# Summarization of Multimodal Presentations with Vision-Language Models: Study of the Effect of Modalities and Structure

THÉO GIGANT, Université Paris-Saclay, CNRS, CentraleSupelec, Laboratoire des signaux et systemes, France

CAMILLE GUINAUDEAU, Université Paris-Saclay, CNRS, LISN, France

FRÉDÉRIC DUFAUX, Université Paris-Saclay, CNRS, CentraleSupelec, Laboratoire des signaux et systemes, France

Vision-Language Models (VLMs) can process visual and textual information in multiple formats: texts, images, interleaved texts and images, or even hour-long videos. In this work, we conduct fine-grained quantitative and qualitative analyses of automatic summarization of multimodal presentations using VLMs with various representations as input. From these experiments, we suggest cost-effective strategies for generating summaries from text-heavy multimodal documents under different input-length budgets using VLMs. We show that slides extracted from the video stream can be beneficially used as input against the raw video, and that a structured representation from interleaved slides and transcript provides the best performance. Finally, we reflect and comment on the nature of cross-modal interactions in multimodal presentations and share suggestions to improve the capabilities of VLMs to understand documents of this nature.

## 1 Introduction

Presentations play a central role in academic, business, and technical communication, serving as a widespread means of conveying information. By leveraging multimodal elements such as visual slides and speech, presentations offer a dynamic and engaging way to present ideas and data. Their sheer volume and the frequency with which individuals encounter them make summarization both a practical and a pressing need. Summaries of multimodal presentations allow users to efficiently extract key points, navigate the presentation, and improve retention, without requiring in-depth engagement with the full presentation.

Presentations are multimodal when they include a visual medium, typically a slideshow, in addition to the speech. Meetings, lectures, and conferences commonly follow this format. Multimodal presentation records usually consist of text-heavy video and audio streams. A transcript and key images that span the entire video can be extracted from these records [14]. Using these extracted modalities or the original modalities, one can construct representations of the presentation to be used by VLMs to perform tasks such as summarization. Examples of such unimodal and multimodal representations for a presentation are illustrated in Figure 1. The choice of representation will have different token lengths and result in different downstream performances and costs.

VLMs are usually trained and evaluated on short input, including text and a few images, or short videos [4, 29]. Previous works have highlighted the poor performance of these models with a longer input [43]. Although long context performance on synthetic benchmarks can be inherited from the underlying language model [48], they usually do not reflect performance in real-world scenarios [24]. The raw video representation can be too long and lie significantly outside of the training context size for these models, rendering summarization very challenging [41]. However, we can build a shorter interleaved image-text representation resembling the training data of some VLMs [4, 20, 21].

Authors' Contact Information: Théo Gigant, Université Paris-Saclay, CNRS, CentraleSupelec, Laboratoire des signaux et systemes, Gif-Sur-Yvette, France, theo.gigant@l2s.centralesupelec.fr; Camille Guinaudeau, Université Paris-Saclay, CNRS, LISN, Orsay, France; Frédéric Dufaux, Université Paris-Saclay, CNRS, CentraleSupelec, Laboratoire des signaux et systemes, Gif-Sur-Yvette, France.

Fig. 1. VLMs are able to process a multimodal presentation in various unimodal and multimodal representations.

In this paper, we examine how the choice of representation of multimodal presentations affects their associated cost and performance for VLM-based summarization. We specify different representations for the model, using single or multiple modalities, including slides, text extracted from slides, the entire video, and the transcribed audio, optionally structured by using the time-alignment between these modalities. Furthermore, we perform a fine-grained analysis of the performance of Qwen2-VL [41] with these different input representations to understand the impact of the modalities and the structure information on the generated summaries. Additionally we benchmark multiple recent state-of-the-art open-weights VLMs using an interleaved slides-transcript representation.

The contributions of this paper can be summarized as follows:

- We benchmark open-weights VLMs for summarization of multimodal presentations using an interleaved slides-transcript input.
- We conduct a fine-grained quantitative and qualitative analysis of the cost and performance associated with many unimodal and multimodal input representations in Qwen2-VL.
- We suggest cost-effective representation strategies for summarization of multimodal presentations with VLMs.
- We discuss the challenges emerging from using multimodal representations in VLMs and share suggestions to mitigate the current shortcomings.

## 2 Related Work

### 2.1 Abstractive Summarization

Previous work has explored various approaches to abstractive summarization, *i.e.* efficiently condensing documents in short generated textual summaries. Recently, models based on the transformer encoder-decoder architecture [40] such as Bart [23] or T5 [37] have shown great performance for abstractive summarization of short textual documents.

These architectures have been adapted for long input, by interchanging the attention module with linear scaling alternatives such as local and global attentions [2, 7, 17, 45].

Current Large Language Models (LLMs), based on the transformer decoder architecture, perform multiple tasks, including abstractive summarization, with remarkable success [47]. Their capabilities even scale to a very large input due to efficient attention implementations [10, 11] and context extension methods [35].

However, these methods fall short for summarization of multimodal presentations because they neglect the visual modality, which is often critical to convey information in such contexts.

## 2.2 Vision-Language Models and Video-Language Models

Vision-Language Models (VLMs) [4] are early-fusion multimodal models comprised of a vision encoder [36, 46] connected to a pretrained LLM. They carry out the document understanding and instruction-following capabilities of LLMs to visual or multimodal documents [29].

The training data for recent VLMs include a significant amount of text-heavy images [21] such as document OCR or document visual question answering. VLMs can also be trained with interleaved text and images [4], meaning that they can process an arbitrary number of images in between text tokens.

Some VLMs can also process video input on top of images, optionally using a dedicated video encoder [26]. Video VLMs are usually trained with minutes-long *in-the-wild* videos picturing little to no text [32]. In contrast, the video stream of multimodal presentation records is typically long, redundant, and text-heavy. As such, summarization of multimodal presentations lies out of the training domain of video VLMs, and, by extracting key frames, can be made more similar to the domain for recent VLMs using the interleaved images-text format [4, 20, 21].

Wu et al. [43] highlighted the low performance of open-weights VLMs for long-context interleaved video-language understanding. Zhang et al. [48] suggested that VLMs can inherit long-context performance from their underlying LLMs, exemplifying on a synthetic benchmark. However, Li et al. [24] showed that the performance of LLMs on synthetic benchmarks does not reflect real-world performance.

## 2.3 Abstractive Summarization of text-heavy long multimodal documents

The summarization of long videos [3, 6], and the understanding of documents with a text-heavy visual modality [18] using VLMs have received significant attention in recent research. However, the intersection of these two domains remains relatively understudied. Some early efforts [14, 28] have benchmarked approaches using single-modality or out-of-domain baselines, revealing critical limitations of existing methods in effectively performing abstractive summarization of text-heavy long multimodal documents, such as those found in multimodal presentations.

TIB [14] is a dataset containing 9,103 multimodal scientific presentations of multiple domains, with abstracts as reference summaries. Each presentation was processed to extract the transcript and key frames. The transcripts were used to benchmark text-only abstractive summarization models, and key frames were provided to use with upcoming long-context VLMs.

VISTA [28] is a similar benchmark, consisting of 18,599 records of multimodal presentations from AI conferences, with their abstracts. The authors tested LLMs with transcript or OCR input, as well as VLMs with video input, and advocated for explicit planning to improve the quality and factual consistency of summaries.

Compared with these, our work constitutes a more comprehensive investigation into VLMs' ability to generate summaries from multimodal content, with a focus on the effect of input document representation.

## 3   Methodology

### 3.1   Dataset

We constitute a benchmark derived from the TIB dataset [14], initially comprised of 9103 records of scientific presentations, with transcribed speech, extracted key frames and author-provided abstract. We chose TIB over VISTA because of the availability of pre-processed key frames and transcripts, and the longer length of the presentations.

The original TIB dataset contains reference summaries that were author-provided, so the abstract style, quality, and content are varied. In order to create a higher-signal subset, we filter the top quality abstracts using a LLM-as-a-judge method. Based on the title and abstract, we prompt SmolLM2 [5] to score the abstracts from 1 to 9 and keep only the references scoring 9.

Gigant et al. [14] stated that the extracted key frames in the TIB dataset often contain screen-recorded or live-captured slideshows, but can also be a live recording of the speaker. To filter out the latter, we use SigLip [46] to perform zero-shot classification on key frames to identify those that contain a slide and those that do not. We discard the key frames that are not slides, and the presentations with less than 3 slides after this stage.

After processing, our benchmark consists of 822 presentations, including transcribed speech and slides with time codes. Using these time codes, we can create an interleaved sequence of slides followed by transcription of the simultaneous speech. For the rest of this work, we will refer to this format as "interleaved" or "structured". Otherwise, we default to a sequence that includes all slides, followed by the transcribed speech. We call this default representation "unstructured".

We share the filtered benchmark[1], as well as a tool to visualize the interleaved representation of the multimodal presentations[2] on the Hugging Face Hub.

