



Predicting and Optimizing Human Behavior and Communication

By

Yaman K Singla

With the supervision of

Prof Rajiv Ratn Shah, IIIT Delhi
Prof Changyou Chen, State University of New York at
Buffalo

Submitted To:

Indraprastha Institute of Information Technology Delhi
State University of New York at Buffalo

January, 2024

If the brain were so simple we could understand it, we would be so simple we couldn't. - Emerson Pugh

Also,

*The sad fact is that we're actually much better at planning the flight path of an interplanetary rocket (rocket science) than we are at managing the economy, merging two corporations, or even predicting how many copies of a book will sell (behavior prediction). So why is it that rocket science **seems** hard, whereas problems having to do with people - which arguably are much harder - **seem** like they ought to be just a matter of common sense?* - Duncan J. Watts

but

Nothing in Nature is random. A thing appears random only through the incompleteness of our knowledge (ignorance). - Baruch Spinoza

and

Timendi causa est nescire. (Ignorance is the cause of fear.)

and

What would life be if we had (fear and) no courage to attempt anything?
- Vincent Van Gogh

Acknowledgments

कार्यदोषोपहतस्वभावः
पृच्छामि त्वां धर्मसमूढचेताः ।
यच्छ्रेयः स्यान्निश्चितं ब्रूहि तन्मे
शिष्यस्तेऽहं शाधि मां त्वां प्रपन्नम् । | BG 2:7||

मयि सर्वाणि कर्माणि संन्यस्याध्यात्मचेतसा ।
निराशीर्निर्ममो भूत्वा युध्यस्व विगतज्वरः । | BG 3:30||

नैव किञ्चित्करोमीति युक्तो मन्येत तत्त्ववित् ।
पश्यञ्शृणवन्स्पृशञ्जिघ्रन्नशननगच्छन्स्वपञ्चश्वसन् ॥
प्रलपन्विसृजन्मृष्टन्मिष्टन्मिष्टन्पि ।
इन्द्रियाणीन्द्रियार्थेषु वर्तन्त इति धारयन् । | BG 5:8-9||

This is an attempt to capture and thank those who have shaped my journey. Albeit, due to indirect and latent relations, this list will remain non-exhaustive despite my best attempts, it is still an attempt worth making.

In no particular order (with names of the organizations where I met these giants): Rajiv Ratn Shah (IIIT-D), Changyou Chen (SUNY-Buffalo), Ranjeeta Rani (GMPS), Roger Zimmerman (National University of Singapore), Debanjan Mahata (Bloomberg), Jessy Junyi Li (University of Texas at Austin), Balaji Krishnamurthy (Adobe MDSR), Jayakumar Subramanian (Adobe MDSR), Amanda Stent (Bloomberg), Anil Seth (FIITJEE), Dhruva Sahrawat (IIIT-D), Yifang Yin (National University of Singapore), Mika Hama (Second Language Testing Institute), Payman Vafaee (Columbia University), Pankaj Bansal (Adobe), Mohit Srivastava (Adobe), Gaurav Jain (Adobe), Shubham Yadav (NSIT), Rohit Jain (NSIT), Mohd Khwaja Salik (NSIT), Pratham Nawal (NSIT), Mayank Singh (NSIT), Somesh Singh (BITS-Pilani Goa), Aanisha Bhattacharyya (Adobe MDSR), Varun Khurana (IIIT-D), Rita Yadav (GMPS), Prabha Sinha (GMPS), Neeta Pandit (GMPS), Geetha Nair (GMPS), Swami Sarvapriyananda (Ramakrishna Mission), and finally my parents and my brother, Aman Singla.

Hopefully, I can return whatever I have gathered from these people back to the earth.

Contents

| | |
|--|-----------|
| Acknowledgments | ii |
| 1 Introduction: The Two Cultures of Social Science | 1 |
| 2 Explaining Behavior: Persuasion Strategies | 6 |
| 2.1 Related Work | 10 |
| 2.2 Generic Taxonomy of Persuasion Strategies | 12 |
| 2.3 Persuasion Strategy Corpus Creation | 13 |
| 2.3.1 Persuasion Strategy Dataset For Image Advertisements | 13 |
| 2.3.2 Persuasion Strategy Dataset For Video Advertisements | 17 |
| 2.4 Modeling: Persuasion Strategy Prediction | 18 |
| 2.4.1 Modelling Persuasion Strategy For Image Advertisements | 18 |
| 2.4.2 Modelling Persuasion Strategy For Video Advertisements | 26 |
| 2.5 Conclusion | 41 |
| 3 Content and Behavior Models | 42 |
| 4 Generating Content Leading to Optimal Behavior | 43 |
| Bibliography | 60 |

Chapter 1

Introduction: The Two Cultures of Social Science

Communication includes all of the procedures by which one mind may affect another [1]. This includes all forms of expression, such as words, gestures, speech, pictures, and musical sounds. Communication can be seen as being composed of seven parts (Fig. 1.1): (the communicator, message, time of message, channel, receiver, time of receipt, and effect). Different fields deal with different parts of communication. I will give a broad overview of these fields in the upcoming paragraphs, but two streams have emerged broadly in behavioral sciences: explanation and prediction of behavior (receiver effect) [2–4].

