
On abstractive and extractive summarization of instructional video transcripts using BERT

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 The overflow of video content in the Internet (from YouTube, MOOCs, news
2 portals) necessitates automated summarizations of data. In our paper, we study
3 extractive and abstractive summarization for of instructional videos. Previously,
4 natural language processing efforts have been focused to meticulously curated
5 datasets far removed from textual inconsistencies that are inherent to videos. Our
6 work on text preprocessing allows to extend the approach summarization of auto-
7 generated amateur video transcripts. Next, we apply state-of-the-art pretrained
8 BERT transformer models to the problem and evaluate the efficiency of training
9 and fine tuning with datasets from WikiHow, How2 videos, and CNN. The results
10 are evaluated using ROUGE, Content F1, and blind assessments by human experts.

11 1 Introduction

12 According to Forbes, more than 500 million hours of videos are watched on YouTube every day and
13 a lot of time is wasted watching videos that are not useful. Video content is rapidly growing and will
14 remain the mainstream for sharing information in future. In this project YAVA (“Your Active Virtual
15 Audience”) we are aspiring to make online exchanges of information between people via audio or
16 video more efficient and enjoyable.

17 There have been a lot of research efforts recently focused on video summarization [e.g. see [Cai
18 et.al.], [Shemer et.al.], [Kaufman et.al.]]. The known methods work by extracting the most important
19 segments and concatenating them together. However, it has been demonstrated that a lot of the
20 time the result is not substantially better and sometimes even worse than random selection of video
21 fragments ([Mayu et.al.]).

22 Summarization in the area of multimodal video processing tackles the problem from a different angle.
23 Instead of producing summary by converting a long video into a short video, it extracts signals from
24 it speech-to-text, facial expressions, spectrogram of speaker’s voice; etc. (see [Samanth et.al.]). and
25 processes them separately produces a short text (an abstract of what it is about). This method has a
26 few advantages:

- 27 • We get access to a set of existing models for text summarization, substantially more mature
28 than those for videos (e.g. [Subramanian et.al.]).
- 29 • We can leverage existing text summarization datasets, which are more easily available, than
30 video datasets (e.g. [Mahnaz et.al.]).
- 31 • Processing texts during algorithm training takes less computational power than processing
32 videos.
- 33 • Arguably, a text summary of a long video is even better for the viewer than a short video,
34 especially from the perspective of a person who needs it to decide whether to watch the

35 full video. It doesn't consume the network bandwidth, doesn't require audio equipment
36 or noise-free environment, takes less device energy to reproduce (especially important for
37 mobile devices) , and the viewer can consume it at their own pace. You can skim the text in
38 any order, any time.

39 The models for these purposes have been developed better than for processing video as a whole (e.g.
40 see [Jaejin et.al.]), and that's why this approach referred to as "multimodal" summarization looks
41 very promising to us and has recently received a lot of attention from other researchers (e.g. see
42 [Palaskar et.al.], [Tripathi et.al.]).

43 The focus of our research is on how-to/instructions videos. According to [https://mediakix.com/
44 blog/most-popular-youtube-videos/](https://mediakix.com/blog/most-popular-youtube-videos/), this type of video is one of the most popular on youtube
45 these days. Also, viewers of such videos are interested in getting a tangible outcome, as compared to
46 viewers of entertainment or sports videos, therefore adding a summary will add more value. which
47 we will use for training purposes. Pioneering efforts in this area have been done by [Palaskar et.al.]
48 based on dataset of how2 videos [Sanabria et.al.]. We plan to improve on their results by taking
49 advantage of "WikiHow: A Large Scale Text Summarization Dataset" [Mahnaz et.al.], improving
50 the models, and applying more advanced techniques to evaluation of output. Why is it important /
51 challenging? We foresee many applications of this approach, especially in education and business,
52 where even minor improvements in information processing may make big differences when applied
53 at scale to online meetings, virtual classrooms and other forms of human interactions via video.

54 Summarizing content is challenging even for a human. The rules of identifying what's important and
55 what can be omitted are subjective, changeable and very hard to formalize. While watching a long
56 video conference, participants often get tired and lose attention. Finally, a lot depends on the context.
57 Yet, as hard as it is, most people get it, and this skill improves through a lot of learning and practice.
58 It gives us hope that training machines to help facilitate this process is both possible and useful.

