

---

# Summarization of instructional video transcripts using BERT

---

Alexandra Savelieva\*

Bryan Au-Yeung\*

Vasanth Ramani\*

\*equal contribution

## Abstract

In this paper, we study abstractive summarization among a variety of “How-to” instructional videos and various written texts. Unlike traditional video summarization which focuses on condensing select video frames, our work uses step by step learning from a combination of news stories, WikiHow articles, and video transcripts. We showcase how a top performing document-level encoder based on BERT can boost the fluency and generalizability of summaries across a wide variety of instructional text and videos. In addition to our fine tuning and order preserving training methods, we present a novel dataset with over 5,000 transcripts extracted and constructed from open-domain videos from YouTube and the HowTo100Million Dataset. Our video dataset spans a variety of categories and is highly diverse in length and style. We demonstrate that our model is highly generalizable and produces summaries comparable to human written texts. To score the semantic adequacy of our abstract summaries, we use Content F1, Meteor, and human evaluations.

## 1 Introduction

Demand for generating keywords and descriptions for user-generated instructional videos has increased as online platforms boost video marketing. Currently, many descriptions for video content are human written and configured to maximize results through search engine optimization. However, descriptions do not always provide clear information on video content and sometimes fail to capture the most important parts of the video. Our research model abstracts from video transcripts to create short descriptions with improved semantic qualities to enhance user engagement and experience.

In our work, we trained a BERT-based model that generates summaries from ASR (speech-to-text) scripts of comparable quality to human written descriptions of instructional YouTube videos. The rest of this paper is divided in the following sections.

- We have curated and publish an aggregated data set of text and summaries formatted out of WikiHow articles, How2 videos, and CNN/Daily Mail stories;
- We follow by fine tuning existing BERT-based text summarization models to make them applicable to auto-generated scripts from instructional videos;
- Improve upon automated metrics [Chin-Yew Lin] used in the evaluation of summaries using a semantic assessment framework proposed by previous research.

## 2 Prior work

Prior to 2014, summarization was centered on extracting lines from single documents using statistical models and neural networks had limited success [6, 7]. The work on sequence to sequence models

from Sutskever et al. and Cho et al. opened up new possibilities for neural networks in natural language processing. From 2014 to 2015, LSTMs (a variety of RNN) became the dominant approach that achieved state of the art results. They became successful in tasks such as speech recognition, machine translation, parsing, image captioning, etc. This paved the way for abstractive summarization, which began to score competitively against extractive summarization. In 2017, Attention is all you need [8] provided a solution to the ‘fixed length vector’ problem, enabling neural networks to focus on important parts of the input for prediction tasks. Transformers with attention became more dominant for certain tasks [9].

In abstractive video summarization, models which incorporate variations of LSTM and deep layered neural networks have become state of the art performers. However, generating summaries from conversational texts in videos are still difficult. The deficiency of human annotated data has limited the amount of benchmarked datasets available for such research. Additionally, most work in the field of video summarization has traditionally focused on the isolation and concatenation of important video frames using natural language processing techniques. There are often inconsistencies and stylistic changes in spoken language that are difficult to translate into written text. In this work, we approach video summarizations by extending top performing single-document text summarization models [12] to narrated instructional videos.

### 3 Methodology

Initial exploratory data analysis indicates that post processing improves summarization performance. For example, in one of the sample videos in our test data set closed captioning confuses the speaker’s words “*how you get a text from a YouTube video*” for “*how you get attacks from a YouTube video*”. Our work iterates through the process described in the following sections.

#### 3.1 Data Collection

We hypothesized that the more labeled summarization data we bring, the more our model will benefit in the training process in terms of generalizability. The datasets illustrate different summary styles that range from single sentence phrases to short paragraphs.