Historically, social scientists have sought explanations of human behavior that can provide interpretable causal mechanisms behind human functioning. A few prominent examples are Milgram's [5] and Asch's [6] experiments on persuasion, explaining the causal mechanism of obedience to authority. The approach of theorizing has worked in physical sciences where the data is plentiful, and theories make unambiguous predictions but have not been too successful in *predicting* social outcomes in behavioral sciences [7–9]. In fact, many studies have shown that expert human opinions fare similar to non-experts (*e.g.*, predicting economic and political trends [8] and societal change: [9]), and the opinion of non-expert population is roughly the same as a random coin toss in predicting behavior (*e.g.*, predicting cascades [10] or image memorability [11]). At the same time, causal mechanisms have their own merits; most notably, they help decision-makers (often humans) to make intuitive sense of the situation and make their next decision based on it.

In parallel, due to the availability of human behavior data at scale, researchers in machine learning are showing a growing interest in traditionally social scientific topics, such as messaging strategies leading to persuasion [12–15], information diffusion [16, 17], and most importantly, prediction

and predictability of human behavior [18, 19]. Machine learning approaches bring with them the culture of (training and) testing their models on large real-world datasets and pushing the state-of-the-art in terms of predictive accuracies; at the same time, often, ML approaches can only be operated as black boxes with no direct mechanism to explain predictions [20, 21].

In the prediction community, different subfields have emerged dealing with the different parts of the problem of optimization of human behavior. For instance, advertisement personalization studies how to optimize (choose) *receiver* for a given message [22], and recommendation systems study how to *choose content* from a set of pre-decided contents for a given receiver to elicit a certain effect [23]. A popular problem within the prediction community is the effect prediction problems, for example, clickthrough (CTR) prediction [24], Twitter cascade prediction [16, 17], sales prediction [18, 25], content memorability prediction [26–28], *etc.* There are also works to optimize the time of the message to elicit certain effect [28, 29]. Therefore, we see that all the factors of communication are studied independently in their own light with the aim of achieving the desired effect.

Effect (or behavior) over a content can also enable us to understand about the content, the communicator, the receiver, or the time. Therefore, efforts have also been made to extract information about the content itself from the behavior it generates. For instance, using keystroke movements [30] and eye movements to improve natural language processing [31, 32]. Similarly, the fields of human alignment and reinforcement learning with human feedback (RLHF) try to use human behavioral signals of likes, up-votes, downloads, and annotations of a response's helpfulness to improve content generation - both text [28, 33–36] and images [37–40].

In the more traditional social science and computational social science cultures, research is carried out to discover causal effects and model them. For instance, propaganda and mass communication studies [41–44] try to understand the culture, time, authors, recipients in a non-invasive manner using the messages exchanged, and persuasion studies [45, 46] where the persuasion strategy present in the content is identified and correlated with (un)successful efforts of persuasion.

A common theme that runs through both research cultures in behavioral sciences is the intent to control behavior. Explanation and prediction are intermediate steps to control and hence optimize behavior. Optimizing behavior means to fulfill the communicator's objectives by controlling the other six parts of the communication process (Fig. 1.1). Due to the problem space being large, the solution needs a general understanding of human behavior as opposed to being domain-specific.

The characteristic that marks the digital age is the prevalence of human behavioral data in huge repositories. This data is *big* (allowing to model heterogeneity), *always-on* (allowing to look in the past as well as live measurements), observational (as opposed to reactive), but also *incomplete*

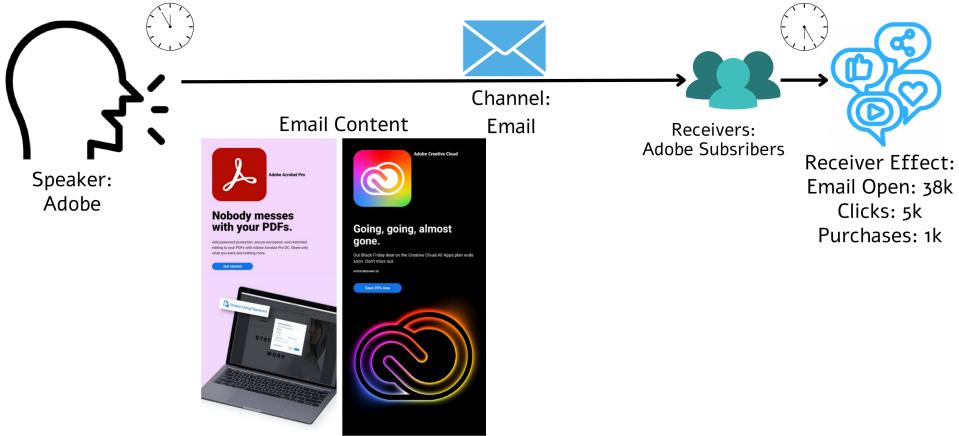


Figure 1.1: Communication process can be defined by seven factors: Communicator, Message, Time of message, Channel, Receiver, Time of receipt, and Effect. Any message is created to serve an end goal. For marketers, the end goal is to bring in the desired receiver effect (behavior) (like clicks, purchases, likes, and customer retention). The figure presents the key elements in the communication pipeline - the marketer, message, channel, receivers, and finally, the receiver effect.

