ma, Jiajun Zhang, and Chengqing Zong. 2018. Read, watch, listen, and summarize: Multi-modal summarization for asynchronous text, image, audio and video. *IEEE Transactions on Knowledge and Data Engineering* 31, 5 (2018), 996–1009.

[30] Zechao Li, Jinhui Tang, Xueming Wang, Jing Liu, and Hanqing Lu. 2016. Multimedia news summarization in search. *ACM Transactions on Intelligent Systems and Technology (TIST)* 7, 3 (2016), 1–20.

[31] Jindřich Libovický and Jindřich Helcl. 2017. Attention Strategies for Multi-Source Sequence-to-Sequence Learning. In *Proceedings of the 55th Annual Meeting of the ACL (Volume 2: Short Papers)*. ACL, Vancouver, Canada, 196–202. https://doi.org/10.18653/v1/P17-2031

[32] Jingzhou Liu, Dominic J. D. Hughes, and Yiming Yang. 2021. Unsupervised Extractive Text Summarization with Distance-Augmented Sentence Graphs. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval* (Virtual Event, Canada) *(SIGIR '21)*. Association for Computing Machinery, New York, NY, USA, 2313–2317. https://doi.org/10.1145/3404835.3463111

[33] Nayu Liu, Xian Sun, Hongfeng Yu, Wenkai Zhang, and Guangluan Xu. 2020. Multistage Fusion with Forget Gate for Multimodal Summarization in Open-Domain Videos. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Association for Computational Linguistics, Online, 1834–1845. https://doi.org/10.18653/v1/2020.emnlp-main.144

[34] Sheng Liu, Zhou Ren, and Junsong Yuan. 2018. SibNet: Sibling Convolutional Encoder for Video Captioning. In *Proceedings of the 26th ACM International Conference on Multimedia* (Seoul, Republic of Korea) *(MM '18)*. ACM, New York, NY, USA, 1425–1434. https://doi.org/10.1145/3240508.3240667

[35] Yang Liu. 2019. Fine-tune BERT for extractive summarization. *arXiv preprint arXiv:1903.10318* (2019).

[36] Kyle Lo, Lucy Lu Wang, Mark Neumann, Rodney Kinney, and Daniel Weld. 2020. S2ORC: The Semantic Scholar Open Research Corpus. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, Online, 4969–4983. https://doi.org/10.18653/v1/2020.acl-main.447

[37] Yuning Mao, Ming Zhong, and Jiawei Han. 2022. CiteSum: Citation Text-guided Scientific Extreme Summarization and Domain Adaptation with Limited Supervision. https://doi.org/10.48550/ARXIV.2205.06207

[38] Rada Mihalcea and Paul Tarau. 2004. TextRank: Bringing Order into Text. In *Proceedings of the 2004 Conference on EMNLP*. ACL, Barcelona, Spain, 404–411. https://www.aclweb.org/anthology/W04-3252

[39] Jonghwan Mun, Minsu Cho, and Bohyung Han. 2017. Text-guided Attention Model for Image Captioning. In *AAAI*. 4233–4239.

[40] Ramesh Nallapati, Feifei Zhai, and Bowen Zhou. 2017. SummaRuNNer: A Recurrent Neural Network Based Sequence Model for Extractive Summarization of Documents. In *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence* (San Francisco, California, USA) *(AAAI'17)*. AAAI Press, 3075–3081.

[41] Ramesh Nallapati, Bowen Zhou, Cicero dos Santos, Çağlar Guİ‡lçehre, and Bing Xiang. 2016. Abstractive Text Summarization using Sequence-to-sequence RNNs and Beyond. In *Proceedings of The 20th SIGNLL Conference on Computational Natural Language Learning*. ACL, Berlin, Germany, 280–290. https://doi.org/10.18653/v1/K16-1028

[42] Shruti Palaskar, Jindřich Libovický, Spandana Gella, and Florian Metze. 2019. Multimodal Abstractive Summarization for How2 Videos. In *Proceedings of the 57th Annual Meeting of the ACL*. ACL, Florence, Italy, 6587–6596. https://doi.org/10.18653/v1/P19-1659