- **CNN/Daily Mail dataset** [?]: CNN and Daily Mail includes a combination of news articles and story highlights written with an average length of 119 words per article and 83 words per summary.
- **Wikihow dataset**: a large scale text dataset containing over 200,000 single document summaries. Wikihow is a consolidated set of recent ‘How To’ instructional texts compiled from wikihow.com, ranging from topics such as ‘How to deal with coronavirus anxiety’ to ‘How to play Uno.’ These articles vary in size and topic but are structured to drive instructions across to the user. The first sentences of each paragraph within the article are concatenated to form a summaries.
- **How2 Dataset**: This YouTube compilation has videos (8,000 videos - approximately 2,000 hours) averaging 90 seconds long and 291 words in transcript length. It includes human written summaries which video owners were instructed to write summaries to maximize the audience. Summaries are two to three sentences in length with an average length of 33 words.

Despite the development of instructional datasets such as Wikihow and How2, advancements in summarization have been limited by the availability of human annotated transcripts and summaries. Such datasets are difficult to obtain and expensive to create, often resulting in repetitive usage of singular-task and highly structured data. As seen in the How2 dataset, videos with a certain length and structured summary are used for training and testing.

We introduce a new dataset, obtained from several How To and Do-It-Yourself YouTube playlists and video sampling from the published HowTo100Million Dataset. The HowTo100Million Dataset is a large scale dataset of over 100 million video clips taken from narrated instructional videos across 140 categories. Our dataset incorporates a sample across all categories and utilizes the natural language annotations from automatically transcribed narrations provided by YouTube.

Table 1: DataSet details

Dataset Size	5,195 videos (Youtube: 1,809. HowTo100Million: 3,386)
YouTube Min / Max Length	4 / 1,940 words
YouTube Average Length	259 words
HowTo100Million Sample Min / Max Length	5 / 6,587 words
HowTo100Million Sample Average Length	859 words

### 3.2 Preprocessing

Due to diversity and complexity of our input data, we built a preprocessing pipeline for aligning the data to a common format. We observed issues with lack of punctuation, incorrect wording, and extraneous introductions which impacted model training. With these challenges, our model misinterpreted text segment boundaries and produces poor quality summaries. In exceptional cases, the model failed to produce any summary. In order to maintain the fluency and coherency in human written summaries, we cleaned and restored sentence structure using Spacy as shown in the figure below. 1.

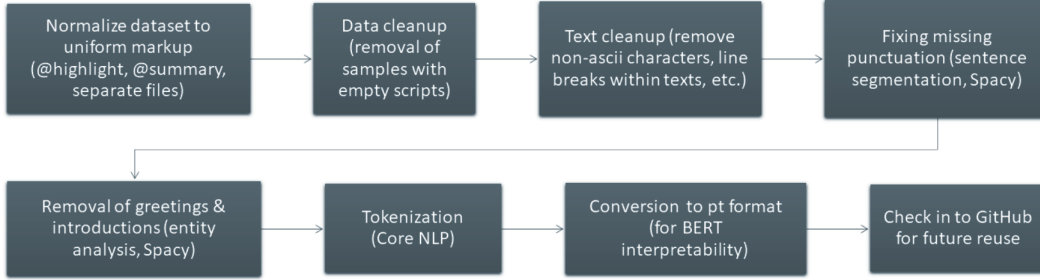


Figure 1: Preprocessing.

The differences in conversational style of the video scripts and transcribed news stories (on which the models were pretrained) impacted the quality of the model output. In our initial application of the extractive summarization model pretrained on CNN/DM dataset, stylistic errors manifested in a distinct way. The model considered initial introductory sentences to be important in generating summaries (this phenomena is referred to by [15] as N-lead, where N is the number of important first sentences). Our model generated short, simple worded summaries such as "hi!" and "hello, this is <first and last name>". In order to improve the quality and score of our model output, we used entity detection spacy and nltk to remove introductions from the CNN/Daily Mail input to our summarization model. We did not anonymize entities in our datasets. We split sentences and tokenized using the Stanford Core NLP toolkit on all datasets and preprocessed the data in the same method used by (See et. al.).

### 3.3 Summarization models

We utilized the BertSum models created by Yang et al[1] for our research. This includes both Extractive and Abstractive summarization models, which employs a documentation level encoder based on Bert. Novelty in transformer architecture comes from a pretrained BERT encoder with a randomly initialized Transformer decoder. The model uses two different learning rates: a low rate for the encoder and a separate higher rate for the decoder to enhance learning.

