

ANSEL Photobot: A Robot Event Photographer with Semantic Intelligence

Dmitriy Rivkin, Gregory Dudek, Nikhil Kakodkar,
David Meger, Oliver Limoyo, Michael Jenkin, Xue Liu,
Francois Hogan

I. ABSTRACT

Our work examines the way in which large language models can be used for robotic planning and sampling, specifically the context of automated photographic documentation. Specifically, we illustrate how to produce a photo-taking robot with an exceptional level of semantic awareness by leveraging recent advances in general purpose language (LM) and vision-language (VLM) models. Given a high-level description of an event we use an LM to generate a natural-language list of photo descriptions that one would expect a photographer to capture at the event. We then use a VLM to identify the best matches to these descriptions in the robot’s video stream. The photo portfolios generated by our method are consistently rated as more appropriate to the event by human evaluators than those generated by existing methods.

II. INTRODUCTION

This paper presents ANSEL (Appropriate sNap SElection) Photobot, the world’s first semantically-aware robot photographer that can take photos across multiple domains starting only with high-level English task descriptions. It is implemented using publicly available language and language-vision models with no fine-tuning. Event photographers are expected to obey social conventions in their photos, and which photos are appropriate to take are highly dependent on the nature event and the activities that are likely to occur. For example, at the christening of a large ship, a key activity was historically the breaking of a champagne bottle on the bow of the ship (and, in fact, special extra-fragile bottles are available for this purpose). The ability to have robots act on these conventions has only been unlocked in the past couple of years due to advances in general purpose LMs [1], [2], [3] and VLMs [4], [5]. LMs contain much world knowledge and can perform common sense reasoning but operate in a abstract language space, while VLMs can ground the language to the robot’s sensory reality. Though their application to robotics is in the early stages, they are already beginning to unlock new robotics capabilities in the domains of manipulation [6], [7] and navigation [8].

The bulk of the prior work in this domain focuses on exploiting the language models’ understanding of the world of places and objects, while our work highlights their capabilities in the domain of social convention. We capture a video stream from a robot present at a party, then extract a fixed number of the most contextually appropriate images from it (Figure 1). We refer to this set of images as a “portfolio.” In order to evaluate our method, we created

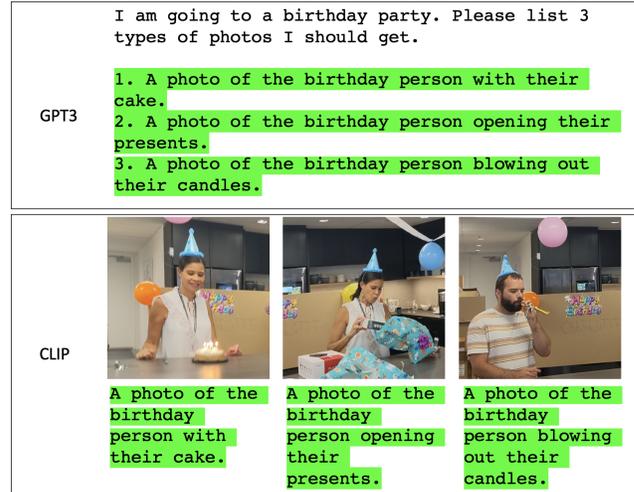


Fig. 1: Our method gets good ideas for photos from an LM (GPT3 [1], top), then finds them in the robot’s visual data stream using a VLM (CLIP [4], bottom). Green highlighting indicates text generated by GPT3. The mismatch between the rightmost caption and image illustrates some of the challenges involved in this approach, and is discussed further in the text below.

robot-centric video recordings of three social events, and generated 9-image portfolios of each using both our method and a modern video summarization technique (CA-SUM, [9]). We then asked workers to perform pairwise comparisons between portfolios, ranking which they believe is more appropriate given the event description. Our method was consistently rated as more appropriate than those created by CA-SUM.

The primary contributions of this work are:

- 1) We propose an approach to the development of robot event photographers with an unprecedented level of general-purpose semantic awareness.
- 2) We describe the design of a functioning system to embody this idea and we evaluate its performance.
- 3) We discuss some of the trade-offs and design considerations.

Since the design of all aspects of a novel functions system like this is necessarily complex and multi-layered, we focus here in the novel aspects of the overall design and use of an LM, and space limitations preclude a detailed description of the hardware and embodiment.

III. RELATED WORK

A. Robot photographers

Robot photographers can be used both to collect personal snapshots, or to document various events activities such as an ongoing parade or an environmental phenomenon [10]. When the problem includes the more comprehensive motion of documenting some phenomenon, the term “photo documenter” or “photo-documentarian” might be more apt (technically, the term “documentarian” also refers to photographer specializing in producing a factual record”).