59 Also, evaluating the quality of summaries and obtaining benchmarks is another problem. As shown in
60 research [Mayu et.al.], engaging human experts for evaluation of results is expensive, while automated
61 techniques lack depth. We will use a combination of both techniques to maximize the quality of
62 results.

63 In our work, we are exploring transferability of modern text summarization techniques to instructional
64 videos scripts on large annotated data sets that we created by preprocessing YouTube videos and data
65 from other authors. We discuss heuristics that were discovered on this data, impacts on the quality of
66 generated summaries, and propose different ways of improving summarization process to deal with
67 these issues. Finally, we identify promising directions for future research.

68 2 Prior work

69 2.1 Text Summarization

70 Text summarization is the task of generating shorter versions of documents while maintaining
71 important information [need link]. This area of research in the natural language processing community
72 has grown rapidly over the past several years due to its practical applications among various industries
73 such as news, reviews, education. Summarization systems take two general approaches: extractive and
74 abstractive. Extractive summarization provides users with textual summaries that have been copied
75 and concatenated from important parts of a document. It is a reliable task capable of maintaining
76 sentence structure and factual correctness. Abstract summarization generates a summary with content
77 that is not always found in the underlying text. It is a complex task that mimics human summarization
78 by generalizing and paraphrasing key points made in the document. Prior to 2014, summarization
79 was centered on extracting lines from single documents using statistical models and neural networks
80 had limited success[6, 7]. Sutskever et al. and Cho et al work on sequence to sequence models opened
81 up new possibilities for neural networks in natural language processing. From 2014 to 2015, LSTMs
82 (variety of RNN) became the dominant approach that achieved state of the art results. They became
83 successful in tasks such as speech recognition, machine translation, parsing, image captioning, etc. It
84 paved the way for abstractive summarization, which began to score competitively against extractive
85 summarization. In 2017, Attention is all you need [8] provided a solution to the 'fixed length vector'
86 problem, enabling neural networks to focus on important parts of the input for prediction tasks.
87 Transformers with attention became more dominant for certain tasks [9].

88 2.2 Multi-modal Summarization

89 Research surrounding multimedia has improved greatly to bridge the gaps between multi-modal
90 content such as speech, visuals, and text. Summarization has been used in meeting records [10],
91 sports videos [11], news [12], each encapsulating synchronized speech, videos, and subtitles. Video
92 summaries consist of cutting important frames out of the video to create a succinct compact version.
93 More recently, research around multimodal summarization, which combines the textual and visual
94 modalities to align with the video content, have reached an early benchmark [13 - shruti's work].
95 The How2Dataset [5] is a collection of 2,000 hours of instructional videos with English subtitles and
96 crowdsourced Portuguese translations. It covers different how-to domains such as sports, cooking,
97 and education. The dataset has been created to be used as a benchmark for multimodal natural
98 language tasks, used in various competitions and research settings. This How2Dataset precedes
99 more recent work constructing data from instructional web videos in the How2100M [14] dataset.
100 The dataset is large-scale and has 136 million video clips and transcripts of humans performing or
101 describing various tasks, but there are no human annotated summaries.

102 3 Problem Statement

103 In our work we set the following goals:

- 104 • Curate and publish a single source of truth data set of text and summaries aggregated and
105 formatted from WikiHow articles, How2 videos, and CNN stories
- 106 • Apply existing BERT-based text summarization models to make them applicable to auto-
107 generated scripts from instructional videos and generalize them to work on instructional
108 videos
- 109 • Augment ROUGE metrics [Chin-Yew Lin] for evaluation of the results with a framework
110 for formalized expert assessment based on our research and criteria proposed by previous
111 works

112 4 Proof of concept

113 For our confidence about the feasibility of the project, we first ran a series of manual experiment by
114 dumping a few auto-generated scripts YouTube scripts and running them through online summariza-
115 tion services. The first results were very disappointing. However, we noticed that auto-generated
116 scripts don't have punctuation and line breaks don't necessarily correspond to the logical ends of
117 sentences. After fixing these issues, we got meaningful summaries and proceeded to generalizing the
118 approach as follows.