(does not capture all that is happening everywhere everytime in a single repository) and *algorithmically confounded* (generated as a byproduct of an engineering process with a goal) [20]. While the predictive culture has tried to make use of some of this data in the form of social media datasets like Twitter [47, 48] and Instagram [49], Google trends [18, 50], Wikipedia [51–53], shopping websites [54, 55] and other data sources [19, 56, 57], these efforts are limited, in the sense of being dependent on one or a few chosen platforms, able to answer a limited set of questions, and restricted by access to private data. We want a model that can understand (predict and explain) *human behavior in general* as opposed to modeling a particular effect (retweet prediction) on a particular platform (*e.g.* Twitter) for a certain type of users. This problem carries parallels with the problem being solved in the natural language processing (NLP) community, where supervised models in NLP are limited by the amount of supervision available and being able to answer one question (for which the supervised model was trained). The problem was solved by developing Large Language Models (LLMs), which are general purpose models capable of *understanding language*, and hence can solve natural language tasks like sentiment analysis, question answering, email generation, and language translation in zero-shot (*i.e.* without needing any explicit training for that task) [58–62].

Similarly, how do we develop a model capable of understanding behavior *in general*? With the intent to answer this question, we take motivation



Figure 1.2: Levels of content analysis. The figure lists tasks and their sample outputs arranged in a hierarchy [1]. This is roughly based on levels of language.

from LLMs, where the idea is to train a model on a data-rich task. The task chosen to train LLMs is the next-word prediction, and the dataset is the text collected from the entire internet. The next-word prediction task is a data-rich task that can be trained on the huge text repositories from the internet. The intuition is that two approaches have always worked for neural networks: larger model sizes and more data for training [58, 60, 61, 63]. Going from a few million tokens of text [60, 63] to a trillion tokens [59, 64] leads to an increase in the transfer learning capability leading to performance improvements over a wide variety of natural language tasks.

The digital revolution has provided us with huge repositories of data. We leverage the human behavior repositories available on the internet for this general-purpose human behavior model. The format of this data is the general communication model shown in Fig. 1.1 consisting of communicator, message, time of message, channel, receiver, time of receipt, and effect. Due to the incomplete nature of behavioral repositories, all the factors are usually not always available. However, a subset is always available, and we show that the data scale, along with a large model, helps make a general behavior understanding model [65]. We call this model, Large Content and Behavior Model (LCBM). We show that LCBM can predict behavior, explain it, and generate a message to bring about certain behavior [28, 40, 65].

Are general LLMs unable to solve behavioral problems? A question that arises is whether LLMs, which already learn trillions of text tokens, are able to understand and predict behavior. We investigate that question over several large models, including GPT-3.5 [59], GPT-4 [66], Llama-13 B and Llama-7B [64], and find that they do not seem to have any behavioral capabilities. The reason for this is that large language models only include one factor (message) out of the 7-factor communication model (Fig. 1.1) while considering other parts as “noise” (for instance, see [67, 68]). This systematic purge of communicator, receiver, channel, time, and, most importantly, behavior causes the models not to develop any behavioral capabilities (Level-C of Shannon and Weaver [1]). As an example, Llava [69], a recent large language and vision model (LVM) trained by connecting a vision encoder with a language model, shows that after training on a few hundred thousand instructions, the language model can now “see”, and is able to answer

questions on the images. However, the questions all lie in the first two levels of content analysis shown in Fig. 1.2. The reason is that the instructions used to align the image encoder with the downstream LLM all lie in the first two levels while ignoring the last two. In the upcoming chapters, we explore how we can train a general behavior model and how including back the other factors of communication in training data help in understanding human behavior.

Outline for the upcoming chapters: Following the two traditions of social sciences, in this work, in Chapter-2, I start with a more traditional approach to behavior explanation, where I cover the first works on extracting persuasion strategies in advertisements (both images and videos) [13, 15]. The contributions of these works include constructing the largest set of generic persuasion strategies based on theoretical and empirical studies in marketing, social psychology, and machine learning literature and releasing the first datasets to enable the study and model development for the same. These works have been deployed to understand the correlation between the kinds of marketing campaigns and customer behavior measured by clicks, views, and other marketing key performance indicators (KPIs).