[43] Jielin Qiu, Jiacheng Zhu, Mengdi Xu, Franck Dernoncourt, Trung Bui, Zhaowen Wang, Bo Li, Ding Zhao, and Hailin Jin. 2022. MHMS: Multimodal Hierarchical Multimedia Summarization. https://doi.org/10.48550/ARXIV.2204.03734

[44] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer. *Journal of Machine*

*Learning Research* 21, 140 (2020), 1–67. http://jmlr.org/papers/v21/20-074.html

[45] Ramon Sanabria, Ozan Caglayan, Shruti Palaskar, Desmond Elliott, Loïc Barrault, Lucia Specia, and Florian Metze. 2018. How2: A Large-scale Dataset For Multimodal Language Understanding. In *Proceedings of the Workshop on Visually Grounded Interaction and Language (ViGIL)*. NeurIPS, 26.1–26.12. http://arxiv.org/abs/1811.00347

[46] Abigail See, Peter J. Liu, and Christopher D. Manning. 2017. Get To The Point: Summarization with Pointer-Generator Networks. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. ACL, Vancouver, Canada, 1073–1083. https://doi.org/10.18653/v1/P17-1099

[47] Rajiv Ratn Shah, Yi Yu, Akshay Verma, Suhua Tang, Anwar Dilawar Shaikh, and Roger Zimmermann. 2016. Leveraging multimodal information for event summarization and concept-level sentiment analysis. *Knowledge-Based Systems* 108 (2016), 102–109.

[48] Xiangqing Shen, Bing Liu, Yong Zhou, Jiaqi Zhao, and Mingming Liu. 2020. Remote sensing image captioning via Variational Autoencoder and Reinforcement Learning. *Knowledge-Based Systems* (2020), 105920.

[49] Xiangxi Shi, Jianfei Cai, Jiuxiang Gu, and Shafiq Joty. 2020. Video captioning with boundary-aware hierarchical language decoding and joint video prediction. *Neurocomputing* 417 (2020), 347–356.

[50] Kaiqiang Song, Lin Zhao, and Fei Liu. 2018. Structure-Infused Copy Mechanisms for Abstractive Summarization. In *Proceedings of the 27th International Conference on Computational Linguistics*. Association for Computational Linguistics, Santa Fe, New Mexico, USA, 1717–1729. https://www.aclweb.org/anthology/C18-1146

[51] Reuben Tan, Bryan A. Plummer, Kate Saenko, JP Lewis, Avneesh Sud, and Thomas Leung. 2022. NewsStories: Illustrating Articles with Visual Summaries. In *Computer Vision – ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXVI* (Tel Aviv, Israel). Springer-Verlag, Berlin, Heidelberg, 644–661. https://doi.org/10.1007/978-3-031-20059-5_37

[52] Joseph Tepperman, David Traum, and Shrikanth Narayanan. 2006. Yeah right: Sarcasm recognition for spoken dialogue systems. 1838–1841.

[53] Ilya O Tolstikhin, Bharath K. Sriperumbudur, and Bernhard Schölkopf. 2016. Minimax Estimation of Maximum Mean Discrepancy with Radial Kernels. In *Advances in Neural Information Processing Systems*, D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett (Eds.), Vol. 29. Curran Associates, Inc. https://proceedings.neurips.cc/paper/2016/file/5055cbf43fac3f7e2336b27310f0b9ef-Paper.pdf

[54] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in neural information processing systems*. 5998–6008.

[55] Michael Völske, Martin Potthast, Shahbaz Syed, and Benno Stein. 2017. TL;DR: Mining Reddit to Learn Automatic Summarization. In *Workshop on New Frontiers in Summarization at EMNLP 2017*, Giuseppe Carenini, Jackie Chi Kit Cheung, Fei Liu, and Lu Wang (Eds.). Association for Computational Linguistics, 59–63. https://doi.org/10.18653/v1/W17-4508

[56] Prince Zizhuang Wang and William Yang Wang. 2019. Riemannian Normalizing Flow on Variational Wasserstein Autoencoder for Text Modeling. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. Association for Computational Linguistics, Minneapolis, Minnesota, 284–294. https://doi.org/10.18653/v1/N19-1025

[57] Kuo-Hao Zeng, Tseng-Hung Chen, Juan Carlos Niebles, and Min Sun. 2016. Title Generation for user generated videos. In *European conference on computer vision*. Springer, 609–625.