We used a 4-GPU Linux machine and established a baseline by training our extractive model on 5,000 video samples from the How2 dataset. We used the BERT base uncased with 10,000 steps, which took 12 hours to train. We fine tuned the summarization model inclusive of the BERT layer, testing various epoch sizes. Sequentially, we added the full CNN/Dailymail, entire how2 dataset and 3,097 samples from Wikihow with a 50,000 step size to the training set to get better results.

Finally, we used the Abstractive summarization model and our aggregated dataset of CNN/Dailymail, Wikihow, and How2 datasets with a total of 535,527 examples and trained for 210,000 steps with a training batch size of 50 and 20 epochs. By controlling the order of datasets in which we trained our model, we were able to improve the fluency and length of our summaries. Due to the limitations of computer processing and storage resources, our final model took four days to train. As stated in previous research, the original model contained more than 180 million parameters and used two Adam optimizers with  $\beta_1=0.9$  and  $\beta_2=0.999$  for the encoder and decoder respectively. The encoder used a learning rate of 0.002 and the decoder had a learning rate of 0.2 to ensure that the encoder was trained with more accurate gradients while the decoder became stable.

### 3.4 Scoring of results

Results were scored using ROUGE. While we expected a correlation between good summaries and high ROUGE scores, we observed examples of poor summaries with high scores, such as in Figure 5, and good summaries with low ROUGE scores. To maintain the semantic quality of summaries, we deferred evaluation to human experts.

Additionally, we added Content F1 scoring, a metric proposed by Carnegie Mellon University to focus on the relevance of content. Similar to ROUGE, Content F1 scores summaries with a weighted f-score and a penalty for incorrect word order. It also discounts stop and buzz words that frequently occur in the how-to domain, such as “learn from experts how to in this free online video”.

In addition to automatically calculated scores, it is important to have human judges review the results. We created a framework of criteria for evaluation using Python, Google Forms, and Excel spreadsheets. Summaries for the surveys are randomly sampled from our dataset to avoid biases. In order to avoid asymmetrical information between human versus machine generated summaries, we removed capitalized text. We used two types of questions: A Turing test question for participants to distinguish AI from human-generated descriptions. The second involves selecting quality ratings for summaries. Below are definitions of criteria for clarity:

- Fluency: Does the text have a natural flow and rhythm?
- Usefulness: Does it have enough information to make a user decide whether they want to spend time watching the video?
- Succinctness: Does the text look concise or does it have redundancy?
- Consistency: Are there any ambiguous, confusing or contradicting statements in the text?
- Realisticity: Is there anything that seems far-fetched and bizarre in words combinations, or do the statements look "normal"?

Options for grading of results are as follows: 1: Bad 2: Below Average 3: Average 4: Good 5: Great.

## 4 Experiments and Results

### 4.1 Training

Our baseline results were obtained from applying the state-of-the art extractive BertSum model pretrained on CNN/DailyMail. We hypothesized achieving high performing results with BERT, but received low scores. Summaries generated from the model were incoherent, repetitive, and uninformative. Despite poor performance, the model performed better in the health subdomain within how-to videos. We explained this as a symptom of heavy coverage in news reports generated by CNN/Daily Mail. We realized that extractive summarization is not the strongest model for our goal: most youtube videos are presented with a casual conversational style, while summaries have higher formality. We pivoted to abstractive summarization to improve performance.

In order to create a generalizable abstractive model, we first trained on a large corpus of news. This allowed our model to understand structured texts. We then introduced Wikihow, which exposes the model to the how-to domain. Finally, we train and validate on the how-to dataset, narrowing the focus of the model to a selectively structured format.