Until now, work in the field of robot photography has been focused primarily on snapping good photos, where the “goodness” of the photos was evaluated individually either by volunteers on a Likert scale as in [11], [12], [13], or a with an existing photo composition quality model [14]. All of these works present hardware designs as well as the accompanying software.

Earlier work, such as [12] and [11], focus on the tasks of detecting people (e.g. using skin detection), navigating to them, and then composing the photograph using heuristics such as the “rule of thirds.” Later works employed more machine-learning focused techniques. For example, [14] used deep reinforcement learning to learn a template-matching policy to capture photos of people in desired poses. [13].

Finally, and perhaps most similar to our work in terms of experimental setup, Newbury *et al.* considered a robot that could navigate along a pre-programmed trajectory and collect large numbers of photographs of people [13]. These photos are then analyzed post-hoc to detect the best ones using a learned photo quality predictor. The authors of this paper note that this closely resembles the behavior of real modern-day photographers who have access to large memory cards.

Robot photographers have been considered in several contexts, from those that optimize image quality from an essentially fixed viewpoint, to those that track and follow a subject of interest. Applications for such robots range from those that might optimize portraiture, to automatically record complex procedures such as surgical procedures or military operations. Some authors have focused to the local actions of a robot photographer such as framing the shots and observing the rules associated with capturing a good composition [15]. A larger-scale class of problems relate to planning the trajectories of one or more cameras for visual coverage of some phenomenon [16], including maintaining visual content with the subject being recorded [17], [18], planning paths that allow an escaping photographic subject to be captured [19], or encircling a subject with a swarm of robot paparazzi [20], [21], [22].

B. Video Summarization

Video summarization is the problem of extracting summaries from videos, either in the form of a set of short clips or a collection of images [23], [24]. Summarization can be based on simple methods like uniform sampling, image content, geometric considerations (where it has some

connections to SLAM), and combinations of location, sensor data, and scene content [25], [26].

Google Inc released a commercial product called “Clips” that was a stationary camera that automatically selected when to take short video sequences based on on-device face recognition (note that Google’s AI-enabled Clips product is distinct from the CLIP software system from OpenAI that we use in this paper). Unfortunately, the lack of viewpoint variability and other pragmatic considerations limited the acceptance of the camera [27] and thus some innovators have even considered endowing it with mobility [28].

Query-guided video summarization, first defined in [29], is similar, except that a user query is also provided. For example, if the video was a tour of a national park, and the query was “water,” then the resultant summary should contain a representative subset of the different scenes of water in the video. Examples include [30], [29], [31].

While the video summarization task, and especially the query guided version, superficially seem quite similar to the robotic photography task as defined in our work, there is a key difference. Namely, our task requires inferring the appropriate photographs to get based on a high-level description, while query-guided image summarization is more object focused. For example, for a video of a birthday party, our task would ask “get all the pictures one would expect from a birthday party”, while query guided video summarization would ask “get a representative selection of all the pictures of cake.”

IV. PRELIMINARIES

Pretrained LMs and VLMs are key to our method. We consider approaches based on the Transformer [32] architecture, scaled up to very large model sizes trained on very large amounts of data and available for public use. We access GPT3 [1] through its web interface, and run CLIP [4] on our own server using publicly available checkpoints. In this section we briefly review each of these models, as well as CA-SUM, a generic video summarization technique which we use as baseline.

A. Transformers

Encoder and decoder models have emerged as a flexible approach for numerous language learning tasks [33], [34], allowing reasoning to occur in a latent embedding space rather than the complex tokens representing raw language. Concomitantly, attention has emerged as a key tool in the recurrent processing of text, allowing long-term correlations to be captured [35], [36]. The use of (masked) multi-head self-attention in both the encoder and decoder, known as Transformers [32], combine these advances, allowing attention to summarize information across both temporal scales and levels of abstraction. Transformers provide opportunities for parallelization, reducing training duration by several orders of magnitude while improving performance, which has been a key to scaling to internet-sized datasets, making Transformers the architecture of choice for modern semantic models.

B. GPT3

The scaling potential of Transformers allows training on language datasets containing billions of tokens. The GPT-2 [37] approach demonstrated the power of unsupervised training for multi-task learning of many related language tasks, with a model of 1.5B parameters, achieving unprecedented performance upon publication. GPT3 [1] further scaled-up GPT-2s architecture to 175B parameters by exploring few-shot training on these many tasks, training on the Common Crawl [38] dataset of roughly 1.0T words from the internet.