119 5 Methodology

120 From the initial exploration and data analysis we saw that in the process of applying existing
121 summarization models to Youtube video scripts we will deal with challenges imposed by parsing
122 speech-to-text output add more complexity to text summarization. For example, in one of the sample
123 videos in our test data set closed captioning confuses the speaker's words "*how you get a text from*
124 *a YouTube video*" for "*how you get attacks from a YouTube video*". So, our work includes several
125 iterations of the process described below:

- 126 • Collection and aggregation of data from multiple sources (HowTo video scripts, WikiHow,
127 CNN stories, YouTube)
- 128 • Preprocessing of video scripts to make them fit the text summarization models (e.g. errors in
129 word recognition, lack of punctuation in closed captioning, getting rid of special characters
130 etc.)
- 131 • Aligning inputs aggregated from multiple sources to common format
- 132 • Selection, deployment, training, fine-tuning of text summarization models
- 133 • Running experiments applying models to the data
- 134 • Evaluation of the outputs using ROUGE metrics and human expert judgement

135 5.1 Collection

136 Bryan or Alexandra - write about data sources, sizes in samples and bytes, metadata

137 5.2 Preprocessing

138 The work we have done at this stage is based on the heuristics observed during evaluation of results.
139 We expected the differences in conversational style of the video scripts and writtent text of CNN stories
140 (on which the models were pretrained) will impact quality of the output. In our first experiments,
141 it manifested in a very distinct way. The model considered the first one-two sentences to be very
142 important for summaries, and we ended up with getting many summaries looking like "hi!" and
143 "hello, this is <first and last name>". It inspired us for implementing an improvement by using entity
144 detection `spacy` and `nltk` to remove introduction from the text that we feed to summarization model.

145 5.3 Input format

146 Bryan - write about differences in format of CNN , wikiHow, howTo scripts and how they were
147 converted

148 5.4 Summarization models

149 Bryan or Alexandra

150 5.5 Evaluation

151 Alexandra or Bryan

152 The style files for NeurIPS and other conference information are available on the World Wide Web at

153 <http://www.neurips.cc/>

154 The file `neurips_2020.pdf` contains these instructions and illustrates the various formatting re-
155 quirements your NeurIPS paper must satisfy.

156 The only supported style file for NeurIPS 2020 is `neurips_2020.sty`, rewritten for $\text{\LaTeX} 2_{\epsilon}$.
157 **Previous style files for \LaTeX 2.09, Microsoft Word, and RTF are no longer supported!**

158 The \LaTeX style file contains three optional arguments: `final`, which creates a camera-ready copy,
159 `preprint`, which creates a preprint for submission to, e.g., arXiv, and `nonatbib`, which will not
160 load the `natbib` package for you in case of package clash.

161 **Preprint option** If you wish to post a preprint of your work online, e.g., on arXiv, using the
162 NeurIPS style, please use the `preprint` option. This will create a nonanonymized version of your
163 work with the text “Preprint. Work in progress.” in the footer. This version may be distributed as
164 you see fit. Please **do not** use the `final` option, which should **only** be used for papers accepted to
165 NeurIPS.

166 At submission time, please omit the `final` and `preprint` options. This will anonymize your
167 submission and add line numbers to aid review. Please *do not* refer to these line numbers in your
168 paper as they will be removed during generation of camera-ready copies.

169 The file `neurips_2020.tex` may be used as a “shell” for writing your paper. All you have to do is
170 replace the author, title, abstract, and text of the paper with your own.

171 The formatting instructions contained in these style files are summarized in Sections 6, 7, and 8
172 below.

173 6 General formatting instructions

174 The text must be confined within a rectangle 5.5 inches (33 picas) wide and 9 inches (54 picas) long.
175 The left margin is 1.5 inch (9 picas). Use 10 point type with a vertical spacing (leading) of 11 points.

176 Times New Roman is the preferred typeface throughout, and will be selected for you by default.
177 Paragraphs are separated by $\frac{1}{2}$ line space (5.5 points), with no indentation.