Following this, in Chapter-3, I delve into the more modern approach of behavior prediction and leveraging the huge repositories of behavior data available. First, we propose models to integrate behavior with relatively smaller language models like BERT [58], and show that the resultant models can understand content better than the base models [32]. Then, we propose an approach to integrate behavior and content together as part of a single model. We call these models Large Content and Behavior Models (LCBM) [65]. We show that these models can predict and explain behavior. Next, in Chapter-4, we show that using these models we can generate messages to elicit certain behavior resulting in behavior optimization. We show this both for the domain of text, by taking the illustrative case of memorability and generating content that is more memorable in long-term [28], and images, by generating images that are more performant, *i.e.*, can result in more likes and downloads [40].

Chapter 2

Explaining Behavior: Persuasion Strategies

Modeling what makes an advertisement persuasive, *i.e.*, eliciting the desired response from consumer, is critical to the study of propaganda, social psychology, and marketing. Despite its importance, computational modeling of persuasion in computer vision is still in its infancy, primarily due to the lack of benchmark datasets that can provide persuasion-strategy labels associated with ads. Motivated by persuasion literature in social psychology and marketing, we introduce an extensive vocabulary of persuasion strategies and build the first ad corpus (both image and video) annotated with persuasion strategies. We then formulate the task of persuasion strategy prediction with multi-modal learning. The image dataset also provides image segmentation masks, which labels persuasion strategies in the corresponding ad images on the test split. We publicly release our code and dataset at <https://midas-research.github.io/persuasion-advertisements/>. This chapter is based on two papers I published along with collaborators [13, 15].

Marketing communications is the mode by which companies and governments inform, remind, and persuade their consumers about the products they sell. They are the primary means of connecting brands with consumers through which the consumer can know what the product is about, what it stands for, who makes it, and can be motivated to try it out. To introduce meaning into their communication, marketers use various rhetorical devices in the form of persuasion strategies such as **emotions** (*e.g.*, Oreo’s “Celebrate the Kid Inside”, humor by showing Ronald McDonald sneaking into the competitor Burger King’s store to buy a burger), **reasoning** (*e.g.*, “One glass of Florida orange juice contains 75% of your daily vitamin C needs”), **social identity** (*e.g.*, Old Spice’s “Smell like a Man”), and **impact** (*e.g.*, Airbnb showing a mother with her child with the headline “My home is funding her future”) (Refer to Fig. 2.3 to see these ads). Similarly, even for marketing the same product, marketers use different persuasion



Figure 2.1: Different persuasion strategies are used for marketing the same product (footwear in this example). The strategies are in red words and to be defined by us in the paper.



Figure 2.2: Examples of videos with their annotated persuasion strategies. Relevant keyframes and ASR captions are shown in the figure, along with the annotated strategies. These two videos can be watched at <https://bit.ly/3Ie3JG0>, <https://bit.ly/30gtLwj>.

strategies to target different demographies (see Fig. 2.1). Therefore, recognizing and understanding persuasion strategies in ad campaigns is vitally important to decipher viral marketing campaigns, propaganda, and enable ad-recommendation.

Studying rhetorics of this form of communication is an essential part of understanding visual communication in marketing. Aristotle, in his seminal work on rhetoric, underlining the importance of persuasion, equated studying rhetorics with the study of persuasion* [70]. While persuasion is studied extensively in social science fields, including marketing [71, 72] and psychology [73, 74], computational modeling of persuasion in computer vision is still in its infancy, primarily due to the lack of benchmark datasets that can

*“Rhetoric may be defined as the faculty of discovering in any particular case all of the available means of *persuasion*” [70]



Figure 2.3: Various rhetoric strategies used in advertisements

provide representative corpus to facilitate this line of research. In the limited work that has happened on persuasion in computer vision, researchers have tried to address the question of which image is more persuasive [75] or extracted low-level features (such as emotion, gestures, and facial displays), which indirectly help in identifying persuasion strategies without explicitly extracting the strategies themselves [76]. On the other hand, decoding persuasion in textual content has been extensively studied in natural language processing from both extractive, and generative contexts [12, 14, 77]. This forms the motivation of our work, where we aim to identify the persuasion strategies used in visual content such as advertisements.

The systematic study of persuasion began in the 1920s with the media-effects research by Lasswell [44], which was used as the basis for developing popular models of persuasion, like the Elaboration Likelihood Model (ELM) [74], Heuristic Systematic Model (HSM) [46], and Hovland’s attitude change approach [73]. Laswell in this research broke down communication into five factors by defining communication as an act of *who* said it, *what* was said, in *what* channel it was said, to *whom* it was said, and with what *effect* it was said. Later, this model was used as the basis for developing popular models of persuasion, like Elaboration Likelihood Model [74], Heuristic Systematic Model [46], and Hovland’s attitude change approach [73]. Amongst these, the most widely accepted model of persuasion theory is the Elaboration Likelihood Model (ELM).