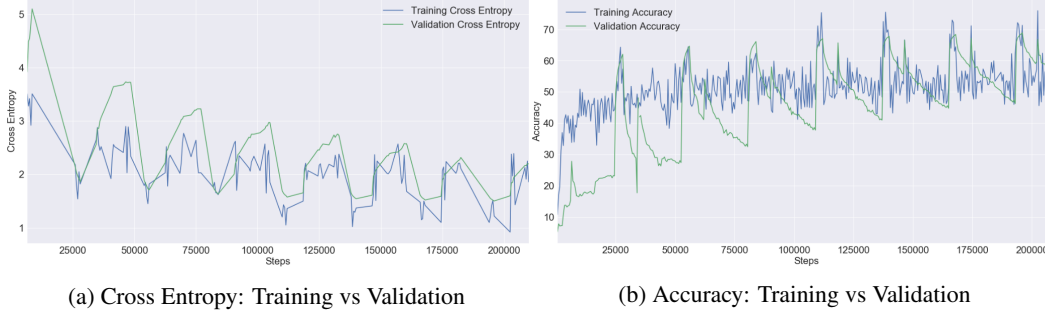


Figure 2: BertSum Abstractive Summarization: Model Performance

The cross entropy chart in the Figure 2a shows that the model is neither overfitting nor underfitting the training data. Good fit is indicated with the convergence of training and validation lines. Figure 2b shows the model’s accuracy metric on the training and validation sets. The model is validated using the How2 dataset against the training dataset. The model improves as expected with more steps.

## 4.2 Evaluation

The BertSum model created by Yang trained on CNN and Daily Mail [Yang] resulted in state of the art scores when applied to samples from those datasets. However, when tested on our How2 Test dataset, it gave very poor performance and a lack of generalization in the model (see Table 2). Looking at the data, we found that the model tends to pick the first one or two sentences for the summary. We hypothesized that removing introductions from the text would help improve ROUGE scores. Our model improved a few ROUGE points after applying preprocessing described in the Section 3.2 above. Another improvement came from adding word deduping to the output of the model, as we observed it occurring on rare words which are unfamiliar to the model. We still did not achieve scores get higher than 22.5 ROUGE-1 F1 and 20 ROUGE-L F1 (initial scores achieved from training with only the CNN/Daily Mail dataset and tested on How2 data). Reviewing scores and texts of individual summaries showed that the model performed better on some topics such as medicine, while scoring lower on others, such as sports. Again, this makes sense for a model that is trained on news: it isn’t reasonable for the model to perform well on esoteric yoga terminology. In our next series of experiments, we used our own dataset for training. Even though the difference in ROUGE scores for the results on [1-3] are not drastically different from [4-5], the quality of summaries from the perspective of human judges is qualitatively different.

Current best results were accomplished by leveraging the full set of labeled datasets (CNN/DM, WikiHow, and How2 videos) with order preserving configuration. Order is very important: as human learner, the model isn’t able to make substantial progress if it switches contexts between tasks of different complexity. The easiest training (CNN/DM) needs to be done first. We then move on to the next step of learning to summarize from WikiHow, which covers more domains and has more complicated, but predictable structure. Only after training textual scripts do we proceed to video scripts, which presents additional challenges of ad-hoc flow and conversational language. We did not see a large impact from spelling errors that frequently occur in ASR-generated scripts without human supervision, but ensuring correct boundaris between sentences by using Spacy to fix punctuation errors made a big difference. Our results for videos have reached the level of the best scores for news [1].

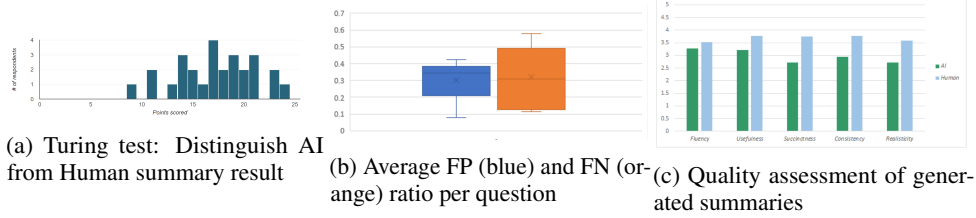


Figure 3: Human evaluation of model-generated summaries in comparison with real video descriptions from YouTube