### 3.2   Models

For a fine-grained analysis of the task at hand, we consider open-weights VLMs that fit on a single A-100 GPU, and are able to process multiple images, interleaved images-text format, videos and long-context input. The outputs are sampled with greedy decoding, with a maximum output size of 1024 tokens. For all models, we adapt the same prompt, using the following instructions:

```
Provide a summary of the presentation in a short paragraph, based on the
transcript and the slides. The summary should introduce the topic, and the speaker
and the presentation context if relevant. It will be written to be informative
about the content for a prospective attendee. {document}
```

The main model we are studying in this work is Qwen2-VL [41], as it checks all the requirements and also allows for arbitrary visual resolution. It also comes in different sizes; specifically, we work with the 2B and 7B instruction-tuned checkpoints. All of these characteristics make it convenient to rely on Qwen2-VL to perform a fine-grained analysis of the impact of the input modalities and structure on the generated summaries. In the benchmark, we include other VLMs capable of handling our multimodal presentations, namely Idefics-2 [21], Idefics-3 [20], and Phi-3.5-vision [1].

We believe that the findings of this work can be generalized to other models, due to the similarity in architecture and training data between current VLMs [21, 34].

### 3.3   Extractive Statistics

We use the transcripts provided in the TIB dataset, and extract words from the slides using an Optical Character Recognition (OCR) model: GOT-OCR2.0 [42]. This allows to compute word-based

---

[1]https://huggingface.co/datasets/gigant/tib-bench
[2]https://huggingface.co/spaces/gigant/slide-presentation-viz

statistics and metrics with respect to the speech by using the transcript, and to the slides by using the OCR, for modality-specific measures of extractiveness and relevance.
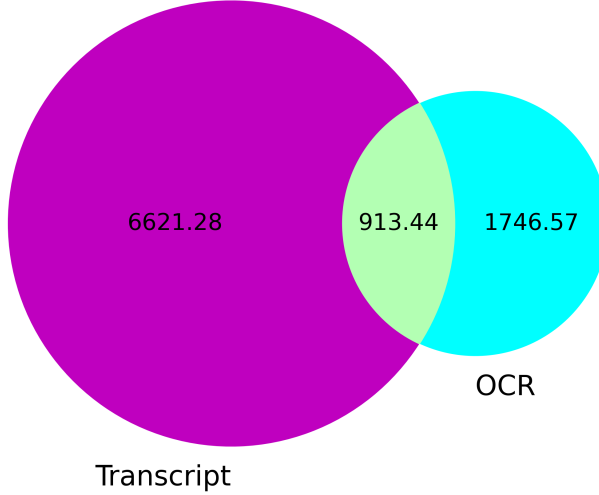


Fig. 2. Average token count in speech transcript and OCR, and their overlap

We compute the average token count in the transcripts, the OCR, and the vocabulary overlap between the two. These statistics are illustrated in Figure 2. We remind the reader that transcriptions from the images and the speech are imperfect and, as a consequence, the overlap might be inaccurate.

The extractive statistics of generated summaries are computed in multiple settings, with respect to the text extracted from the slides, the speech, or their concatenation, called `ocr`, `transcript` and `overall` respectively.

The extractive coverage refers to the vocabulary overlap between the source and the summary. The extractive density is the average length of the extractive fragments from the source in the summaries. A more extractive summary with respect to a source, can have a higher vocabulary overlap *ie* higher coverage, and/or longer extractive fragments, *ie* higher density.

### 3.4 Evaluation Metrics

Many metrics have been proposed for the fine-grained evaluation of summarization systems. However, meta-evaluations have shown that they usually scale poorly for long document summarization [19]. Metrics such as Bart Score [44] and ESTIME [39] rely on encoding the input document in a pretrained model, which has a limited input size, usually less than 1000 tokens. The documents in our experiments exceed on average more than ten times this input size, so we believe that these scores will not be reliable in this setting.

The Rouge score [27] is the standard metric for the evaluation of summaries. It gives an accurate estimate of the overall performance of a system. Rouge-$n$ counts the $n$-gram overlaps between the proposed summary and a reference summary. We will work primarily with unigram and bigram-based Rouge scores, denoted Rouge-1 ($R_1$) and Rouge-2 ($R_2$).

GRUEN ($G$) [51] is a reference-free metric intended to evaluate the linguistic quality of generated texts. It only uses the generated text as input, so it does not need to account for the input being long and multimodal.

The importance-based relevance score ($IbR$) [15] has been designed to provide a reference-free evaluation of the relevance of the summaries. The relevance of a summary is the measure of whether it contains the main ideas from the source document. As this metric relies on word importance scoring based on information retrieval methods, it scales well with long input and has been shown to correlate well with human-evaluated relevance of summaries on long-document summarization datasets. It uses the source as a textual reference and will have to be adapted for use with multimodal input. To this end, we are using extracted text from the speech and/or slides as sources. Similarly to the extractive statistics, we define three modes for this metric, in accordance to the source input, called `ocr`, `transcript` and `overall`. These input-modality-specific metrics allow us to conduct a more fine-grained analysis of the generated summaries with respect to the input modalities.

The Rouge scores ($R_1$ and $R_2$) provide us with an estimate of the overall performance of the summarization systems. GRUEN ($G$) scores the linguistic qualities of the systems. Finally, the multimodality-adapted importance-based relevance scores ($IbR_{transcript}$, $IbR_{ocr}$ and $IbR_{overall}$) give us an estimate of how well the model gathers the important information from the document, attributed to each input modality.

## 4 Benchmark

We benchmark some VLMs suited for the multimodal summarization task. The selected models are instantiated in the `bfloat16` format if possible; otherwise, we use 4-bit quantization and optionally truncate the input to 16k tokens to avoid out-of-memory errors.

The VLMs we are benchmarking are the following:

- Idefics-2 [21] is based on the Mistral-7B LLM, connected to the Siglip image encoder. It was trained with context length up to 2k tokens.
- Idefics-3 [20] is based on the Llama-3 LLM, connected to the Siglip image encoder. It was trained with a context length up to 10K tokens.
- Phi-3.5-vision [1] is based on the Phi-3.5-mini LLM, connected to the Clip image encoder. It was trained with a context length up to 128k tokens.
- Qwen2-VL [41] is based on the Qwen2 LLM, connected to a custom image encoder. It was trained with a context length up to 16k tokens.

We evaluate the summaries generated by these models with the collection of metrics described in Section 3.4. The results are reported in Table 1.

Idefics-2 has been trained with an interleaved image-text format that does not exceed 2k tokens. Even if it is technically capable of processing longer documents, we observe an inability to generate coherent summaries with significantly longer input. In fact, the GRUEN ($G$) score shows that the linguistic quality of the summaries generated by Idefics-2 is low, revealing the difficulty of this model in producing coherent text in such long context situations.

The other models are trained on longer context and all produce summaries of higher linguistic quality. The relevance of the generated summaries, as measured by the importance-based relevance ($IbR$) scores, highlights the differences between these models in the understanding of the presentations.

Table 1. Results of open-weight VLMs on TIB-bench

| Model | #Params | Visual token budget | $R_1$ ↑ | $R_2$ ↑ | $G$ ↑ | $IbR_{transcript}$ ↑ | $IbR_{ocr}$ ↑ | $IbR_{overall}$ ↑ |
|---|---|---|---|---|---|---|---|---|
| Idefics-2 [21] | 8B | 64 | 4.9 | 0.3 | 36.1 | 1.2 | 0.6 | 1.2 |
| Idefics-3<sup>q</sup> [20] | 8B | 169 | 12.4 | 1.9 | 48.0 | 21.2 | 7.9 | 16.3 |
| Phi-3.5-vision<sup>q†</sup> [1] | 4B | 144 | 17.1 | 2.9 | 51.6 | 29.7 | 8.9 | 20.2 |
| Qwen2-VL [41] | 2B | 512 | **27.1** | **5.4** | **57.1** | **53.8** | **17.2** | **33.4** |

<sup>q</sup> denotes models with weights quantized at 4-bits

<sup>†</sup> denotes models whose input was truncated to 16k tokens

From Idefics-2 to Idefics-3 multiple changes can explain the relative improvement in this task, including a longer context training, a more capable backbone LLM, and more visual tokens per image. Laurençon et al. [21] observed no gain from adding visual tokens per image after 64, while Laurençon et al. [20] observed that the OCR task benefited from more visual tokens per image. Similarly to OCR, the summarization of multimodal presentations is based on text-heavy images. In Section 5.2, we investigate the scaling of performance with more visual tokens per image to see if these observations translate to this task.

Even with a smaller parameter count, Phi-3.5-vision scores significantly better than Idefics-3, most likely due to its training data including videos and more long context examples. Between Idefics-3 and Phi-3.5-vision, the difference in performance seems to be linked to the language model more than the vision encoder, probably because of better long context understanding.

Qwen2-VL scores even better than Phi-3.5-vision, with a smaller parameter count, likely due to various improvements including a larger visual token budget and a larger and more diverse training dataset. Contrary to Phi-3.5-vision, it can process the whole non-truncated input with non-quantized weights.

The importance-based relevance (*IbR*) scores reveal that with these models, visual and textual understanding are improved jointly. A new VLM release usually relies simultaneously on better training data and a better backbone LLM, causing both visual and language understanding to improve.