These models of persuasion posit a dual process theory that explains attitude and behavior change (persuasion) in terms of the following major factors: stimuli (messages), personal motivation (the desire to process the message), capability of critical evaluation, and cognitive busyness. These factors could be divided into cognitive, behavioral, and affective processes of attitude change. Thus, a person may begin liking a new political candidate because she just donated \$100 to the campaign (behavior-initiated change), because the theme music in a recently heard commercial induced a general pleasantness (affect-initiated change), or because the person was impressed with the candidate’s issue positions (cognitive initiated change). Similarly, if a person already likes a political candidate he may agree to donate money to the campaign (behavioral influence), may feel happiness upon meeting the

candidate (affective influence), and may selectively encode the candidate’s issue positions (cognitive influence) [74].

ELM posits that when facing a message from a persuader, the persuadee reacts by using the two information processing channels: central processing or peripheral processing. When the persuadee processes information centrally, the cognitive responses, or elaborations, will be much more relevant to the information, whereas when processing peripherally, the individual may rely on heuristics and other rules of thumb when elaborating on a message. The factors which influence how and how much one will elaborate the persuasive message is given by the message type, personal motivation, and other factors presented in the ELM. Being at the high end of the elaboration continuum, people assess object-relevant information in relation to schemas that they already possess, and arrive at a reasoned attitude that is supported by information [78].

In this chapter, we build on these psychological insights from persuasion models in sociology and marketing and study the message strategies that lead to persuasion. We codify, extend, and unify persuasion strategies studied in the psychology and marketing literature into a set of 20 strategies divided into 9 groups (see Fig. 2.4, Table 2.2): *Authority and Credibility*, *Social Identity and Proof*, where cognitive indirection in the form of group decisioning and expert authority is used for decisions, *Value and Impact Formulation* where logic is used to explain details and comparisons are made, *Reciprocity*, *Foot in the door*, *Overcoming Resistance* where social and cognitive consistency norms are harnessed to aid decision-making, *Scarcity*, *Anthropomorphism* and *Emotion* where information is evaluated from the lenses of feelings and emotions. In addition to introducing the most extensive vocabulary for persuasion strategies, we make a superset of persuasion strategies presented in the prior NLP works, which introduced text and domain-specific persuasion tactics, thus making large-scale understanding of persuasion across multiple contexts comparable and replicable.

Constructing a large-scale dataset containing persuasion strategies labels is time-consuming and expensive. We leverage active learning to mitigate the cost of labeling fine-grained persuasion strategies in advertisements. We first introduce an attention-fusion model trained in a multi-task fashion over modalities such as text, image, and symbolism. We use the action-reason task from the Pitts Ads dataset [79] to train the model and then annotate the raw ad images from the same dataset for persuasion strategies based on an entropy based active learning technique.

To sum up, our contributions include:

1. We construct the largest set of generic persuasion strategies based on theoretical and empirical studies in marketing, social psychology, and machine learning literature.
2. We introduce the first dataset for studying persuasion strategies in advertisements. This enables initial progress on the challenging task of auto-



Figure 2.4: Persuasion strategies in advertisements. Marketers use both text and vision modalities to create ads containing different messaging strategies. Different persuasion strategies are constituted by using various rhetorical devices such as slogans, symbolism, colors, emotions, allusion.

matically understanding the messaging strategies conveyed through visual advertisements. We also construct a prototypical dataset containing image segmentation masks annotating persuasion strategies in different segments of an image.

3. We formulate the task of predicting persuasion strategies with a multi-task attention fusion model.
4. We conduct extensive experiments on the released corpus, showing the effect of different modalities on identifying persuasion strategies, correlation between strategies and topics and objects with different strategies.

2.1 Related Work

How do messages change people’s beliefs and actions? The systematic study of persuasion has captured researchers’ interest since the advent of mass influence mechanisms such as radio, television, and advertising. Work in persuasion spans across multiple fields, including psychology, marketing, and machine learning.

Persuasion in Marketing and Social Psychology: Sociology and communication science has studied persuasion for centuries now starting from the seminal work of Aristotle on rhetoric. Researchers have tried to construct and validate models of persuasion. Due to space constraints, while we cannot cover a complete list of literature, in Section 2.2, we list the primary studies which originally identified the presence and effect of various persuasion tactics on persuadees. We build on almost a century of this research and crystallize them into the persuasion strategies we use for anno-

tation and modeling. Any instance of (successful) persuasion is composed of two events: (a) an attempt by the persuader, which we term as the persuasion strategies, and (b) subsequent uptake and response by the persuadee [80, 81]. In this work, we study (a) only while leaving (b) for future work. Throughout the rest of the paper, when we say persuasion strategy, we mean the former without considering whether the persuasion was successful or not.

Persuasion in Machine Learning: Despite extensive work in social psychology and marketing on persuasion, most of the work is qualitative, where researchers have looked at a small set of messages with various persuasion strategies to determine their effect on participants. Computational modeling of persuasion is still largely lacking. In the limited work in computational modeling of persuasion, almost all of it is concentrated in the NLP literature, with only very few works in computer vision. Research on persuasion in NLP under the umbrella of argumentation mining is broadly carried out from three perspectives: extracting persuasion tactics, studying the effect of constituent factors on persuasion, and measurement of persuasiveness nature of content. A few examples of research studies that annotate persuasive strategies in various forms of persuader-persuadee interactions like discussion forums, social media, blogs, academic essays, and debates are [80, 82, 83]. We use these and other studies listed in Section 2.2 to construct our vocabulary of persuasion strategies in advertisements.