Table 2: Comparison of results

Experiment		Rouge-1	Rouge-L	Content-F1
Model	Pretraining Data			
1. BertSum	CNN and Daily Mail	18.08	18.01	26.0
2. BertSum with preprocessing	CNN and Daily Mail	20.51	18.86	-
3. BertSum with pre- and postprocessing	CNN and Daily Mail	22.47	20.07	-
4. BertSum	How-To, WikiHow, CNN and Daily Mail	24.4	21.45	-
5. BertSum with postprocessing	How-To, WikiHow, CNN and Daily Mail	26.32	22.47	-
6. BertSum with no shuffling and more training data	How-To, WikiHow, CNN and Daily Mail	48.26	44.02	36.4
7. BertSum[Yang]	CNN/DailyMail	43.23	39.63	-
8. Model [13]	How-To	59.3	59.2	48.9
9. Pointer Generator+Coverage	WikiHow	28.53	26.54	-
10. Lead 3	CNN/DM	40.34	36.57	-

From our initial incomprehensive results, we achieved fluent and understandable video descriptions which give a clear idea about the content. Our scores did not surpass scores from other researchers despite employing BERT. However, when looking at comparisn of the texts, our summaries appear to be more fluent and useful in content for users looking at summaries in the how-to domain. Some examples are given in the appendix:

Based on these observations, we decided that the model performed strongly enough for us move to the final stages. We leveraged independent experts and evaluated the quality of our summaries in comparison to descriptions that users provide for their videos on Youtube. We recruited a diverse group of 30+ volunteers (31 have responded at the time of writing this paper) to blindly evaluate a set of 25 randomly selected video summaries that were generated by our model and descriptions of videos from our curated dataset. We had two types of questions: one, a version of famous Turing test, was a challenge to distinguish AI from human-curated descriptions and used the framework described in Section 3.4. The aggregated results for both evaluations are in Figures 3a - 3c. We observe zero perfect scores on Turing test answers. Results included many false positives and false negatives. The quality of our model output is comparable to average YouTube summaries. As expected, the fluency of our summaries is comparable to human-curated text. "Realistic" text is the main growth opportunity, because the abstractive model is prone to generating incoherent sentences that are grammatically correct.

## 5 Conclusion

The contributions of our research are addressing multiple issues that we identified in pursuit of generalizing BertSum model for summarization of instructional video scripts throughout the whole training process.

- We complemented existing labeled summarization datasets with autogenerated instructional video scripts and human-curated descriptions.
- We explored how different combinations of training data and parameters impact the training performance of BertSum abstractive summarization model.
- We came up with novel preprocessing steps for auto-generated closed captioning scripts before summarization.
- We generalized BertSum abstractive summarization model to autogenerated instructional video scripts with the quality level that is close to randomly sampled descriptions created by Youtube users.
- We designed and implemented a new framework for blind unbiased review that produces more actionable and objective scores, augmenting ROUGE, BLEU and Content F1.

All the artifacts listed above are available in to our repository for the benefit of future researchers <sup>1</sup>. Overall, the results we obtained by now on amateur narrated instructional videos make us believe that we were able to come up with a trained model that generates summaries from ASR (speech-to-text) scripts of competitive quality to human-curated descriptions on YouTube.

## References

- [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.

---

<sup>1</sup><https://github.com/alebryvas/berk266/> - it’s not a public repository yet, but we can provide access upon request

- [11] Dian Tjondronegoro, Xiaohui Tao, Johannes Sasongko, and Cher Han Lau. 2011. Multi-modal summarization of key events and top players in sports tournament videos. In Applications of Computer Vision (WACV), 2011 IEEE Workshop on, pages 471–478. IEEE
- [12] Alexander M. Rush, Sumit Chopra, and Jason Weston. 2015. A neural attention model for abstractive sentence summarization. In Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing, pages 379–389. Association for Computational Linguistics.
- [13] Shruti Palaskar, Jindrich Libovicky, Spandana Gella, Florian Metze. 2019. Multimodal Abstractive Summarization for How2 Videos. In Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, pages 6587–6596. Association for Computational Linguistics.
- [14] Antoine Miech, Dimitri Zhukov, Jean-Baptiste Alayrac, Makarand Tapaswi, Ivan Laptev, Josef Sivic. 2019. HowTo100M: Learning a Text-Video Embedding by Watching Hundred Million Narrated Video Clips. In ICCV 2019. <https://arxiv.org/abs/1906.03327>.
- [15] Wikihow: A large scale text summarization dataset. M Koupaei, WY Wang. arXiv preprint arXiv:1810.09305, 2018. 13, 2018 =====