C. CLIP

The internet contains a wealth of image data accompanied by associated text, creating the opportunity for training paired vision and language models at scale. We use the particularly successful CLIP [4] approach for our work, which has been trained on the WebImageText (WIT) dataset, containing 400M images from 500K language queries. For scalable performance, CLIP jointly trains a visual encoder and a Transformer language model to optimize a metric-learning objective. In other words, CLIP is trained to embed text and images into a common embedding space using a dataset of image/caption pairs, where the loss is minimized by increasing the cosine similarity between image-caption pairs that appear in the dataset while minimizing those of random image-caption pairings. Although the WIT dataset contains poor label accuracy compared to previous vision-specific datasets, its large size, when utilized by CLIP’s modern VLM architecture, has produced a semantically aware visual understanding model which is state-of-the-art on many tasks.

D. CA-SUM

CA-SUM [23] is a recently proposed unsupervised video summarization technique which we use as a baseline to compare our results against. This method tries to reason about the uniqueness and diversity of the video frames and presents a technique for a concentrated attention mechanism represented by a block diagonal sparse attention matrix.

Given a set of video frames the method first converts these into feature representations using an off-the-shelf pre-trained model. These features form the inputs to a self-attention pipeline. The derived attention matrix A is used as input for the concentrated attention mechanism. This mechanism produces a block diagonal sparse matrix B that concentrates the information about the uniqueness and diversity of each frame within the block. The uniqueness u_t for each frame is represented by the entropy over the entire frame sequence i.e. the entropy of each row of the attention matrix A . While the diversity d_t is represented by the mean of the cosine similarity of each frame within this block to the ones that lie outside the block.

The final output of the network produces a set of frame-level scores that represents each frames’ importance. At inference, given a temporal segmentation of the video (obtained using the KTS algorithm [39]), the importance of each segment is inferred by averaging the scores of the frames

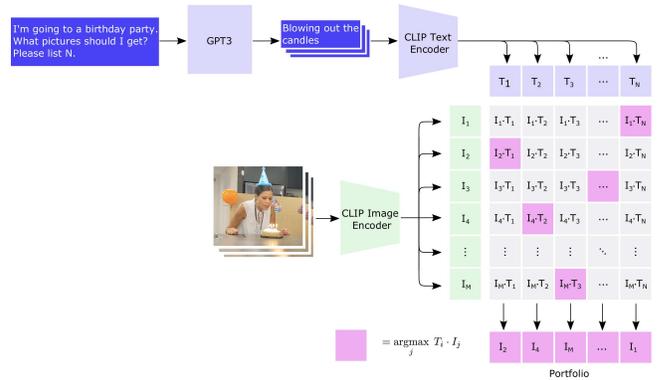


Fig. 2: Summary of our method. GPT3 is queried with a high level event description and asked to generate some phrases describing photographs that the photographer should try to capture. We then map these phrases and the images captured by the robot’s camera into a common embedding space using CLIP. Next, we compute the pairwise dot product between all the phrase and and all the image embeddings. Finally, we apply the argmax operation over images for each phrase in order to find which image best describes each phrase. These images are output as the portfolio.

within this segment. Finally, by limiting the length of the video summary to be 15% of the total length of video, the video summary is generated by solving the Knapsack problem.

V. METHOD

Our approach is to leverage GPT3’s world knowledge to extract textual task descriptions of the stereotypical photographs one would expect to see at an event given a prompt which describes the event at high level. An example of this approach in action is given in Figure 1. In that example, the phrases returned by GPT3 include specific objects and concepts that can easily be identified in photographs (“person,” “cake,” “candles,” “presents”). There are also activities which are often more challenging, but still possible, to recognize (“blowing out,” “opening”). Rather than parsing these concepts out explicitly and attempting to detect them, we simply embed the entire phrase output by GPT3 using CLIP’s text encoder. Next, we also embed all the frames of the collected video using CLIP’s image encoder. We then compute pairwise cosine similarities between each phrase embedding and each image embedding. Finally, we construct our portfolio by, for each GPT3-generated phrase, picking the image with the highest cosine similarity. This process is summarized in Figure 2. The rest of this section describes some nuances of the implementation.

A. Prompt engineering

The term “prompt engineering” refers to the design of prompts to the language models to get them to output reasonable results. In the case of our approach this is doubly relevant since we have two language-consuming models, GPT3 and CLIP.

```

<fullQuery> ::= <example><query>
<query> ::= <intro> <demand> <constraint><EOL>
<intro> ::= "I am going to <eventName>."
<demand> ::= "Please name" <quantity> "different types
of photos of the event that I should get."
<constraint> ::= "Please focus on the content
of the photos rather than their composition."
<example> ::= <exampleQuery><EOL><resultList> | ""
<exampleQuery> ::= <exampleIntro><demand><constraint>
<exampleIntro> ::= "I am going to <exampleEventName>"
<photoIdeaList> ::= <photoIdea>|<photoIdea><photoIdeaList>

```

Fig. 3: BNF grammar for constructing GPT3 queries which create good ANSEL event plans. A photoIdea is an individual phrase such as “A photo of the birthday cake”. eventName is the name of the event you would like suggestions for, and exampleEventName is that for a different event that is used to construct the example. The same exampleEvent and photoIdeaList can be reused for priming across different values of eventName.