[58] Aston Zhang, Yi Tay, Shuai Zhang, Alvin Chan, Anh Tuan Luu, Siu Cheung Hui, and Jie Fu. 2021. Beyond Fully-Connected Layers with Quaternions: Parameterization of Hypercomplex Multiplications with $1/n$ Parameters. https://doi.org/10.48550/ARXIV.2102.08597

[59] Jingqing Zhang, Yao Zhao, Mohammad Saleh, and Peter Liu. 2020. Pegasus: Pre-training with extracted gap-sentences for abstractive summarization. In *International Conference on Machine Learning*. PMLR, 11328–11339.

[60] Zhengkun Zhang, Xiaojun Meng, Yasheng Wang, Xin Jiang, Qun Liu, and Zhenglu Yang. 2021. UniMS: A Unified Framework for Multimodal Summarization with Knowledge Distillation. https://doi.org/10.48550/ARXIV.2109.05812

[61] Zijian Zhang, Chang Shu, Youxin Chen, Jing Xiao, Qian Zhang, and Lu Zheng. 2021. ICAF: Iterative Contrastive Alignment Framework for Multimodal Abstractive Summarization. https://doi.org/10.48550/ARXIV.2108.05123

[62] Hao Zheng and Mirella Lapata. 2019. Sentence centrality revisited for unsupervised summarization. *arXiv preprint arXiv:1906.03508* (2019).

[63] Luowei Zhou, Chenliang Xu, and Jason J Corso. 2018. Towards Automatic Learning of Procedures From Web Instructional Videos. In *AAAI Conference on Artificial Intelligence*. 7590–7598. https://www.aaai.org/ocs/index.php/AAAI/AAAI18/paper/view/17344

[64] Junnan Zhu, Haoran Li, Tianshang Liu, Yu Zhou, Jiajun Zhang, and Chengqing Zong. 2018. MSMO: Multimodal Summarization with Multimodal Output. In *Proceedings of the 2018 Conference on EMNLP*. ACL, Brussels, Belgium, 4154–4164. https://doi.org/10.18653/v1/D18-1448

## A  EXPERIMENTATION AND DEPLOYMENT

We discuss the experimentation environment and the deployment parameters in this section.

### A.1  Experimentation Details

We ran several sets of experiments over our proposed mTLDRgen model to figure out the most optimal set of hyperparameters. Table 7 describes the most optimal hyperparameters for mTLDRgen over the mTLDR dataset. We initialize the mTLDRgen weights using the pre-trained BART and then train the network further over the mTLDR data. During training, we set the learning rate as $3e − 5$ and set the lambda value as 18. The gradients are accumulated for five iterations, and tri-gram blocking is used to penalize the decoder.

**Table 7: HyperParameters used to train mTLDRgen.**

| Parameter | Value |
| --- | --- |
| Epochs | 55 |
| Accumulate gradient steps | 5 |
| Ranking loss margin | 0.001 |
| MLE weight | 0.1 |
| Warmup steps | 10000 |
| Max learning rate | 3e-5 |
| Max source length | 512 |
| Training max summary length | 36 |
| Testing max summary length | 40 |
| Num of Beams | 4 |
| GPU | 2 X A6000 |
| VMemory | 48GB |

We train the model till loss converges and the validation accuracy does not improve for five continuous iterations. Figure 3 shows the variation of loss till 2000 iterations.
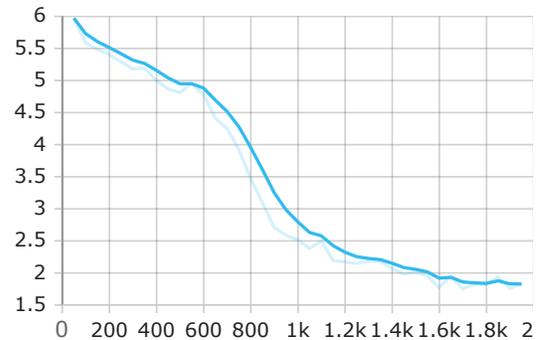