## 6 Appendix

### 6.1 Model details

Extractive summarization is generally a binary classification task with labels indicating whether sentences should be included in the summary. Abstractive summarization, on the other hand, requires language generation capabilities to create summaries containing novel words and phrases not found in the source text.

The architecture in the Figure 4 shows the BERTSUM model. It uses a novel documentation level encoder based on BERT which can encode a document and obtain representation for the sentences. CLS token is added to every sentence instead of just 1 CLS token in the original BERT model. Abstractive model uses an encoder-decoder architecture, combining the same pretrained BERT encoder with a randomly initialized Transformer decoder. The model uses a special technique where the encoder portion is almost kept same with a very low learning rate and a separate learning rate is used for the decoder to make it learn better.

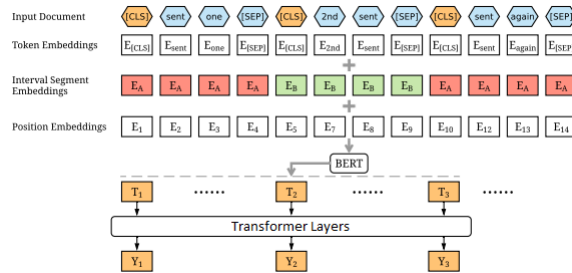


Figure 4: BERTSUM Architecture. From [Yang et. al.]



## 6.2 Illustrated Example of why ROUGE metrics is not sufficient

```
*****
Reference: now that you have spent the time cleaning your oven learn how to keep it clean with expert tips in this free h
ow to video on how to better clean your oven

Hypothesis: make sure your oven is clean .<q>clean your oven .<q>make sure you want to clean the oven with a towel .<q>ge
t your food .<q>put your food in your baking soda and water .<q>do n't go to the kitchen .

rouge-1:      P: 29.55      R: 40.62      F1: 34.21
rouge-2:      P:  6.98      R:  9.68      F1:  8.11
rouge-3:      P:  2.38      R:  3.33      F1:  2.78
rouge-4:      P:  0.00      R:  0.00      F1:  0.00
rouge-l:      P: 24.16      R: 31.50      F1: 27.34
rouge-w:      P: 14.23      R:  9.78      F1: 11.59
*****
```

Figure 5: An example where ROUGE metric is confusing.

## 6.3 Examples of Comparison of our model output vs Benchmark and reference summaries

- Summary 1: growing rudbeckia requires full hot sun and good drainage. grow rudbeckia with tips from a gardening specialist in this free video on plant and flower care. care for rudbeckia with gardening tips from an experienced gardener
- Benchmark 1: growing black - eyed - susan is easy with these tips, get expert gardening tips in this free gardening video .
- Reference 1: growing rudbeckia plants requires a good deal of hot sun and plenty of good drainage for water . start a rudbeckia plant in the winter or anytime of year with advice from a gardening specialist in this free video on plant and flower care
- Summary 2: camouflage thick arms by wearing sleeves that are not close to the arms and that have a line that goes all the way to the waist. avoid wearing jackets and jackets with tips from an image consultant in this free video on fashion. learn how to dress for fashion modeling
- Benchmark 2: hide thick arms and arms by wearing clothes that hold the arms in the top of the arm. avoid damaging the arm and avoid damaging the arms with tips from an image consultant in this free video on fashion .
- Reference 2: hide thick arms by wearing clothes sleeves that almost reach the waist to camouflage the area .conceal the thickness at the top of the arms with tips from an image consultant in this free video on fashion.