For all these models, the task of summarization of multimodal presentations is most likely out-of-domain, which would benefit models such as Qwen2-VL, trained on a broad range of tasks, including analogous ones such as long document and video understanding.

In the following experiments, we are further exploring the capabilities of Qwen2-VL and using its video understanding capabilities and dynamic visual token budget to provide a fine-grained analysis of the impact and scaling of different visual and multimodal settings.

## 5 Fine-grained Analysis

### 5.1 Impact of the structure and modalities

Qwen2-VL [41] can process text, images, and/or videos as input. It includes a dynamic resolution mechanism that enables it to process images of varying resolution into different numbers of visual tokens. We experiment with different numbers of visual tokens, and call "visual token budget" the maximum number of visual tokens per image. The whole budget is used for square images of sufficient resolution; in other cases, the actual number of visual tokens used is smaller.

In the following experiments, visual token budget varies from 2 to 512. Given the various shapes and resolutions of the slides, the actual number of visual tokens per image is dynamic and ranges on average from 1 to 183. The measured averages are reported in Table 2. At 512 tokens per image,

we already exceed the maximum resolution of the extracted slides in the dataset, explaining the larger relative gap between the budget and the actual visual tokens used. The videos are kept at the original resolution, which is higher on average.

Table 2. Measured average number of visual tokens for different visual token budgets, and resulting input length for slides and video sampled at 0.1 or 1 fps. Settings where the input length exceeds the model context size are italicized.

| Visual token budget | Average visual tokens per slide | Total tokens (slides) | Total tokens (video@0.1fps) | Total tokens (video@1fps) |
| --- | --- | --- | --- | --- |
| 2 | 1.0 | 0.1k | 0.1k | 1k |
| 8 | 6.0 | 0.2k | 0.6k | 6k |
| 64 | 50.5 | 1k | 3k | *31k* |
| 128 | 119.5 | 4k | 13k | *134k* |
| 512 | 182.8 | 6k | *53k* | *538k* |

In slide-based multimodal presentations, the video stream is highly redundant and consists mostly of slides spanning time segments to cover it. We conjecture that the visual stream can be preprocessed to keep only a few key frames without significant loss of information. The method for obtaining these key frames is described in the paper that introduces the TIB dataset [14] and focuses on retrieving the slides from the video stream.

In order to validate this assumption, we feed the model with either the full video or the extracted key frames. We conduct experiments with different sample rates: 0.1 frames per second, following Liu et al. [28], or 1 frame per second. The average input token length rapidly grows very large with more tokens per frame, as detailed in Table 2.
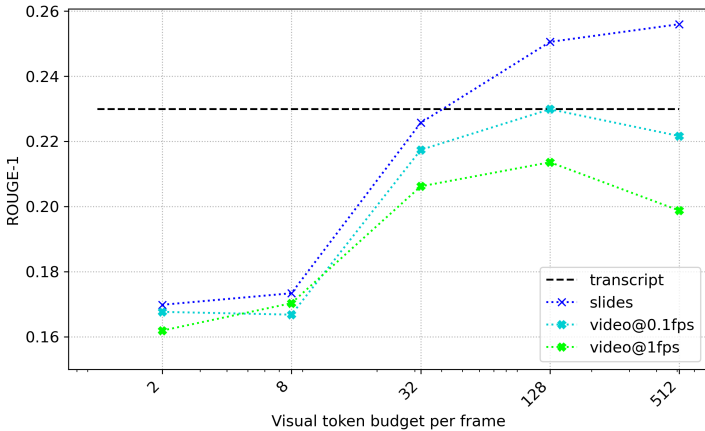


Fig. 3. Unimodal performance at different token budgets

Figure 3 shows the performance of the model with input of a single modality: video, slides or transcript. At all visual token budgets per frame, the slides input outperforms the video input, all while consisting of much fewer frames (*cf* Table 2). The use of key frame input demonstrates a notable advantage with significantly improved performance and substantially lower computational

cost, as evidenced in Table 2. The key frame input with high per slide budget also compares favorably with the transcript input.

As a way to represent the multimodal input for the model, we consider two options: concatenating the text after the slides or interleaving the tokens according to the temporal information gathered during preprocessing. In the interleaved format, the visual tokens for a slide are directly followed by the tokens of the transcribed simultaneous speech. Thus the images are near the related text in the sequence. The other option is a sequence of all slides followed by the transcribed speech. In this default representation, we do not leverage any knowledge about the structure of the presentation.

We conjecture that the interleaved format gives useful information about the structure of the presentation to the model. In order to verify and quantify this claim, we feed the model interleaved or concatenated modalities and evaluate the generated summaries. We call these experiments *structure* and *no structure* respectively in Figure 4.

Due to the scaling of the Transformer computation with sequence length [40], a larger visual token budget mechanically increases the computation for a given document. The additional tokens can also be repurposed by the model for internal computations [12], in a phenomenon called "register" tokens. In an attempt to decouple this additional computation with the benefit of visual information, we conducted an ablation study by replacing the slides with placeholder black images. When slides are not included in the input, these placeholder slides are added to keep the same total computation budget, and watch for an effect similar to the "register" tokens discussed by Darcet et al. [12].
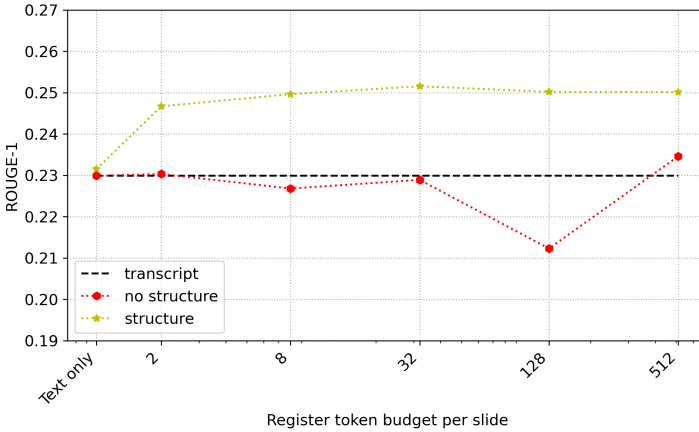


Fig. 4. The addition of structure improves the Rouge score compared to the transcript alone

The resulting experiments enable us to watch for the outcome of the structure information for summarization, decoupled from the effect of additional computations. As illustrated in Figure 4, scaling the register token budget for the placeholder slides does not show a trend on the measured Rouge score, but random variation around the mean that we believe is an artifact of the added noise. However, the structure information offers a significant improvement compared to the text-only baseline.

## 5.2 Model and visual token budget scaling

Figure 5 shows the scaling of the Rouge score with varying visual token budgets for various input representations. With low visual token budgets, visual information is not as important as structure
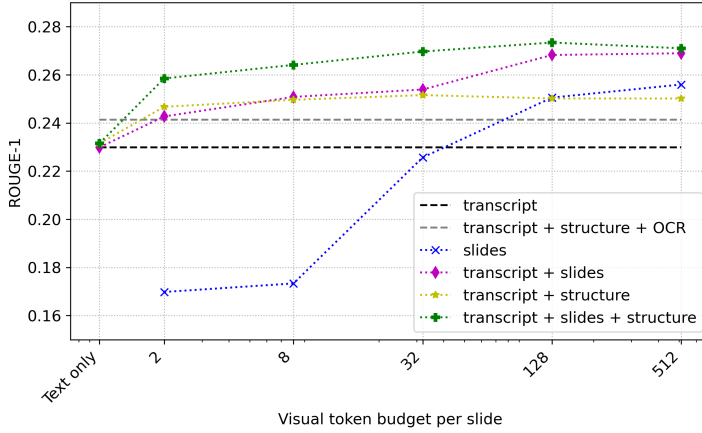
Fig. 5. Rouge score with different visual token budgets

or text information. However, when the budget exceeds 32 tokens per slide, the visual information becomes the predominant modality. At low budgets, the structure information vastly increases the results, but its importance lessens with more budget, to the point where it seems to converge with the unstructured text and image input at 512 tokens per slide.