1) *CLIP*: In prior work, researchers found that it was possible to use prompt engineering to employ CLIP as an Imagenet classifier. There, the authors found that prompt ensembling significantly increased classification performance. For example, instead of just comparing image embeddings to the text embeddings for “a cat”, they would average the embeddings for “A painting of a cat,” “A cartoon cat,” “A realistic cat”, etc. [4]. We experimented briefly with this kind of ensembling, using the same ensemble of modifiers as [4], but did not find that it improved performance on our tasks. Instead, we directly fed the outputs of GPT3 into the CLIP text encoder.

We did find, however, that query phrases including descriptions of the composition (most commonly “close up” and “wide shot”) tended to lead to particularly poor photo selections. As such, we rejected any phrase sets that GPT3 proposed that had any of the terms [close-up, closeup, close up, wide shot] within them. When a phrase set was rejected, another would be sampled. We also explicitly asked GPT3 to focus on content over composition (see Figure 3 for details).

2) *GPT3*: We use Backus–Naur form (BNF) [40] to facilitate and formalize the prompt engineering for our approach, but note that prompt engineering requires some degree of human intercession to optimize queries for different domains. For the event domains evaluated in this paper, our prompts have the form presented in in Figure 3. This form creates a concrete specification for the desired outputs and allows for the priming of GPT3 using example input-output pairs, following the common pattern used in works such as [5].

Despite its apparent simplicity, obtaining good results can be elusive and depends on queries that are both precise enough to be effective, yet general enough to be understood by the language model and provide actionable results. We found that using the words “photos of the event” rather than just “photos” caused GPT3 to more reliably output the proper kinds of outputs. Without this, it would sometimes return different types of photos (landscape, HDR, black and white, etc). We also found that it was able to reliably generate lists of the requested length, and would always return responses

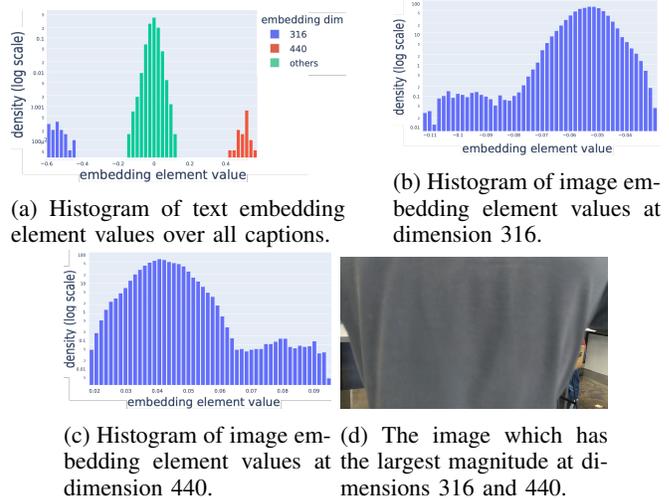


Fig. 4: Two dimensions of the CLIP text embeddings have very large values, which tends to favor particularly bad images.

in the enumerated format shown in Figure 1, which made the outputs easy to parse. Without explicitly specifying the length of the list, the model tends to only output a few (2-3) types of photographs. In addition, based on the observation that CLIP often provides highly undesirable results when presented with directives related to composition, we ask the model to focus on content.

Finally, we observe that GPT3 behaves somewhat stochastically – regenerating answers using the same prompt often leads to different phrasing. For example, it sometimes uses the term “the birthday person” and other times it will be “the birthday boy/girl.”

B. Dropping large embedding dimensions

During our experiments, we noticed that CLIP embeddings have two dimensions with values which are much larger than those of all the others. These are 316 and 440 for all text embeddings (Figure 4a), and 413 and 644 for all image embeddings (not shown). These large dimensions have a significant impact on the results; frames which have particularly large values on embedding dimensions 316 and 440 will have much higher cosine similarities to *all* query phrases. Figures 4b and 4c show the distribution of the values of the image embedding dimensions 316 and 440 respectively. In both of these we observe a long tail of values that are significantly higher than the mean, meaning these frames should be selected particularly often. When we examine these frames with extremely large values of these two dimensions we find two interesting things. First, the same frame has the highest magnitude of both dimensions 316 and 440. Second, this frame, shown in Figure 4d, is of extremely poor quality.