178 The paper title should be 17 point, initial caps/lower case, bold, centered between two horizontal
179 rules. The top rule should be 4 points thick and the bottom rule should be 1 point thick. Allow $\frac{1}{4}$ inch
180 space above and below the title to rules. All pages should start at 1 inch (6 picas) from the top of the
181 page.

182 For the final version, authors' names are set in boldface, and each name is centered above the
183 corresponding address. The lead author's name is to be listed first (left-most), and the co-authors'
184 names (if different address) are set to follow. If there is only one co-author, list both author and
185 co-author side by side.

186 Please pay special attention to the instructions in Section 8 regarding figures, tables, acknowledgments,
187 and references.

188 **7 Headings: first level**

189 All headings should be lower case (except for first word and proper nouns), flush left, and bold.

190 First-level headings should be in 12-point type.

191 **7.1 Headings: second level**

192 Second-level headings should be in 10-point type.

193 **7.1.1 Headings: third level**

194 Third-level headings should be in 10-point type.

195 **Paragraphs** There is also a `\paragraph` command available, which sets the heading in bold, flush
196 left, and inline with the text, with the heading followed by 1 em of space.

197 **8 Citations, figures, tables, references**

198 These instructions apply to everyone.

199 **8.1 Citations within the text**

200 The `natbib` package will be loaded for you by default. Citations may be author/year or numeric, as
201 long as you maintain internal consistency. As to the format of the references themselves, any style is
202 acceptable as long as it is used consistently.

203 The documentation for `natbib` may be found at

204 `http://mirrors.ctan.org/macros/latex/contrib/natbib/natnotes.pdf`

205 Of note is the command `\citet`, which produces citations appropriate for use in inline text. For
206 example,

207 `\citet{hasselmo}` investigated\dots

208 produces

209 Hasselmo, et al. (1995) investigated...

210 If you wish to load the `natbib` package with options, you may add the following before loading the
211 `neurips_2020` package:

212 `\PassOptionsToPackage{options}{natbib}`

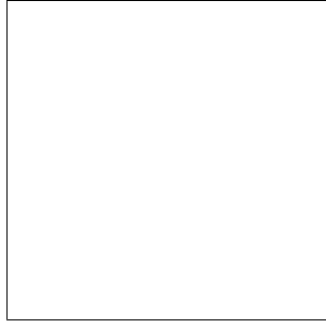


Figure 1: Sample figure caption.

213 If `natbib` clashes with another package you load, you can add the optional argument `nonatbib`
214 when loading the style file:

```
215 \usepackage[nonatbib]{neurips_2020}
```

216 As submission is double blind, refer to your own published work in the third person. That is, use “In
217 the previous work of Jones et al. [4],” not “In our previous work [4].” If you cite your other papers
218 that are not widely available (e.g., a journal paper under review), use anonymous author names in the
219 citation, e.g., an author of the form “A. Anonymous.”

220 8.2 Footnotes

221 Footnotes should be used sparingly. If you do require a footnote, indicate footnotes with a number¹
222 in the text. Place the footnotes at the bottom of the page on which they appear. Precede the footnote
223 with a horizontal rule of 2 inches (12 picas).

224 Note that footnotes are properly typeset *after* punctuation marks.²

225 8.3 Figures

226 All artwork must be neat, clean, and legible. Lines should be dark enough for purposes of reproduction.
227 The figure number and caption always appear after the figure. Place one line space before the figure
228 caption and one line space after the figure. The figure caption should be lower case (except for first
229 word and proper nouns); figures are numbered consecutively.

230 You may use color figures. However, it is best for the figure captions and the paper body to be legible
231 if the paper is printed in either black/white or in color.

232 8.4 Tables

233 All tables must be centered, neat, clean and legible. The table number and title always appear before
234 the table. See Table 1.

235 Place one line space before the table title, one line space after the table title, and one line space after
236 the table. The table title must be lower case (except for first word and proper nouns); tables are
237 numbered consecutively.

238 Note that publication-quality tables *do not contain vertical rules*. We strongly suggest the use of the
239 `booktabs` package, which allows for typesetting high-quality, professional tables:

```
240 https://www.ctan.org/pkg/booktabs
```

241 This package was used to typeset Table 1.

¹Sample of the first footnote.