Comparing the transcript only with the OCR-augmented textual baseline in Figure 5, we see that access to the text printed on the slide improves the performance. However, the slides can be more informative than the text printed on them. For example, the formatting of the text, the colors, and the illustrations can convey meaning. In fact, the summaries generated with the visual input greatly exceed the textual baseline augmented with OCR, suggesting that the textual content of the slides covers only a small amount of the overall semantics they provide.
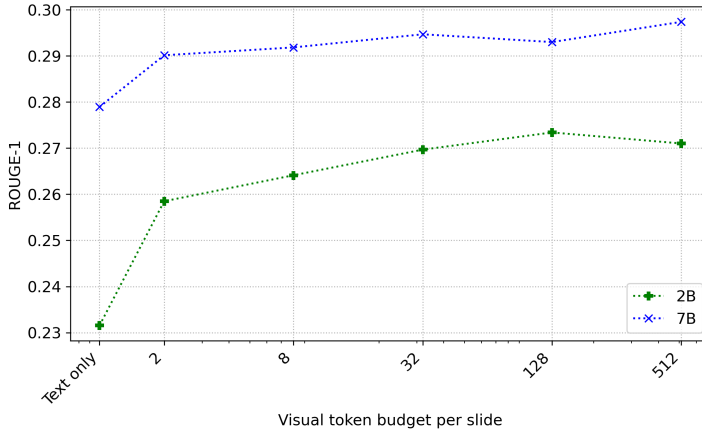


Fig. 6. Rouge score with varying visual token budgets for different model sizes

As illustrated in Figure 6, the improvement related to the scaling of the model from 2 billion to 7 billion parameters exceeds the increase caused by the additional visual and structure information. The larger model still benefits from the additional structure and modality.
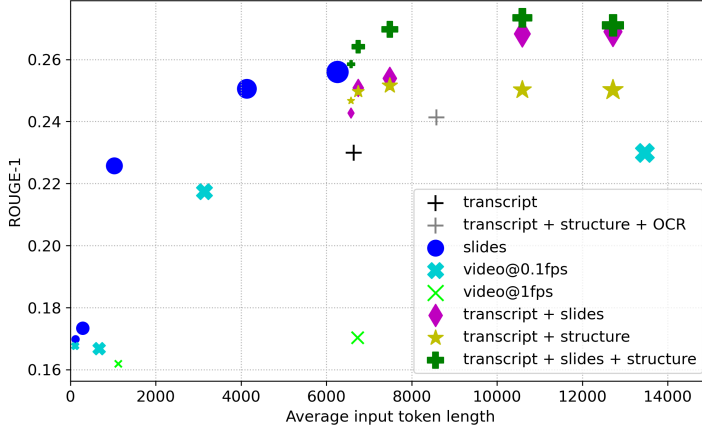


Fig. 7. Slides-only and structured multimodal representations constitute input-length-optimal representations for different input lengths. Bigger icons depict higher visual token budget.

## 5.3 Cost-effective input settings

The unimodal and multimodal input settings can have very different input lengths. In order to devise input-length-effective settings, we represent the Rouge score against the average input token length in Figure 7. If we only account the cost of input size with the goal of improving Rouge score, the strategies resulting in the best tradeoffs draw a set of the input-length-optimal solutions, called Pareto frontier.

With an average input token length shorter than 6k, using only the slides as input is the most efficient strategy. For average token lengths larger than 6k, structured multimodal input with interleaved slides and transcript is the most effective. These two strategies draw the Pareto frontier of the most efficient strategies with any input token length. The text-only solutions that use transcript and/or OCR are significantly below this Pareto frontier, as are the solutions that use the whole video.

## 5.4 Extractiveness and relevance of the extractive fragments with respect to input modalities

Figures 8a and 9a show the extractive coverage and density in unimodal and multimodal input settings, with respect to either the text extracted from the slides, from the speech or their concatenation, called ocr, transcript and overall modes respectively. Compared to single-modality input, adding an additional modality results in less extractive summaries with respect to the original modality. Note that the overall mode is significantly biased towards the text extracted from speech, because of the imbalance illustrated in Figure 2, and as a consequence its change will be mostly correlated to the transcript mode.

Figure 10a shows the importance-based relevance score in unimodal and multimodal settings. Similar scores are observed in the ocr and transcript modes between the unimodal or multimodal inputs, although the summaries were shown to be less extractive. This score computes an overlap
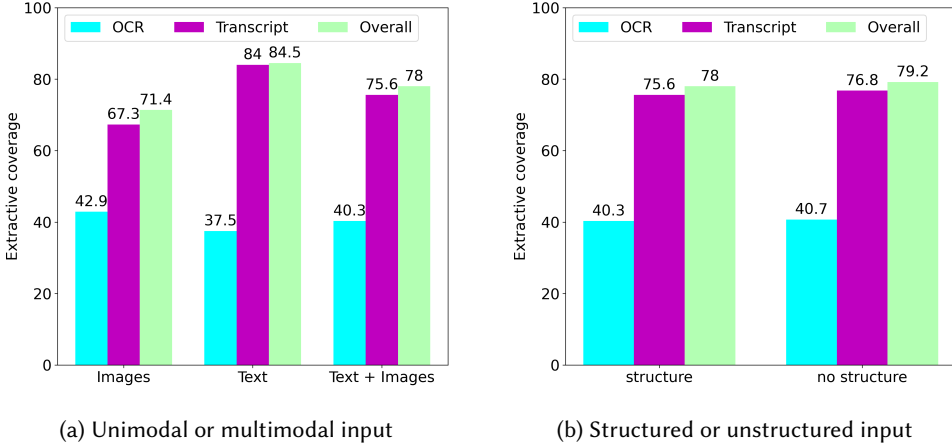
(a) Unimodal or multimodal input                                    (b) Structured or unstructured input

Fig. 8. Extractive coverage of the predicted summaries with reference to input modality



(a) Unimodal or multimodal input                                    (b) Structured or unstructured input
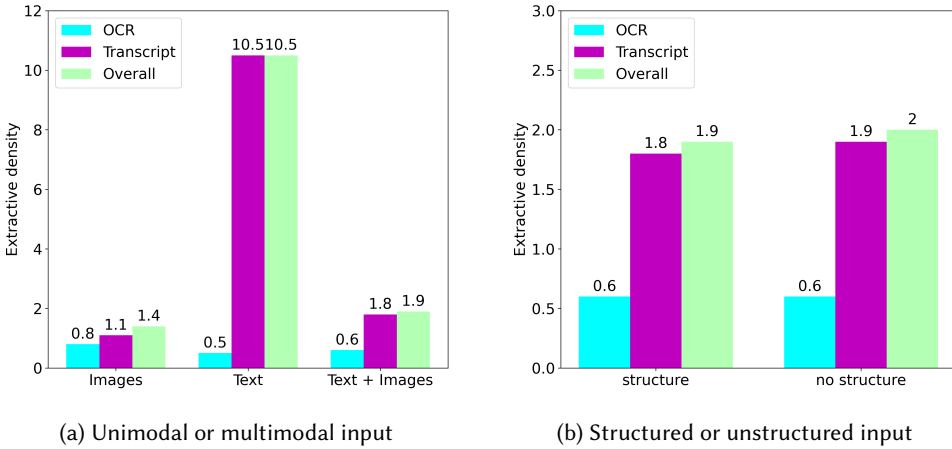
Fig. 9. Extractive density of the predicted summaries with reference to input modality

of $n$-grams between the summary and the source, weighted by the importance of the $n$-grams. A less extractive summary will mechanically have a lower score. Therefore, the extractive fragments in the multimodal setting are more relevant on average compared to the single-modality settings.

Figure 11 illustrates this change in model behavior when scaling the visual token budget in the structured multimodal representation. The summaries become less extractive with respect to the speech (via transcript), while keeping a similar relevance, meaning that the extractive fragments are more relevant on average. At the same time, the summaries become more extractive and more relevant with respect to visual content (via OCR). Overall, the summaries are less extractive and more relevant.

We observe that with sufficient visual token budget, the structure and unstructured multimodal input seemed to converge to comparable performance. The extractive statistics and the importance metric pictured in Figures 8b, 9b, and 10b, corroborate these findings. With large visual token
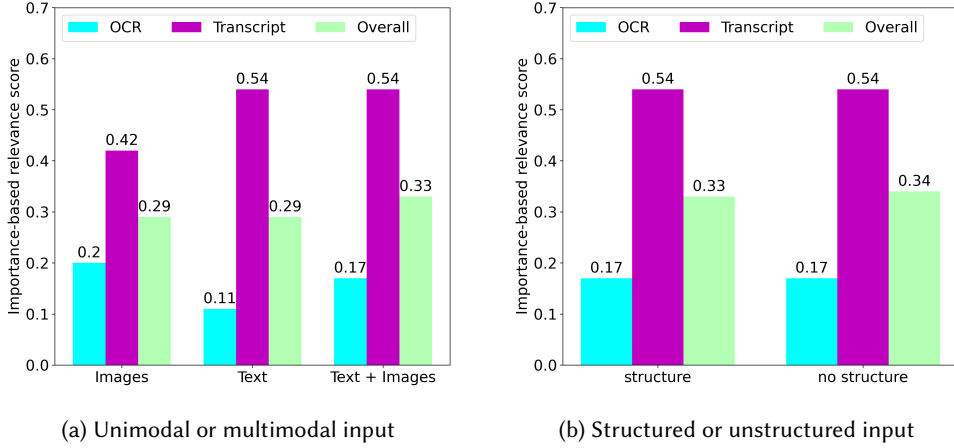
(a) Unimodal or multimodal input

(b) Structured or unstructured input

Fig. 10. Importance-based relevance score of the predicted summaries with reference to input modality
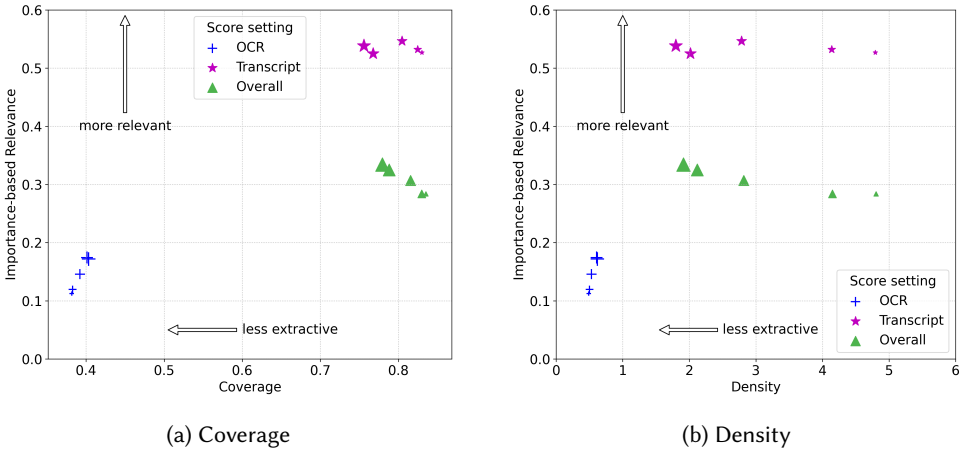


(a) Coverage

(b) Density

Fig. 11. Scaling the visual token budget balances the extractive source from transcript to OCR and improves the relevance with respect to the OCR and the overall document. Bigger icons depict higher visual token budgets.

budget, the model using structured or unstructured input behaves similarly with respect to the extractiveness and the importance of the extractive fragments.