We have not seen these issues previously reported, and they are certainly not addressed in the original work which presented CLIP. We hypothesise that this quirk has not had a significant impact on the majority of other works using

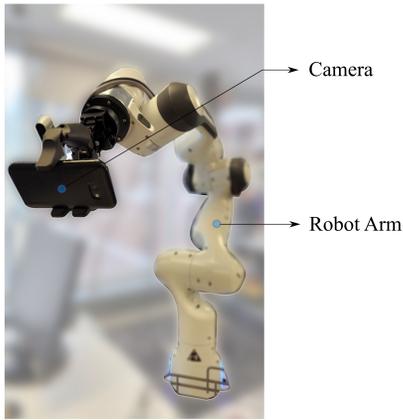


Fig. 5: ANSEL hardware embodiment: Robot arm holding the cell phone camera used to record data.

CLIP because they tend to work with more highly “curated” datasets (like web data and youtube videos) than data coming from real-world robots, so they encounter fewer, or even none, “garbage” frames like those in 4d. We found that zeroing any dimensions with magnitude greater than 0.3 and then re-normalizing the embeddings was an effective method for reducing the dominance of these “garbage” frames.

C. Faces

In view of our objective of finding satisfying photographs, we used human subjects and “A/B” trials as the key objective assessment mechanism for our approach vis-a-vis alternative algorithmic approaches. In pre-existing automated photograph acquisition systems, although they are generally not driven by semantic criteria, let alone general-purpose ones, attention to human faces is a common attribute. For example, the Google Clips camera uses human faces to (sometimes) trigger video clip recording. With this in mind, we crop the photographs we acquire around human faces that appear in the image. If no human faces are detected, we filter out the image.

Faces are detected using the Multi-Task Cascaded Convolutional Neural Networks (MTCNN) [41] algorithm that detects faces using a cascade of $h=$ three convolutional neural networks, implemented using Tensorflow. The bounding box for the ensemble of all faces is extracted, and expanded by a small margin (10 *per cent*) and extended downward to the bottom of the image plane (size of detector for the extent of bodies is used). The image is then zoomed and clipped around this bounding box to provide a more human-centric view of the scene.

VI. EVALUATION

To validate and evaluate our approach, we deployed our ANSEL Photobot system using a Franka Emika robot arm holding a cell phone used as a video recording device (shown in Fig. 5). Using an off-the-shelf cell phone camera allowed us to achieve high image quality and a degree of portability

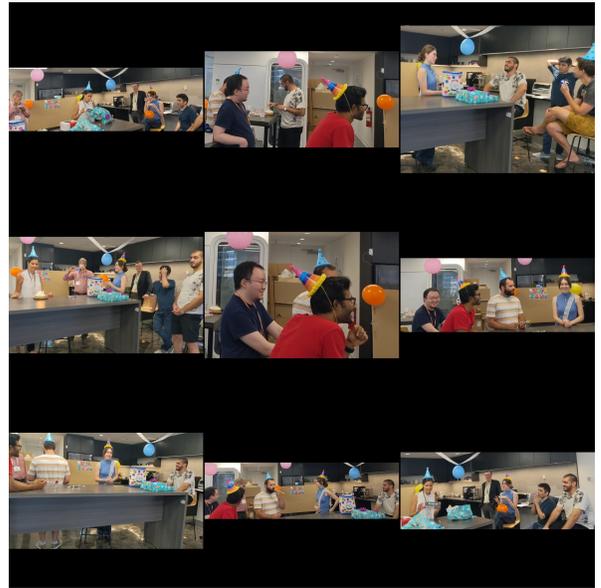


Fig. 6: Portfolio of photos from a birthday party scenario generated using ANSEL.

and hardware independence, but notably core processing operations are performed off board (ie. in “the cloud”).

In this study, we had the robot collect video of 3 simulated social events, a birthday party, a wine tasting event, and a painting class. These descriptions (“a birthday party”, “a wine tasting event”, and “a painting class”) were used as eventNames in the query construction. Note that while we experimented briefly with the use of examples in priming GPT3, and preliminary results seemed promising, the results presented in this work omit the use of these priming examples.

The events were “simulated” in the sense that these were actual events, but they were convened for the purposes of this paper in an office environment by people pretending to experience social activities (although they did, in fact, experience the events in their full form including actually making paintings, eating cake, and learning about tastes). For each event, we generate a photographic summary of nine photos¹, stitched together into a collage (an example is shown in Figure 6, using both our method and a recent generic video summarization technique (CA-SUM) which serves as a baseline [23]).

We slightly adapted CA-SUM to match our problem setup. The original method selects short video clips (usually 0.5-2 seconds) and limits the total duration of these clips to 15% of the original video length. We change this to limit the number of clips to 9, and selected the frame from the center of each clip. The model we used was trained on the TV-Sum dataset [29]. This dataset contains videos of similar length to our (15-30 minutes), along with user annotations which were used to perform model selection. We apply the same

¹CLIP hyperparameters: model=ViT-L/14@336px. GPT3 hyperparameters: model=text-davinci-002, temperature=0.7, max_tokens=2000, top_p=1, frequency_penalty=0, presence_penalty=0

face-cropping procedure to the frames selected by CA-SUM as we do for our method.