²As in this example.

Table 1: Sample table title

Part		
Name	Description	Size (μm)
Dendrite	Input terminal	~ 100
Axon	Output terminal	~ 10
Soma	Cell body	up to 10^6

9 Final instructions

Do not change any aspects of the formatting parameters in the style files. In particular, do not modify the width or length of the rectangle the text should fit into, and do not change font sizes (except perhaps in the **References** section; see below). Please note that pages should be numbered.

10 Preparing PDF files

Please prepare submission files with paper size “US Letter,” and not, for example, “A4.”

Fonts were the main cause of problems in the past years. Your PDF file must only contain Type 1 or Embedded TrueType fonts. Here are a few instructions to achieve this.

- You should directly generate PDF files using `pdflatex`.
- You can check which fonts a PDF files uses. In Acrobat Reader, select the menu Files>Document Properties>Fonts and select Show All Fonts. You can also use the program `pdf fonts` which comes with `xpdf` and is available out-of-the-box on most Linux machines.
- The IEEE has recommendations for generating PDF files whose fonts are also acceptable for NeurIPS. Please see <http://www.emfield.org/icuwb2010/downloads/IEEE-PDF-SpecV32.pdf>
- `xfig` “patterned” shapes are implemented with bitmap fonts. Use “solid” shapes instead.
- The `\bbold` package almost always uses bitmap fonts. You should use the equivalent AMS Fonts:

```
\usepackage{amsfonts}
```

followed by, e.g., `\mathbb{R}`, `\mathbb{N}`, or `\mathbb{C}` for \mathbb{R} , \mathbb{N} or \mathbb{C} . You can also use the following workaround for reals, natural and complex:

```
\newcommand{\RR}{\mathbb{R}} %real numbers
\newcommand{\Nat}{\mathbb{N}} %natural numbers
\newcommand{\CC}{\mathbb{C}} %complex numbers
```

Note that `amsfonts` is automatically loaded by the `amssymb` package.

If your file contains type 3 fonts or non embedded TrueType fonts, we will ask you to fix it.

10.1 Margins in L^AT_EX

Most of the margin problems come from figures positioned by hand using `\special` or other commands. We suggest using the command `\includegraphics` from the `graphicx` package. Always specify the figure width as a multiple of the line width as in the example below:

```
\usepackage[pdftex]{graphicx} ...
\includegraphics[width=0.8\linewidth]{myfile.pdf}
```

See Section 4.4 in the `graphics` bundle documentation (<http://mirrors.ctan.org/macros/latex/required/graphics/grfguide.pdf>)

A number of width problems arise when L^AT_EX cannot properly hyphenate a line. Please give LaTeX hyphenation hints using the `\-` command when necessary.

Broader Impact

The contribution of our research is three-fold:

- We created and published a data set of how-to videos with time-tagged scripts, machine-generated summaries
- We generalized existing text summarization models to the scripts extracted from the videos [Sanabria et.al.]
- We augmented ROUGE metrics [Chin-Yew Lin] for evaluation of the results with a framework for formalized expert assessment based on our research and criteria proposed by previous works

At a high level, we hope that our analysis of transferability of summarization techniques from text to videos will have both practical and theoretical impacts by helping identify promising directions for future research.

References

References follow the acknowledgments. Use unnumbered first-level heading for the references. Any choice of citation style is acceptable as long as you are consistent. It is permissible to reduce the font size to small (9 point) when listing the references. **Note that the Reference section does not count towards the eight pages of content that are allowed.**

@article{DBLP:journals/corr/abs-1810-09305, author = Mahnaz Koupaee and William Yang Wang, title = WikiHow: A Large Scale Text Summarization Dataset, journal = CoRR, volume = abs/1810.09305, year = 2018, url = <http://arxiv.org/abs/1810.09305>, archivePrefix = arXiv, eprint = 1810.09305, timestamp = Wed, 31 Oct 2018 14:24:29 +0100, biburl = <https://dblp.org/rec/journals/corr/abs-1810-09305.bib>, bibsource = dblp computer science bibliography, <https://dblp.org>

[1] Yang Liu, Mirella Lapata. Text Summarization with Pretrained Encoders. (2019) URL. <https://arxiv.org/abs/1908.08345v2>

[2] 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)*, pages 1073-1083.