## 5.5 Qualitative Analysis

During manual inspection of the generated summaries, we notice different behaviors with different input settings for the model. For instance, summaries generated with a multimodal input tend to be more comprehensive than those generated with slides or transcripts alone.

As illustrated in Table 3, the visual information allows the model to correct misspelled words from the transcript, often proper nouns. With the example from Table 4 we see that it also includes contextual information that is not always found in the transcript, such as the name of the speaker and the venue and event where the presentation is taking place. Conversely, the speech tends to provide

more detailed information compared to the slides, as displayed in Table 5. The multimodal input seems to frequently trigger an unintended "bullet point" mode that was observed less frequently with unimodal inputs. This behaviour is shown in Table 6.

Table 3. Summaries generated from the presentation "Clojure: a language for problem solvers", display the spell-checking abilities of the additional modalities

| Summary with transcript input | Summary with slides input | Summary with multimodal input |
|---|---|---|
| The presentation by Stuart Halloway is about closure, a programming language for problem solvers. The speaker introduces the topic of closure and its importance in software development. He discusses the characteristics of closure, such as being knowledgeable, flexible, powerful, and smart. He also mentions the challenges of working with transient data structures and the need for persistent data structures. Halloway explains how closure uses persistent data structures to solve these challenges and how it provides direct access to the platform. He also discusses the benefits of using closure, such as being able to build systems quickly and having a powerful set of persistent data structures. The presentation ends with a summary of the benefits of using closure and a call to action for attendees to learn more about the language. | The presentation is about Clojure, a language designed for problem solvers. The speaker, Stuart Halloway, introduces Clojure and its features, such as its ability to handle transience and persistence, and its use in various programming contexts. The presentation includes slides that explain the language's syntax and features, as well as examples of how it can be used to solve problems. The slides also show how Clojure can be used to implement interfaces and work with collections, directories, files, XML, JSON, and ResultSets. The presentation ends with a question and answer session and a list of resources for further learning. | The presentation by Stuart Halloway is about Clojure, a programming language for problem solvers. The main focus of the presentation is to explain the benefits of using Clojure for software development, particularly in terms of flexibility, power, and knowledge. Halloway discusses the importance of having knowledge in software development and how Clojure provides a way to build systems that are knowledgeable and powerful. He also emphasizes the importance of having flexible software systems that can adapt to changing requirements. The presentation includes slides that explain the different types of data structures in Clojure, such as vectors, maps, sets, and lists, and how they can be used to build powerful and flexible software systems. Halloway also discusses the benefits of using Clojure for developing large-scale systems, such as the ability to build systems that are composable and can be easily modified and extended. The presentation also includes a demonstration of Clojure in action, using a code example to show how it can be used to build a simple program. Overall, the presentation provides a comprehensive overview of Clojure and its benefits for software development. |

Table 4. Summaries generated from the presentation "Keyword Args — the killer Ruby feature you aren't using", illustrating the additional context when leveraging the slides

| Summary with transcript input | Summary with slides input | Summary with multimodal input |
|---|---|---|
| The speaker is introducing the topic of keyword arguments in Ruby programming. They explain that keyword arguments are a feature in Ruby that allows for more flexible and functional programming. They provide examples of how keyword arguments can be used to pass arguments to functions and how they can be used to make functions more flexible. The speaker also discusses the benefits of using keyword arguments over positional arguments, such as making functions more clear and flexible. | The presentation is about RubyConf Cincinnati 2016, featuring a talk by Guyren G Howe titled "Keyword Args - the ikller Ruby feature you aren't using." The speaker discusses the use of keyword arguments in Ruby, highlighting their benefits such as making code more intuitive, clearer, more flexible, and more functional. The presentation includes slides that explain the different types of keyword arguments, such as positional, required, and optional arguments, and how they can be used to make code more expressive. The speaker also demonstrates how keyword arguments can be used to create more flexible and functional code. The presentation concludes with a discussion on dependency injection and context funneling, which are other Ruby features that are not commonly used. | The presentation by Guyren G Howe at RubyConf Cincinnati 2016 focuses on keyword arguments in Ruby programming. The speaker explains that keyword arguments are a feature in Ruby that allows for more flexible and functional programming. They are particularly useful for passing positional arguments, optional arguments, and keyword arguments to functions. The speaker provides examples of how keyword arguments can make code more clear, more flexible, and more functional. He also discusses the benefits of using keyword arguments over positional arguments, such as the ability to pass in optional arguments and the ability to override default values. The presentation concludes with a summary of the key points and a call to action for attendees to consider using keyword arguments in their Ruby code. |

## 6    Discussion

We think that the behaviors observed in the qualitative analysis are a consequence of the training data domain of VLMs. Specifically, the underrepresentation of conflict and mismatch between the textual and visual input.

Qwen2-VL, similarly to most VLMs, is trained with data covering multiple tasks that rely mostly on translating the information from one modality to another, *eg* captioning or visual question answering. McKinzie et al. [34] observed that common VLM evaluation is focused on captioning problems, with 3 out of 8 benchmarks in their study being captioning. They discriminate the

Table 5. Summaries generated from the presentation "UrbanSim2: Simulating the Connected Metropolis", illustrating the additional context when leveraging the transcript: the slides mention Paul Waddell but the speaker is Fletcher Foti as said in the presentation

| Summary with transcript input | Summary with slides input | Summary with multimodal input |
| --- | --- | --- |
| Fletcher Foti is a PhD in city planning who is now a chief data scientist at Synthesity. He is presenting on UrbanSim, a simulation tool that simulates the connected metropolis. He explains that UrbanSim is an agent-based simulation of regional real estate markets, and that it captures individual households, individual jobs, and the decisions that individual households make. He also mentions that UrbanSim is an open-source project that is BSD licensed. | The presentation is titled "What is UrbanSim?" and is part of the FOSS4G PDX 2014 conference. The speaker, Paul Waddell, is a professor at the University of California, Berkeley. The presentation covers the purpose of UrbanSim, which is an agent-based simulation of regional real estate markets. The simulation is written in Java in 2000 and Python in 2006. The code is implemented to use Pandas this year. The presentation also discusses the regional task, which includes reducing greenhouse gas emissions, housing the region's population at all income levels, and embodying local visions. The speaker also mentions the Open Source Stack, which includes tools like Pandas, Python, and Scikit. | The presentation by Fletcher Foti focused on UrbanSim, a statistical model used for simulating the connected metropolis. The speaker explained that UrbanSim is used to forecast change in population, employment, and the built environment over a 30-year planning horizon for a given metropolitan region. The model is agent-based, simulating individual households, jobs, and land use decisions. The speaker also highlighted the open-source nature of the project, with UrbanSim being written in Java in 2000 and Python in 2006. The model was re-implemented in 2013 on new scientific Python tools, specifically Pandas. The speaker also discussed the open-source stack, which includes Python, Pandas, and other libraries. The presentation included a demo of the Urban Data Science Toolkit, which includes tools for network analysis, data visualization, and web mapping. The speaker also mentioned the Sync City project, which is a large-scale transportation model required by federal law. |

Table 6. Summaries generated from the presentation "Genode meets the Pinephone", exemplifying the unprompted use of bullet point lists with the multimodal input

| Summary with transcript input | Summary with slides input | Summary with multimodal input |
| --- | --- | --- |
| The talk was about the last year's journey to bring together the GenoOS framework and the Pinephone. The speaker, Norman Feske, from GenoDepps, the developer behind the GenoOS framework, gave a presentation on the topic. | The presentation is titled "Genode meets the Pinephone" and is given by Norman Feske. The presentation covers the motivation behind the development of Genode, which is a microkernel-based operating system. The presentation also discusses the corporate motives of the developers and the sentiments shared by members of civil society. The development story of Genode is presented, including the boot steps, device access, and cascaded authorities. The presentation also covers the development of Genode on the Pinephone, including the device tree, device drivers, and the Pinephone as a development platform. The presentation ends with a preview of the Genode OS on the Pinephone and a list of goals for 2022. | The presentation by Norman Feske at FOSDEM 2022 focuses on the development of GenoOS, a microkernel-based operating system for smartphones. The presentation covers the motivation behind the project, the development process, and the current status of the project. Key points include: - The project aims to combine the GenoOS framework with the PinePhone, a smartphone that uses a custom microkernel. - The development process involves booting the device, creating a custom Linux kernel, and porting the kernel to the PinePhone. - The development team has made progress in porting the kernel and has successfully run a small application on the PinePhone. - The team is working on improving performance and adding new features, such as persistent storage and mobile data connectivity. - The presentation also includes a demonstration of running a small application on the PinePhone and a preview of the Sculpt OS on the PinePhone. - The team plans to continue working on the project in 2022, with a focus on video telephony, persistent storage, mobile data connectivity, power management, and performance improvements. |

multimodal data used to train VLMs between captioning image-text pairs and interleaved image-text documents. They describe their interleaved image-text data as longer than the captioning data and with more diverse text with less relevance on average to images.