Other studies focus on factors such as argument ordering [84, 85], target audience [86], and prior beliefs [87] for their effect in bringing about persuasion. Studies such as [88, 89] also try to measure persuasiveness and generate persuasive content. The generation of persuasive (textual) messages has been studied [90] and, in particular, a novel ML method for learning user model tailored persuasion strategy has also been proposed [91, 92].

As one of the first works in the limited work in the computer vision domain, Joo *et al.* [76] introduced syntactical and intent features such as facial displays, gestures, emotion, and personality, which result in persuasive images. Their analysis was done on human images, particularly politicians, during their campaigns. Their work on political campaigners is more restrictive than general product and public-service advertisements. Moreover, they deal with low-level features such as gestures and personality traits depicted through the face, which are important for detecting persuasion strategies but are not persuasion strategies themselves. Recently, Bai *et al.* [75] studied persuasion in debate videos where they proposed two tasks: debate outcome prediction and intensity of persuasion prediction. Through these tasks, they predict the persuasiveness of a debate speech, which is orthogonal to the task of predicting the strategy used by the debater. Other similar works which discuss persuasiveness of images and videos are [93, 94].

2.2 Generic Taxonomy of Persuasion Strategies

This section introduces the generic taxonomy of persuasive strategies, their definitions, examples, and connections with prior work. Representative literature from a) SPM: Social Psychology and Marketing, b) ML: Machine Learning

1. Authority and Credibility: SPM:[5, 74, 95–99] ML:[77, 80, 100–102]
 - (a) **Guarantees:** Guarantees reduce risk and people try out such products more often.
 - (b) **Authority:** Authority indicated through expertise, source of power, third-party approval, credentials, and awards
 - (c) **Trustworthiness:** Trustworthiness indicated honesty and integrity of the source through tropes like years of experience, “trusted brand”, numbers and statistics
2. Social Identity and Proof: SPM:[103–107] ML: [80, 88, 100, 102, 108–112]
 - (a) **Social Identity:** *Normative* influence, which involves conformity with the positive expectations of “another”, who could be “another person, a group, or one’s self” (includes self-persuasion, fleeting attraction, alter-casting, and exclusivity)
 - (b) **Social Proof:** *Informational influence* by accepting information obtained from others as evidence about reality, *e.g.*, customer reviews and ratings
3. Reciprocity: SPM:[96, 113–116] ML:[77, 80, 84, 88, 100]
 - (a) **Reciprocity:** By *obligating* the recipient of an act to repayment in the future, the rule for reciprocity begets a sense of future obligation, often unequal in nature
4. Foot in the door: SPM: [96, 117, 118] ML:[83, 119, 120]
 - (a) **Foot in the door:** Starting with small requests followed by larger requests to facilitate compliance while maintaining *cognitive coherence*.
5. Overcoming Resistance: SPM:[121–123] ML:{None}
 - (a) **Overcoming Resistance:** Overcoming resistance (reactance) by postponing consequences to the future, by focusing resistance on realistic concerns, by forewarning that a message will be coming, by acknowledging resistance, by raising self-esteem and a sense of efficacy.
6. Value and Impact Formulation: SPM:[124–129] ML:[130, 131]
 - (a) **Concreteness:** Using concrete facts, evidence, and statistics to appeal to the logic of consumers
 - (b) **Anchoring and Comparison:** A product’s value is strongly influenced by what it is compared to.
 - (c) **Social Impact:** Emphasizes the importance or bigger (societal) impact of a product

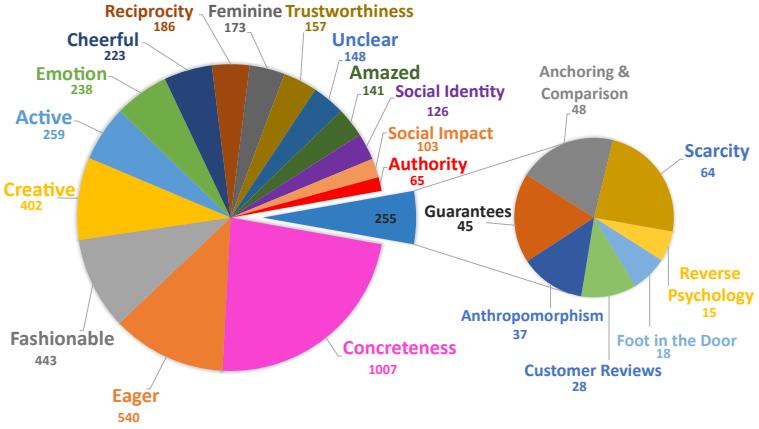


Figure 2.5: Distribution of Persuasion Strategies in the image persuasion strategy dataset. The top-3 strategies are Concreteness, Eager, and Fashionable.