We then surveyed 10 individuals to ascertain which collage was more semantically relevant to the each event type. We simultaneously presented them with collages for each method, and asked them to pick their preferred one. In order to prime the individuals to focus on semantics over composition, we first asked them to generate a list of nine photo types that would be appropriate to the event type before doing the evaluation. We also showed the users GPT3’s suggestions after they had created their own and asked them to score each set of phrases (both their own and the ones from GPT3) on a scale of 0-10.

VII. RESULTS

All participants found the quality of the targets for photo subjects to be of high quality, often preferring them in retrospect over labels the subject themselves has suggested². After seeing the GPT3 labels, subject’s scoring of their own labels had an average value of 7.4/10 while their scores for the labels from GPT3 were 7.0/10 (see Figure 8 for more details). Clearly, GPT3 provided photographic directives at a high semantic level that were as good as those generated by the humans, and in the opinion many subjects even better than their own.

ANSEL outperforms CA-SUM in 2/3 events (birthday and paint), while tying it on the wine tasting, as illustrated by Figure 7. This confirms that ANSEL is capable of identifying the stereotypical, appropriate photographs associated with different types of social events. A potential explanation as to why ANSEL wins so decisively in two events and yet ties the third might be found in Figure 8. In the case of the wine tasting, on average people actually tended to prefer GPT3’s captions to their own, which may indicate that they had weaker priors about the expected photos for this event type. Weaker expectations on the photo types could explain the apparent indifference to summarization method encountered in the wine event.

VIII. CONCLUSION

We have described an approach that uses an impressive level of implicit semantic knowledge embedded in a large language model to generate a robot planning systems that operates at a highly abstract level. We introduced ANSEL Photobot, a robot photographer that can identify key events from a video stream using a combination of LM and VLM models by adapting the semantic specifications to the event of interest. Results show that our approach consistently generates photo portfolios that are consistently rated as more appropriate than video summarizing baselines by human evaluators. In the context of photographer/documentarian robotics specifically, we believe that there are opportunities

²Full list GPT3-generated suggestions for wine tasting event: People discussing the wine. The different types of wine on display. People mingling and networking. The venue of the event. The buffet or food that is being served. The decorations or theme of the event. The staff serving the wine. The guests enjoying themselves

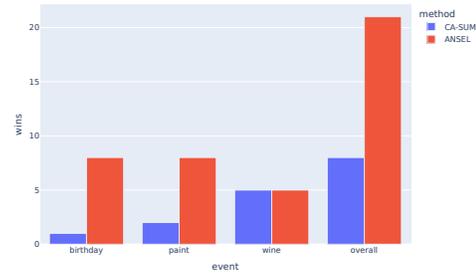


Fig. 7: Number of wins for ANSEL and CA-SUM on each of the events, as well as the aggregated value combining all events. The ANSEL algorithm is consistently preferred in 2/3 events, and ties in only one case.

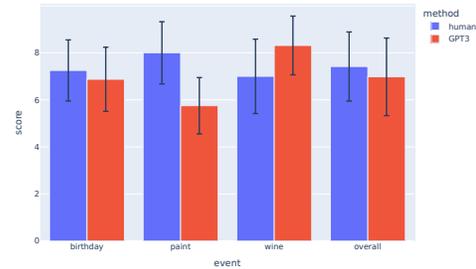


Fig. 8: Scores assigned by users to their own plan sets (human) and those generated by GPT3. Error bars represent standard deviation. Users were told to assign scores of 0-10 to each of the caption sets, so both systems were ranked as very good and also statistically indistinguishable.

to improve the quality of the photo portfolios by incorporating image quality metrics as well (such as symmetry, affect, sharpness). These methods are beyond the scope of this paper as such enhancements would have made it more difficult to compare the results of the semantic filtering against baseline methods.

While it appears that GPT3 can produce near human level plans for the task examined in this paper, the challenge of grounding these concepts in the visual reality of robots remains significant. In Fig. 1 for example, it is clear that the subject is not “the birthday person blowing out their candles.” However, the person in the image is clearly blowing into a tube, and in fact it might be fair to say that this image matches the concept of “blowing” much better than a person leaned over some candles. Furthermore, the person is wearing a party hat, and the words “happy birthday” are clearly visible. It appears that CLIP is not capable of balancing this opulence of “blowing-ness” and “birthday-ness” against the dearth of candles. Assuming that VLMs continue to evolve as rapidly as they have in the last several years, we believe that these kinds of errors will soon become a thing of the past, and that LMs and VLMs will have a transformative impact on robotics.