We believe that the interleaved image-text that is extracted from multimodal presentations describes a different flavor of cross-modal interactions compared to the interleaved image-text documents typically found in VLMs' training data. The latter's interleaved image-text documents mainly come from articles, in which the images are usually used for merely illustrative purposes.

### 6.1 Non redundancy and conflict between modalities

In multimodal presentations, information can be split between modalities with limited redundancy, as illustrated in Tables 4 and 5. In order to create a comprehensive summary, a VLM should be able to cross-reference modalities. Some visual question answering benchmarks [9, 22] are designed to require cross-referencing between modalities. They show the limited performance of VLMs in this task and their dependence on external tools such as information retrieval pipelines.

Additionally, mismatched or conflicting information can be presented across modalities. For example, when the speaker verbally states a fact that contradicts the content displayed on a slide, due to a typographical error, a last-minute change (*cf* Table 5), or an *erratum*. The contradiction can also be sourced from inaccuracies during preprocessing, *eg* the transcription (*cf* Table 3). In our understanding of the training data domains, VLMs are not substantially faced with this type of data during training, and therefore are not explicitly trained to handle these situations. Liu et al. [31] explored the conflict between a VLM's internal knowledge and its visual input. They proposed a prompting strategy as a way to steer the VLM into trusting vision over internal knowledge. Sarcasm detection in VLMs [8] is another line of work that represents a typical example of conflicting modalities. In these cases, the task implicitly defines a *stronger* modality, on which trust should be focused. In multimodal presentations, errors can appear in both modalities, even in the same sample. Consequently, no overall *stronger* modality can be chosen. Liu et al. [30] described *weak* modalities as modalities that are very noisy but could provide relevant information on a per sample basis. They advocate for a multiplicative method rather than additive to merge modalities in a late-fusion scenario, effectively leveraging the lower uncertainty, or entropy, associated with weak modalities in relevant situations. Modern VLMs are early-fusion models. Therefore it is not trivial to adapt this uncertainty-based method to merge decisions based on each modality.

### 6.2 Visual and written concepts entanglement in the visual modality

The visual modality contains information that can be depicted or written. Goh et al. [16] showed that printing textual information over the input on multimodal vision models can steer the output, acting as "typographical attacks". We believe that the unprompted "bullet point" mode (*cf* Table 6) which follows the format usually seen in slides, emerging in multimodal contexts, could derive from this phenomenon. Materzyńska et al. [33] further studied the entanglement of visual and written concepts in CLIP. They hypothesized that it is a consequence of the prevalence of text next to the concept they represent in the training data, such as signs and labels. Overall, text written in the visual modality, characteristic of the text-heavy visual stream of multimodal presentations, has an unpredictable effect on the model behavior. It can be used as a way to interfere with the instructions via typographical attacks or be confused with visual concepts.

### 6.3 VLMs are mostly not agnostic to the input modalities

Similar information can be provided to VLMs through visual or textual modalities, in practice resulting in very different outcomes. Liang et al. [25] observed a modality gap in the embeddings of vision-language encoders such as CLIP. This gap is characterized by the representation of an image being closer to representations of different images than to the representation of its associated caption in the shared representation space. Zhu et al. [50] observed conflicts in the knowledge embedded in the model parameters depending on the input modality, even when the model accurately captions the images. Alternatively, Templeton et al. [38] managed to extract monosemantic features from a VLM, some of which were multimodal and multilingual. These features are directions in the activation space that represent concepts that are agnostic to language and modality. Zhang et al. [49] described a phenomenon dubbed cross-modal consistency associated with the performance

gap of a VLM on the same task depending on the input modality. In their experiments, this gap is substantial and results in the VLM being inconsistent with different input modalities. Deng et al. [13] observed a "Blind Faith in Text" effect where VLMs tend to trust textual input when it is inconsistent with the visual data, or incorrect. In contrast, in our experiments, the model was able to correct errors in the transcript by relying on visual information (*cf* Table 3), probably because it was informed in its prompt that the text is a transcript, therefore implicitly error prone.

These works illustrate some of the inherent challenges in the abstractive summarization of multimodal presentations. The internal representation of VLMs can contain modality-agnostic features but is still overall vastly inconsistent between and inside modalities.

## 7 Conclusion

Available state-of-the-art open-weights Vision-Language Models are still overall quite unsuccessful at summarizing long multimodal presentations. We find some of them, notably Qwen2-VL, which was likely trained on adjacent tasks that involve long-range text-heavy interleaved image-text documents, to express promising downstream performance for multimodal summarization.

Through quantitative and qualitative analyses, we show the benefit of representing the input presentation record as interleaved slides and transcript. This structured multimodal representation results in higher scores for VLMs and extends the Pareto frontier of input-length-optimal methods beyond the use of slides alone. The use of both modalities offers many benefits, such as providing additional information and enabling spell checking through cross-referencing modalities. Multimodal systems produce summaries that are less extractive with respect to speech or text displayed on slides. They also include more relevant elements of the presentation, sourced from either one or both modalities. However, discrepancies between transcript and slides pose challenges for the VLM and occasionally exhibit impromptu side effects.

We discuss some of the challenges that arise from interactions between conflicting modalities in multimodal presentations. Consequently, we advocate for the inclusion of more diverse cross-modality interactions in the training data mixtures for VLMs, specifically documents involving mismatched or conflicting modalities. Such data is commonly encountered in real-world scenarios such as multimodal presentations, where the visual content displayed is often commented upon, expanded, or even criticized in the accompanying speech. Data of this nature can also be sourced from reviews and commentaries, or created synthetically by altering a modality [13].

By incorporating more documents with conflicting modalities into datasets and benchmarks, future research can better understand and mitigate the effects of cross-modal interactions in VLMs. This would pave the way for the development of more robust, reliable, and trustworthy multimodal systems.

## References

[1] Marah I. Abdin, Sam Ade Jacobs, Ammar Ahmad Awan, Jyoti Aneja, Ahmed Awadallah, Hany Awadalla, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Harkirat S. Behl, Alon Benhaim, Misha Bilenko, Johan Bjorck, Sébastien Bubeck, Martin Cai, Caio César Teodoro Mendes, Weizhu Chen, Vishrav Chaudhary, Parul Chopra, Allie Del Giorno, Gustavo de Rosa, Matthew Dixon, Ronen Eldan, Dan Iter, Amit Garg, Abhishek Goswami, Suriya Gunasekar, Emman Haider, Junheng Hao, Russell J. Hewett, Jamie Huynh, Mojan Javaheripi, Xin Jin, Piero Kauffmann, Nikos Karampatziakis, Dongwoo Kim, Mahoud Khademi, Lev Kurilenko, James R. Lee, Yin Tat Lee, Yuanzhi Li, Chen Liang, Weishung Liu, Eric Lin, Zeqi Lin, Piyush Madan, Arindam Mitra, Hardik Modi, Anh Nguyen, Brandon Norick, Barun Patra, Daniel Perez-Becker, Thomas Portet, Reid Pryzant, Heyang Qin, Marko Radmilac, Corby Rosset, Sambudha Roy, Olatunji Ruwase, Olli Saarikivi, Amin Saied, Adil Salim, Michael Santacroce, Shital Shah, Ning Shang, Hiteshi Sharma, Xia Song, Masahiro Tanaka, Xin Wang, Rachel Ward, Guanhua Wang, Philipp Witte, Michael Wyatt, Can Xu, Jiahang Xu, Sonali Yadav, Fan Yang, Ziyi Yang, Donghan Yu, Chengruidong Zhang, Cyril Zhang, Jianwen Zhang, Li Lyna Zhang, Yi Zhang, Yue Zhang, Yunan Zhang, and Xiren Zhou. 2024. Phi-3 Technical Report: A Highly Capable Language Model Locally on Your Phone. *CoRR* (Jan. 2024). https://openreview.net/forum?id=k05PGWgTv2

[2]  Joshua Ainslie, Santiago Ontanon, Chris Alberti, Vaclav Cvicek, Zachary Fisher, Philip Pham, Anirudh Ravula, Sumit
      Sanghai, Qifan Wang, and Li Yang. 2020. ETC: Encoding Long and Structured Inputs in Transformers. In *Proceedings of
      the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, Bonnie Webber, Trevor Cohn, Yulan
      He, and Yang Liu (Eds.). Association for Computational Linguistics, Online, 268–284. doi:10.18653/v1/2020.emnlp-
      main.19

[3]  Toqa Alaa, Ahmad Mongy, Assem Bakr, Mariam Diab, and Walid Gomaa. 2024. Video Summarization Techniques: A
      Comprehensive Review. doi:10.48550/arXiv.2410.04449 arXiv:2410.04449 [cs].