**Figure 3: Loss momentum of mTLDRgen over the mTLDR dataset.**

### A.2  Deployment

The web API for mTLDRgen has been created using the Flask [5] framework. As Flask only allows single-user access at a given time, the

---
[5] https://flask.palletsprojects.com/en/2.2.x/

API is running on top of the Gunicorn [6] framework, allowing multiple users to access the API at the same time. At a particular instance, the API is able to handle a load of four concurrent requests without major variation in inference time. The trained `mTLDRgen` model is hosted separately as an API running on Docker, giving the users option to either generate TLDR summaries using the web interface or directly by calling the `mTLDRgen` using any programming language.

## B  HUMAN EVALUATION SETUP

We evaluate the generated summaries over four [14] parameters – Informativeness, Fluency, Coherence and Relevance. Figure **??** shows the form utilised by human evaluators to benchmark the generated summaries against the competing baselines.

(1) Informativeness: The generated summary should house a certain level of information. The information can be in direct correlation with the source document or the target summary.

(2) Fluency: It encapsulates how the individual sentence stands in the generated summary. Every sentence in the summary should be grammatically and syntactically correct and should have no capitalization or punctuation errors.

(3) Coherence: It analyzes how the summary as a whole makes sense. The summary should be human-readable and should make sense contextually.

(4) Relevance: It computes how much information from the source document is available in the generated summary. The information on the generated summary should only come from the source document; any information generated outside the source document is termed as a hallucinating summary.

---

[6]https://gunicorn.org/

## Human Evalution of abstractive text summarization

Please read the abstract, watch the video and rate the generated summaries over these four parameters -- Informativeness, Fluency, Coherence, and Relevance

Informativeness: The generated summary should house a certain level of information. The information can be in direct correlation with the source document or the target summary.

Fluency: It encapsulates how the individual sentence stands in the generated summary. Every sentence in the summary should be grammatically and syntactically correct and should have no capitalization or punctuation errors.

Coherence: It analyzes how the summary as a whole makes sense. The summary should be human-readable and should make sense contextually.

Relevance: It computes how much information from the source document is available in the generated summary. The information on the generated summary should only come from the source document; any information generated outside the source document is termed as a hallucinating summary.

yashkumaratri@gmail.com (not shared) Switch accounts

Abstract :

Neural network models have shown their promising opportunities for multi-task learning, which focus on learning the shared layers to extract the common and task-invariant features. However, in most existing approaches, the extracted shared features are prone to be contaminated by task-specific features or the noise brought by other tasks. In this paper, we propose an adversarial multi-task learning framework, alleviating the shared and private latent feature spaces from interfering with each other. We conduct extensive experiments on 16 different text classification tasks, which demonstrates the benefits of our approach. Besides, we show that the shared knowledge learned by our proposed model can be regarded as off-the-shelf knowledge and easily transferred to new tasks. The datasets of all 16 tasks are publicly available at http://nlp.fudan.edu.cn/data/.

Video: Adversarial Multi-task Learning for Text Classification —- Pengfei Liu, Xipeng Qiu and Xuanjing Huang on Vimeo

Your answer

Summary 1: we propose an adversarial multi-task framework, where we conduct experiments demostrating private feature space do not interefere with eachother. The model is regarded as off the shelf and transferred to new task.

Summary 2: In this paper, we propose an multi-task framework, the shared and private latent feature spaces not interfere with each other. We conduct experiments on 16 text classification tasks.

Which will you prefer ?

○ Summary 1

○ Summary 2

○ Neither

○ Both

Informativeness score for the preferred summary [1-5]

Your answer

Fluency score for the preferred summary [1-5]

Your answer

Coherency score for the preferred summary [1-5]

Your answer

Relevance score for the preferred summary [1-5]

Your answer

**Figure 4: Human Evaluation form for collection feedback over the mTLDRgen and baseline generated summaries.**