REFERENCES

- [1] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, *et al.*, “Language models are few-shot learners,” in *Advances in Neural Information Processing Systems*, H. Larochelle, M. Ranzato,

- R. Hadsell, M. Balcan, and H. Lin, Eds., vol. 33. Curran Associates, Inc., 2020, pp. 1877–1901.
- [2] R. Thoppilan, D. De Freitas, J. Hall, N. Shazeer, A. Kulshreshtha, H.-T. Cheng, *et al.*, “Lamda: Language models for dialog applications,” 2022. [Online]. Available: <https://arxiv.org/abs/2201.08239>
 - [3] A. Chowdhery, S. Narang, J. Devlin, M. Bosma, G. Mishra, A. Roberts, *et al.*, “Palm: Scaling language modeling with pathways,” 2022. [Online]. Available: <https://arxiv.org/abs/2204.02311>
 - [4] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, *et al.*, “Learning transferable visual models from natural language supervision,” 2021. [Online]. Available: <https://arxiv.org/abs/2103.00020>
 - [5] J.-B. Alayrac, J. Donahue, P. Luc, A. Miech, I. Barr, Y. Hasson, *et al.*, “Flamingo: a visual language model for few-shot learning,” 2022. [Online]. Available: <https://arxiv.org/abs/2204.14198>
 - [6] M. Ahn, A. Brohan, N. Brown, Y. Chebotar, O. Cortes, B. David, *et al.*, “Do as i can and not as i say: Grounding language in robotic affordances,” in *arXiv preprint arXiv:2204.01691*, 2022.
 - [7] W. Huang, F. Xia, T. Xiao, H. Chan, J. Liang, P. Florence, *et al.*, “Inner monologue: Embodied reasoning through planning with language models,” 2022. [Online]. Available: <https://arxiv.org/abs/2207.05608>
 - [8] D. Shah, B. Osinski, B. Ichter, and S. Levine, “Lm-nav: Robotic navigation with large pre-trained models of language, vision, and action,” 2022. [Online]. Available: <https://arxiv.org/abs/2207.04429>
 - [9] E. Apostolidis, G. Balaouras, V. Mezaris, and I. Patras, “Summarizing videos using concentrated attention and considering the uniqueness and diversity of the video frames,” in *Proceedings of the 2022 International Conference on Multimedia Retrieval*, 2022, pp. 407–415.
 - [10] T. Manderson, J. Li, N. Dudek, D. Meger, and G. Dudek, “Robotic coral reef health assessment using automated image analysis,” *Journal of Field Robotics*, vol. 34, no. 1, pp. 170–187, 2017.
 - [11] M. Zabaraukas and S. Cameron, “Luke: An autonomous robot photographer,” in *2014 IEEE International Conference on Robotics and Automation (ICRA)*, 2014, pp. 1809–1815.
 - [12] Z. Byers, M. Dixon, K. Goodier, C. Grimm, and W. Smart, “An autonomous robot photographer,” in *Proceedings 2003 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS 2003) (Cat. No.03CH37453)*, vol. 3, 2003, pp. 2636–2641 vol.3.
 - [13] R. Newbury, A. Cosgun, M. Koseoglu, and T. Drummond, “Learning to take good pictures of people with a robot photographer,” in *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2020, pp. 11 268–11 275.
 - [14] H. Kang, J. Zhang, H. Li, Z. Lin, T. Rhodes, and B. Benes, “Lerop: A learning-based modular robot photography framework,” 2019. [Online]. Available: <https://arxiv.org/abs/1911.12470>
 - [15] M. Dixon, C. M. Grimm, and W. D. Smart, “Picture composition for a robot photographer,” 2003.
 - [16] E. Bourque and G. Dudek, “Viewpoint selection-an autonomous robotic system for virtual environment creation,” in *Proceedings. 1998 IEEE/RSJ International Conference on Intelligent Robots and Systems. Innovations in Theory, Practice and Applications (Cat. No. 98CH36190)*, vol. 1. IEEE, 1998, pp. 526–532.
 - [17] F. Shkurti and G. Dudek, “Maximizing visibility in collaborative trajectory planning,” in *2014 IEEE International Conference on Robotics and Automation (ICRA)*, 2014, pp. 3771–3776.
 - [18] F. Shkurti, N. Kakodkar, and G. Dudek, “Model-based probabilistic pursuit via inverse reinforcement learning,” in *2018 IEEE International Conference on Robotics and Automation (ICRA)*, 2018, pp. 7804–7811.
 - [19] F. Shkurti and G. Dudek, “Topologically distinct trajectory predictions for probabilistic pursuit,” in *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2017, pp. 5653–5660.
 - [20] M. Jenkin and G. Dudek, “The paparazzi problem,” in *Proceedings. 2000 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS 2000)(Cat. No. 00CH37113)*, vol. 3. IEEE, 2000, pp. 2042–2047.
 - [21] P. Brisset, A. Drouin, M. Gorraz, P.-S. Huard, and J. Tyler, “The paparazzi solution,” in *MAV 2006, 2nd US-European competition and workshop on micro air vehicles*. Citeseer, 2006, pp. pp–xxxx.