[4]  Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch,
      Katherine Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford, Serkan Cabi, Tengda Han, Zhitao Gong, Sina
      Samangooei, Marianne Monteiro, Jacob L. Menick, Sebastian Borgeaud, Andy Brock, Aida Nematzadeh, Sahand
      Sharifzadeh, Mikołaj Bińkowski, Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, and Karén Simonyan. 2022.
      Flamingo: a Visual Language Model for Few-Shot Learning. *Advances in Neural Information Processing Systems* 35 (Dec.
      2022), 23716–23736. https://proceedings.neurips.cc/paper_files/paper/2022/hash/960a172bc7fbf0177ccccbb411a7d800-
      Abstract-Conference.html

[5]  Loubna Ben Allal, Anton Lozhkov, Elie Bakouch, Gabriel Martín Blázquez, Guilherme Penedo, Lewis Tunstall, Andrés
      Marafioti, Hynek Kydlíček, Agustín Piqueres Lajarín, Vaibhav Srivastav, Joshua Lochner, Caleb Fahlgren, Xuan-Son
      Nguyen, Clémentine Fourrier, Ben Burtenshaw, Hugo Larcher, Haojun Zhao, Cyril Zakka, Mathieu Morlon, Colin
      Raffel, Leandro von Werra, and Thomas Wolf. 2025. SmolLM2: When Smol Goes Big – Data-Centric Training of a
      Small Language Model. doi:10.48550/arXiv.2502.02737 arXiv:2502.02737 [cs] version: 1.

[6]  Evlampios Apostolidis, Eleni Adamantidou, Alexandros I. Metsai, Vasileios Mezaris, and Ioannis Patras. 2021. Video
      Summarization Using Deep Neural Networks: A Survey. *Proc. IEEE* 109, 11 (Nov. 2021), 1838–1863. doi:10.1109/JPROC.
      2021.3117472 Conference Name: Proceedings of the IEEE.

[7]  Iz Beltagy, Matthew E. Peters, and Arman Cohan. 2020. Longformer: The Long-Document Transformer. doi:10.48550/
      arXiv.2004.05150 arXiv:2004.05150 [cs].

[8]  Santiago Castro, Devamanyu Hazarika, Verónica Pérez-Rosas, Roger Zimmermann, Rada Mihalcea, and Soujanya
      Poria. 2019. Towards Multimodal Sarcasm Detection (An _Obviously_ Perfect Paper). In *Proceedings of the 57th Annual
      Meeting of the Association for Computational Linguistics*, Anna Korhonen, David Traum, and Lluís Màrquez (Eds.).
      Association for Computational Linguistics, Florence, Italy, 4619–4629. doi:10.18653/v1/P19-1455

[9]  Yang Chen, Hexiang Hu, Yi Luan, Haitian Sun, Soravit Changpinyo, Alan Ritter, and Ming-Wei Chang. 2023. Can
      Pre-trained Vision and Language Models Answer Visual Information-Seeking Questions? https://openreview.net/
      forum?id=3MEV5aIDDq

[10] Tri Dao. 2023. FlashAttention-2: Faster Attention with Better Parallelism and Work Partitioning. doi:10.48550/arXiv.
      2307.08691 arXiv:2307.08691 [cs].

[11] Tri Dao, Dan Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. 2022. FlashAttention: Fast and Memory-Efficient
      Exact Attention with IO-Awareness. *Advances in Neural Information Processing Systems* 35 (Dec. 2022), 16344–
      16359. https://proceedings.neurips.cc/paper_files/paper/2022/hash/67d57c32e20fd0a7a302cb81d36e40d5-Abstract-
      Conference.html

[12] Timothée Darcet, Maxime Oquab, Julien Mairal, and Piotr Bojanowski. 2023. Vision Transformers Need Registers.
      https://openreview.net/forum?id=2dnO3LLiJ1

[13] Ailin Deng, Tri Cao, Zhirui Chen, and Bryan Hooi. 2025. Words or Vision: Do Vision-Language Models Have Blind
      Faith in Text? doi:10.48550/arXiv.2503.02199 arXiv:2503.02199 [cs].

[14] Théo Gigant, Frederic Dufaux, Camille Guinaudeau, and Marc Décombas. 2023. TIB: A Dataset for Abstractive
      Summarization of Long Multimodal Videoconference Records. In *Proceedings of the 20th International Conference on
      Content-based Multimedia Indexing (CBMI '23)*. Association for Computing Machinery, New York, NY, USA, 61–70.
      doi:10.1145/3617233.3617238

[15] Théo Gigant, Camille Guinaudeau, Marc Decombas, and Frederic Dufaux. 2024. Mitigating the Impact of Reference
      Quality on Evaluation of Summarization Systems with Reference-Free Metrics. In *Proceedings of the 2024 Conference
      on Empirical Methods in Natural Language Processing*, Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (Eds.).
      Association for Computational Linguistics, Miami, Florida, USA, 19355–19368. doi:10.18653/v1/2024.emnlp-main.1078

[16] Gabriel Goh, Nick Cammarata †, Chelsea Voss †, Shan Carter, Michael Petrov, Ludwig Schubert, Alec Radford, and
      Chris Olah. 2021. Multimodal Neurons in Artificial Neural Networks. *Distill* 6, 3 (March 2021), e30. doi:10.23915/distill.00030

[17] Mandy Guo, Joshua Ainslie, David Uthus, Santiago Ontanon, Jianmo Ni, Yun-Hsuan Sung, and Yinfei Yang. 2022.
      LongT5: Efficient Text-To-Text Transformer for Long Sequences. In *Findings of the Association for Computational
      Linguistics: NAACL 2022*, Marine Carpuat, Marie-Catherine de Marneffe, and Ivan Vladimir Meza Ruiz (Eds.). Association
      for Computational Linguistics, Seattle, United States, 724–736. doi:10.18653/v1/2022.findings-naacl.55

[18] Geewook Kim, Teakgyu Hong, Moonbin Yim, JeongYeon Nam, Jinyoung Park, Jinyeong Yim, Wonseok Hwang,
      Sangdoo Yun, Dongyoon Han, and Seunghyun Park. 2022. OCR-Free Document Understanding Transformer. In

*Computer Vision – ECCV 2022*, Shai Avidan, Gabriel Brostow, Moustapha Cissé, Giovanni Maria Farinella, and Tal Hassner (Eds.). Springer Nature Switzerland, Cham, 498–517. doi:10.1007/978-3-031-19815-1_29

[19] Huan Yee Koh, Jiaxin Ju, He Zhang, Ming Liu, and Shirui Pan. 2022. How Far are We from Robust Long Abstractive Summarization?. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (Eds.). Association for Computational Linguistics, Abu Dhabi, United Arab Emirates, 2682–2698. doi:10.18653/v1/2022.emnlp-main.172

[20] Hugo Laurençon, Andrés Marafioti, Victor Sanh, and Leo Tronchon. 2024. Building and better understanding vision-language models: insights and future directions. https://openreview.net/forum?id=iSL0FHZStr

[21] Hugo Laurençon, Leo Tronchon, Matthieu Cord, and Victor Sanh. 2025. What matters when building vision-language models? *Advances in Neural Information Processing Systems* 37 (Jan. 2025), 87874–87907. https://proceedings.neurips.cc/paper_files/paper/2024/hash/a03037317560b8c5f2fb4b6466d4c439-Abstract-Conference.html

[22] Paul Lerner, Olivier Ferret, Camille Guinaudeau, Hervé Le Borgne, Romaric Besançon, Jose G. Moreno, and Jesús Lovón Melgarejo. 2022. ViQuAE, a Dataset for Knowledge-based Visual Question Answering about Named Entities. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '22)*. Association for Computing Machinery, New York, NY, USA, 3108–3120. doi:10.1145/3477495.3531753

[23] Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: Denoising Sequence-to-Sequence Pre-training for Natural Language Generation, Translation, and Comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel Tetreault (Eds.). Association for Computational Linguistics, Online, 7871–7880. doi:10.18653/v1/2020.acl-main.703

[24] Tianle Li, Ge Zhang, Quy Duc Do, Xiang Yue, and Wenhu Chen. 2024. Long-context LLMs Struggle with Long In-context Learning. *CoRR* (Jan. 2024). https://openreview.net/forum?id=CFYzMJJB6a

[25] Victor Weixin Liang, Yuhui Zhang, Yongchan Kwon, Serena Yeung, and James Y. Zou. 2022. Mind the Gap: Understanding the Modality Gap in Multi-modal Contrastive Representation Learning. *Advances in Neural Information Processing Systems* 35 (Dec. 2022), 17612–17625. https://proceedings.neurips.cc/paper_files/paper/2022/hash/702f4db7543a7432431df588d57bc7c9-Abstract-Conference.html

[26] Bin Lin, Yang Ye, Bin Zhu, Jiaxi Cui, Munan Ning, Peng Jin, and Li Yuan. 2024. Video-LLaVA: Learning United Visual Representation by Alignment Before Projection. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (Eds.). Association for Computational Linguistics, Miami, Florida, USA, 5971–5984. doi:10.18653/v1/2024.emnlp-main.342

[27] Chin-Yew Lin. 2004. ROUGE: A Package for Automatic Evaluation of Summaries. In *Text Summarization Branches Out*. Association for Computational Linguistics, Barcelona, Spain, 74–81. https://aclanthology.org/W04-1013

[28] Dongqi Liu, Chenxi Whitehouse, Xi Yu, Louis Mahon, Rohit Saxena, Zheng Zhao, Yifu Qiu, Mirella Lapata, and Vera Demberg. 2025. What Is That Talk About? A Video-to-Text Summarization Dataset for Scientific Presentations. doi:10.48550/arXiv.2502.08279 arXiv:2502.08279 [cs].