7. Scarcity: SPM: [132–135] ML:[77, 84, 109]
 - (a) **Scarcity**: People assign more value to opportunities when they are less available. This happens due to psychological reactance of losing freedom of choice when things are less available or they use availability as a cognitive shortcut for gauging quality.
8. Anthropomorphism: SPM:[107, 136, 137] ML:{None}
 - (a) **Anthropomorphism**: When a brand or product is seen as human-like, people will like it more and feel closer to it.
9. Emotion: Aesthetics, feeling and other non-cognitively demanding features used for persuading consumers SPM:[74, 138, 139]
ML:[82, 101, 102, 109, 112, 130, 140]
 - (a) **Amazed**
 - (b) **Fashionable**
 - (c) **Active, Eager**
 - (d) **Feminine**
 - (e) **Creative**
 - (f) **Cheerful**
 - (g) **Further Minor**
10. **Unclear**: If the ad strategy is unclear

2.3 Persuasion Strategy Corpus Creation

2.3.1 Persuasion Strategy Dataset For Image Advertisements

To annotate persuasion strategies on image advertisements, we leverage raw images from the Pitts Ads dataset. It contains 64,832 image ads with labels of topics, sentiments, symbolic references (*e.g.* dove symbolizing peace), and

reasoning the ad provides to its viewers (see Fig 2.10 for a few examples). The dataset had ads spanning multiple industries, products, services, and also contained public service announcements. Through this, they presented an initial work for the task of understanding visual rhetoric in ads. Since the dataset already had a few types of labels associated with the ad images, we used active learning on a model trained in a multi-task learning fashion over the reasoning task introduced in their paper. We explain the model and then the annotation strategy followed in §2.4.

To commence training, we initially annotated a batch of 250 randomly selected ad images with persuasion strategies defined in Section 2.2. We recruited four research assistants to label persuasion strategies for each advertisement. Definitions and examples of different persuasion strategies were provided, together with a training session where we asked annotators to annotate a number of example images and walked them through any disagreed annotations. To assess the reliability of the annotated labels, we then asked them to annotate the same 500 images and computed Cohen’s Kappa statistic to measure inter-rater reliability. We obtained an average score of 0.55. The theoretical maximum of Kappa, given the unequal distribution, is 0.76. In such cases, Cohen [141] suggested that one should divide kappa by its maximum value k/k_{\max} , which comes out to be 0.72. This is a *substantial* agreement. Further, to maintain labeling consistency, each image was double annotated, with all discrepancies resolved by an intervention of the third annotator using a majority vote.

The assistants were asked to label each image with no more than 3 strategies. If an image had more than 3 strategies, they were asked to list the top 3 strategies according to the area covered by the pixels depicting that strategy. In total, we label 3000 ad-images with their persuasion strategies; and the number of samples in train, val, and test split are 2500, 250, and 250, respectively[†]. Fig. 2.5 presents the distribution of persuasion strategies in the dataset. It is observed that concreteness is the most used strategy in the dataset, followed by eagerness and fashion. The average number of strategies in an ad is 1.49, and the standard deviation is 0.592. We find that scarcity (92.2%), guarantees (91.1%), reciprocity (84.4%), social identity (83.3%), and cheerful (83%), are the top 5 strategies, which occur in groups of 2 or 3. We observe that the co-occurrence of these strategies is due to the fact that many of them cover only a single modality (*i.e.*, text or visual), leaving the other modality free for a different strategy. For example, concreteness is often indicated by illustrating points in text, while the visual modality is free for depicting, say, emotion. See Fig. 2.6 for an example, where the image depicting *Authority* also has concreteness strategy in it. Similarly, feminine emotion is also depicted in Fig. 2.1, along with concreteness.