 - [22] S. Baldi, N. Maric, R. Dornberger, and T. Hanne, “Pathfinding optimization when solving the paparazzi problem comparing a and dijkstra’s algorithm,” in *2018 6th International Symposium on Computational and Business Intelligence (ISCBI)*, 2018, pp. 16–22.
 - [23] E. Apostolidis, G. Balaouras, V. Mezaris, and I. Patras, “Summarizing videos using concentrated attention and considering the uniqueness and diversity of the video frames,” ser. ICMR ’22. New York, NY, USA: Association for Computing Machinery, 2022, pp. 407–415. [Online]. Available: <https://doi.org/10.1145/3512527.3531404>
 - [24] J. Fajtl, H. S. Sokeh, V. Argyriou, D. Monekosso, and P. Remagnino, “Summarizing videos with attention,” 2018.
 - [25] A. Flint, C. Mei, I. Reid, and D. Murray, “Growing semantically meaningful models for visual slam,” in *2010 IEEE computer society conference on computer vision and pattern recognition*. IEEE, 2010, pp. 467–474.
 - [26] Y. Girdhar and G. Dudek, “Modeling curiosity in a mobile robot for long-term autonomous exploration and monitoring,” *Autonomous Robots*, vol. 40, no. 7, pp. 1267–1278, 2016.
 - [27] C. Bonnington, “Google clips smart camera isn’t smart enough, but its aims are still worth considering,” *Slate [Online]*, Feb, vol. 27, pp. 1–3, 2018.
 - [28] J. Pierce, “Roomba+ clips cam: Exploring unpredictable autonomy in everyday smart systems,” in *Companion Publication of the 2020 ACM Designing Interactive Systems Conference*, 2020, pp. 317–320.
 - [29] A. Sharghi, J. S. Laurel, and B. Gong, “Query-focused video summarization: Dataset, evaluation, and a memory network based approach,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.
 - [30] M. Narasimhan, A. Rohrbach, and T. Darrell, “Clip-it! language-guided video summarization,” in *Advances in Neural Information Processing Systems*, M. Ranzato, A. Beygelzimer, Y. Dauphin, P. Liang, and J. W. Vaughan, Eds. Curran Associates, Inc., pp. 13 988–14 000.
 - [31] S. Xiao, Z. Zhao, Z. Zhang, Z. Guan, and D. Cai, “Query-biased self-attentive network for query-focused video summarization,” *IEEE Transactions on Image Processing*, vol. 29, pp. 5889–5899, 2020.
 - [32] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, *et al.*, “Attention is all you need,” in *Advances in Neural Information Processing Systems*, I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, *et al.*, Eds. Curran Associates, Inc.
 - [33] R. Jozefowicz, O. Vinyals, M. Schuster, N. Shazeer, and Y. Wu, “Exploring the limits of language modeling,” 2016.
 - [34] Y. Wu, M. Schuster, Z. Chen, Q. V. Le, M. Norouzi, W. Macherey, *et al.*, “Google’s neural machine translation system: Bridging the gap between human and machine translation,” *CoRR*, vol. abs/1609.08144, 2016. [Online]. Available: <http://arxiv.org/abs/1609.08144>
 - [35] D. Bahdanau, K. Cho, and Y. Bengio, “Neural machine translation by jointly learning to align and translate,” 2014. [Online]. Available: <https://arxiv.org/abs/1409.0473>
 - [36] Y. Kim, C. Denton, L. Hoang, and A. M. Rush, “Structured attention networks,” in *International Conference on Learning Representations*, 2017. [Online]. Available: <https://openreview.net/forum?id=HkeE0NvqJg>
 - [37] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, and I. Sutskever, “Language models are unsupervised multitask learners,” 2019.
 - [38] C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, *et al.*, “Exploring the limits of transfer learning with a unified text-to-text transformer,” *Journal of Machine Learning Research*, vol. 21, no. 140, pp. 1–67, 2020. [Online]. Available: <http://jmlr.org/papers/v21/20-074.html>
 - [39] D. Potapov, M. Douze, Z. Harchaoui, and C. Schmid, “Category-specific video summarization,” in *European conference on computer vision*. Springer, 2014, pp. 540–555.
 - [40] B. JW, “c the syntax and semantics of the proposed international algebraic language of the zurich acm-gamm conference,” in *Proceedings of the International Conference of Information Processing UNESCO Paris June*, 1959.
 - [41] K. Zhang, Z. Zhang, Z. Li, and Y. Qiao, “Joint face detection and alignment using multitask cascaded convolutional networks,” *IEEE Signal Processing Letters*, vol. 23, no. 10, pp. 1499–1503, Oct 2016.