[3] Ilya Sutskever, Oriol Vinyals, and Quoc V Le. Sequence to sequence learning with neural networks. *Neural Information Processing Systems*, 2014.

[4] Haoran Li, Junnan Zhu, Cong Ma, Jiajun Zhang, and Chengqing Zong. (2017). Multi-modal summarization for asynchronous collection of text, image, audio and video. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 1092-1102. Association for Computational Linguistics.

[5] Sanabria, R., Caglayan, O., Palaskar, S., Elliott, D., Barrault, L., Specia, L., and Metze, F. How2: A large-scale dataset for multimodal language understanding. *CoRR*, abs/1811.00347, 2018. URL. <https://arxiv.org/abs/1811.00347>

[6] Nenkova, A. (2005). Automatic text summarization of newswire: Lessons learned from the document understanding conference. In *Proceedings of AAAI 2005*, Pittsburgh, USA.

[7] Svore, K., Vanderwende, L., and Burges, C. (2007). Enhancing single-document summarization by combining RankNet and third-party sources. In *Proceedings of the EMNLP-CoNLL*, pages 448-457. [7, 8]

[8] Yu-Hsiang Huang. Attention is all you need - pytorch. <https://github.com/jadore801120/attention-is-all-you-need-pytorch>, 2018.

[9] Nima Sanjabi. Abstractive text summarization with attention-based mechanism. Master's thesis, Universitat Politècnica de Catalunya, July 2018.

[10] Berna Erol, D-S Lee, and Jonathan Hull. 2003. Multimodal summarization of meeting recordings. In *Multimedia and Expo, 2003. ICME'03. Proceedings. 2003 International Conference on*, volume 3, pages III-25. IEEE.

- 325 [11] Dian Tjondronegoro, Xiaohui Tao, Johannes Sasongko, and Cher Han Lau. 2011. Multi-modal summariza-
 326 tion of key events and top players in sports tournament videos. In Applications of Computer Vision (WACV),
 327 2011 IEEE Workshop on, pages 471–478. IEEE
- 328 [12] Alexander M. Rush, Sumit Chopra, and Jason Weston. 2015. A neural attention model for abstractive
 329 sentence summarization. In Proceedings of the 2015 Conference on Empirical Methods in Natural Language
 330 Processing, pages 379–389. Association for Computational Linguistics.
- 331 [13] Shruti Palaskar, Jindrich Libovicky, Spandana Gella, Florian Metze. 2019. Multimodal Abstractive
 332 Summarization for How2 Videos. In Proceedings of the 57th Annual Meeting of the Association for Computational
 333 Linguistics, pages 6587–6596. Association for Computational Linguistics.
- 334 [14] Antoine Miech, Dimitri Zhukov, Jean-Baptiste Alayrac, Makarand Tapaswi, Ivan Laptev, Josef Sivic. 2019.
 335 HowTo100M: Learning a Text-Video Embedding by Watching Hundred Million Narrated Video Clips. In ICCV
 336 2019. <https://arxiv.org/abs/1906.03327>.
- 337 [1] Alexander, J.A. & Mozer, M.C. (1995) Template-based algorithms for connectionist rule extraction. In
 338 G. Tesauero, D.S. Touretzky and T.K. Leen (eds.), *Advances in Neural Information Processing Systems 7*, pp.
 339 609–616. Cambridge, MA: MIT Press.
- 340 [2] Bower, J.M. & Beeman, D. (1995) *The Book of GENESIS: Exploring Realistic Neural Models with the*
 341 *GENeral NEural Simulation System*. New York: TELOS/Springer-Verlag.
- 342 [3] Hasselmo, M.E., Schnell, E. & Barkai, E. (1995) Dynamics of learning and recall at excitatory recurrent
 343 synapses and cholinergic modulation in rat hippocampal region CA3. *Journal of Neuroscience* **15**(7):5249-5262.