[29] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. 2023. Visual Instruction Tuning. *Advances in Neural Information Processing Systems* 36 (Dec. 2023), 34892–34916. https://papers.nips.cc/paper_files/paper/2023/hash/6dcf277ea32ce3288914faf369fe6de0-Abstract-Conference.html

[30] Kuan Liu, Yanen Li, Ning Xu, and Prem Natarajan. 2018. Learn to Combine Modalities in Multimodal Deep Learning. doi:10.48550/arXiv.1805.11730 arXiv:1805.11730 [stat].

[31] Xiaoyuan Liu, Wenxuan Wang, Youliang Yuan, Jen-tse Huang, Qiuzhi Liu, Pinjia He, and Zhaopeng Tu. 2024. Insight Over Sight? Exploring the Vision-Knowledge Conflicts in Multimodal LLMs. doi:10.48550/arXiv.2410.08145 arXiv:2410.08145 [cs].

[32] Ruipu Luo, Ziwang Zhao, Min Yang, Zheming Yang, Minghui Qiu, Tao Wang, Zhongyu Wei, Yanhao Wang, and Cen Chen. 2025. Valley: Video Assistant with Large Language model Enhanced abilitY. doi:10.48550/arXiv.2306.07207 arXiv:2306.07207 [cs].

[33] Joanna Materzyńska, Antonio Torralba, and David Bau. 2022. Disentangling Visual and Written Concepts in CLIP. 16410–16419. https://openaccess.thecvf.com/content/CVPR2022/html/Materzynska_Disentangling_Visual_and_Written_Concepts_in_CLIP_CVPR_2022_paper.html

[34] Brandon McKinzie, Zhe Gan, Jean-Philippe Fauconnier, Sam Dodge, Bowen Zhang, Philipp Dufter, Dhruti Shah, Xianzhi Du, Futang Peng, Anton Belyi, Haotian Zhang, Karanjeet Singh, Doug Kang, Hongyu Hè, Max Schwarzer, Tom Gunter, Xiang Kong, Aonan Zhang, Jianyu Wang, Chong Wang, Nan Du, Tao Lei, Sam Wiseman, Mark Lee, Zirui Wang, Ruoming Pang, Peter Grasch, Alexander Toshev, and Yinfei Yang. 2025. MM1: Methods, Analysis and Insights from Multimodal LLM Pre-training. In *Computer Vision – ECCV 2024*, Aleš Leonardis, Elisa Ricci, Stefan Roth, Olga Russakovsky, Torsten Sattler, and Gül Varol (Eds.). Springer Nature Switzerland, Cham, 304–323. doi:10.1007/978-3-031-73397-0_18

[35] Bowen Peng, Jeffrey Quesnelle, Honglu Fan, and Enrico Shippole. 2023. YaRN: Efficient Context Window Extension of Large Language Models. https://openreview.net/forum?id=wHBfxhZu1u

[36] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. 2021. Learning Transferable Visual Models From Natural Language Supervision. In *Proceedings of the 38th International Conference on Machine Learning*. PMLR, 8748–8763. https://proceedings.mlr.press/v139/radford21a.html ISSN: 2640-3498.

[37] 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

[38] Adly Templeton, Tom Conerly, Jonathan Marcus, Jack Lindsey, Trenton Bricken, Brian Chen, Adam Pearce, Craig Citro, Emmanuel Ameisen, Andy Jones, and others. 2024. Scaling monosemanticity: extracting interpretable features from Claude 3 sonnet.

[39] Oleg Vasilyev and John Bohannon. 2021. ESTIME: Estimation of Summary-to-Text Inconsistency by Mismatched Embeddings. In *Proceedings of the 2nd Workshop on Evaluation and Comparison of NLP Systems*, Yang Gao, Steffen Eger, Wei Zhao, Piyawat Lertvittayakumjorn, and Marina Fomicheva (Eds.). Association for Computational Linguistics, Punta Cana, Dominican Republic, 94–103. doi:10.18653/v1/2021.eval4nlp-1.10

[40] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. 2017. Attention is All you Need. In *Advances in Neural Information Processing Systems*, Vol. 30. Curran Associates, Inc. https://proceedings.neurips.cc/paper/2017/hash/3f5ee243547dee91fbd053c1c4a845aa-Abstract.html

[41] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. 2024. Qwen2-VL: Enhancing Vision-Language Model's Perception of the World at Any Resolution. doi:10.48550/arXiv.2409.12191 arXiv:2409.12191 [cs].

[42] Haoran Wei, Chenglong Liu, Jinyue Chen, Jia Wang, Lingyu Kong, Yanming Xu, Zheng Ge, Liang Zhao, Jianjian Sun, Yuang Peng, Chunrui Han, and Xiangyu Zhang. 2024. General OCR Theory: Towards OCR-2.0 via a Unified End-to-end Model. doi:10.48550/arXiv.2409.01704 arXiv:2409.01704 [cs].

[43] Haoning Wu, Dongxu Li, Bei Chen, and Junnan Li. 2024. LongVideoBench: A Benchmark for Long-context Interleaved Video-Language Understanding. *Advances in Neural Information Processing Systems* 37 (Dec. 2024), 28828–28857. https://proceedings.neurips.cc/paper_files/paper/2024/hash/329ad516cf7a6ac306f29882e9c77558-Abstract-Datasets_and_Benchmarks_Track.html

[44] Weizhe Yuan, Graham Neubig, and Pengfei Liu. 2021. BARTScore: Evaluating Generated Text as Text Generation. In *Advances in Neural Information Processing Systems*, Vol. 34. Curran Associates, Inc., 27263–27277. https://proceedings.neurips.cc/paper_files/paper/2021/hash/e4d2b6e6fdeca3e60e0f1a62fee3d9dd-Abstract.html

[45] Manzil Zaheer, Guru Guruganesh, Kumar Avinava Dubey, Joshua Ainslie, Chris Alberti, Santiago Ontanon, Philip Pham, Anirudh Ravula, Qifan Wang, Li Yang, and Amr Ahmed. 2020. Big Bird: Transformers for Longer Sequences. In *Advances in Neural Information Processing Systems*, Vol. 33. Curran Associates, Inc., 17283–17297. https://papers.neurips.cc/paper_files/paper/2020/hash/c8512d142a2d849725f31a9a7a361ab9-Abstract.html

[46] Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. 2023. Sigmoid Loss for Language Image Pre-Training. In *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*. 11941–11952. doi:10.1109/ICCV51070.2023.01100 ISSN: 2380-7504.

[47] Haopeng Zhang, Philip S. Yu, and Jiawei Zhang. 2024. A Systematic Survey of Text Summarization: From Statistical Methods to Large Language Models. doi:10.48550/arXiv.2406.11289 arXiv:2406.11289 [cs] version: 1.

[48] Peiyuan Zhang, Kaichen Zhang, Bo Li, Guangtao Zeng, Jingkang Yang, Yuanhan Zhang, Ziyue Wang, Haoran Tan, Chunyuan Li, and Ziwei Liu. 2024. Long Context Transfer from Language to Vision. *CoRR* (Jan. 2024). https://openreview.net/forum?id=VZyBATUZjr

[49] Xiang Zhang, Senyu Li, Ning Shi, Bradley Hauer, Zijun Wu, Grzegorz Kondrak, Muhammad Abdul-Mageed, and Laks V. S. Lakshmanan. 2024. Cross-Modal Consistency in Multimodal Large Language Models. doi:10.48550/arXiv.2411.09273 arXiv:2411.09273 [cs].

[50] Tinghui Zhu, Qin Liu, Fei Wang, Zhengzhong Tu, and Muhao Chen. 2024. Unraveling Cross-Modality Knowledge Conflicts in Large Vision-Language Models. doi:10.48550/arXiv.2410.03659 arXiv:2410.03659 [cs].

[51] Wanzheng Zhu and Suma Bhat. 2020. GRUEN for Evaluating Linguistic Quality of Generated Text. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, Trevor Cohn, Yulan He, and Yang Liu (Eds.). Association for Computational Linguistics, Online, 94–108. doi:10.18653/v1/2020.findings-emnlp.9