[†]Table 2.4 shows the detailed distribution of the number of strategies in ads

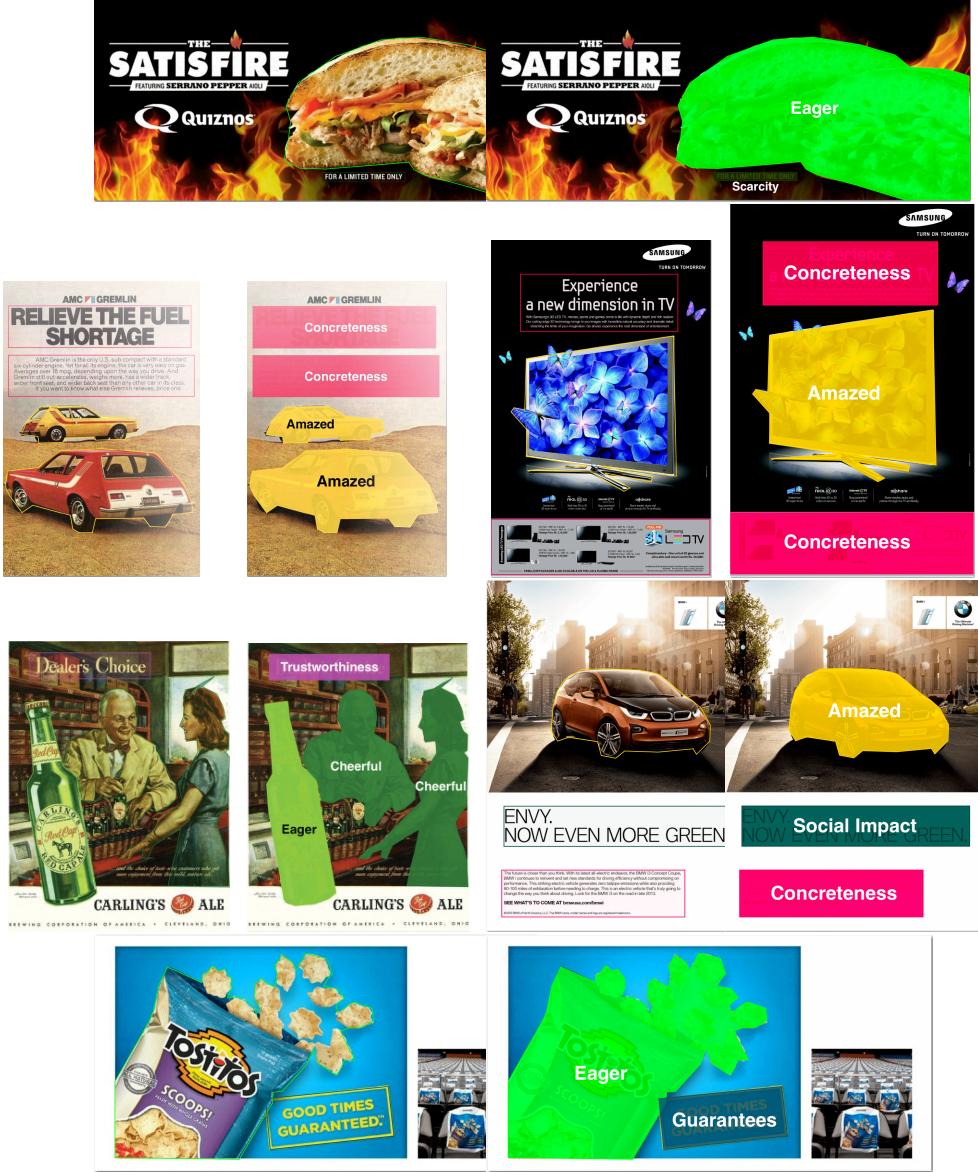


Figure 2.6: Image with a segmentation mask depicting the strategies *Emotion:Cheerful*, *Emotion:Eager* and *Trustworthiness*.

Next, we calculate the Dice correlation coefficient[‡] statistics for pairs of co-occurring persuasion strategies. The top-5 pairs are eager-concreteness (0.27), scarcity-reciprocity (0.25), eager-cheerful (0.19), amazed-concreteness (0.17), and eager-reciprocity (0.17). We find that these correlation values are

[‡]The Dice Coefficient is defined as: $2 * |X \cap Y| / (|X| + |Y|)$, where X and Y are two sets; a set with vertical bars on either side refers to the cardinality of the set, i.e. the number of elements in that set; and \cap refers to the intersection of two sets.



Figure 2.7: Advertisements containing humans and concreteness

not particularly high since marketers seldom use *common pairings* of messaging strategies to market their products. The visual part mostly shows eager strategy in ads; therefore, we find that the text modality becomes free to show other strategies. That is why primarily text-based concreteness, cheerfulness, and reciprocity strategies are present with the visual-based eager strategy in the text modality. Also, primarily vision-based amazement, eagerness, and scarcity (short-text) strategies co-occur with text-based reciprocity and concreteness (*e.g.*, Fig. 2.1).

Next, we calculate the correlation between image topics and objects present with persuasion strategies. We see that the emotion:feminine and emotion:fashionable strategies are most often associated with beauty products and cosmetics ($\text{corr}=0.4256, 0.2891$). This is understandable since most beauty products are aimed at women. We see that the fast-food and restaurant industries often use eagerness as their messaging strategy ($\text{corr} = 0.5877, 0.3470$). We find that the presence of humans in ads is correlated with the concreteness strategy (see Fig 2.7 for a few examples) ($\text{corr}=0.3831$). On the other hand, vehicle ads use emotion:amazed and concreteness ($\text{corr}=0.5211, 0.2412$) (see Fig:2.8 for detailed correlations).

Similar to a low correlation in co-occurring strategies, we find that product segments and their strategies are not highly correlated. This is because marketers use different strategies to market their products even within a product segment. Fig. 2.1 shows an example in which the footwear industry (which is a subsegment of the apparel industry) uses different strategies to market its products. Further, for a batch of 250 images, we also label segmented image regions corresponding to the strategies present in the image. These segment masks were also double-annotated. Fig. 2.6 presents an example of masks depicting parts of the image masked with different persuasion strategies in a drink advertisement.

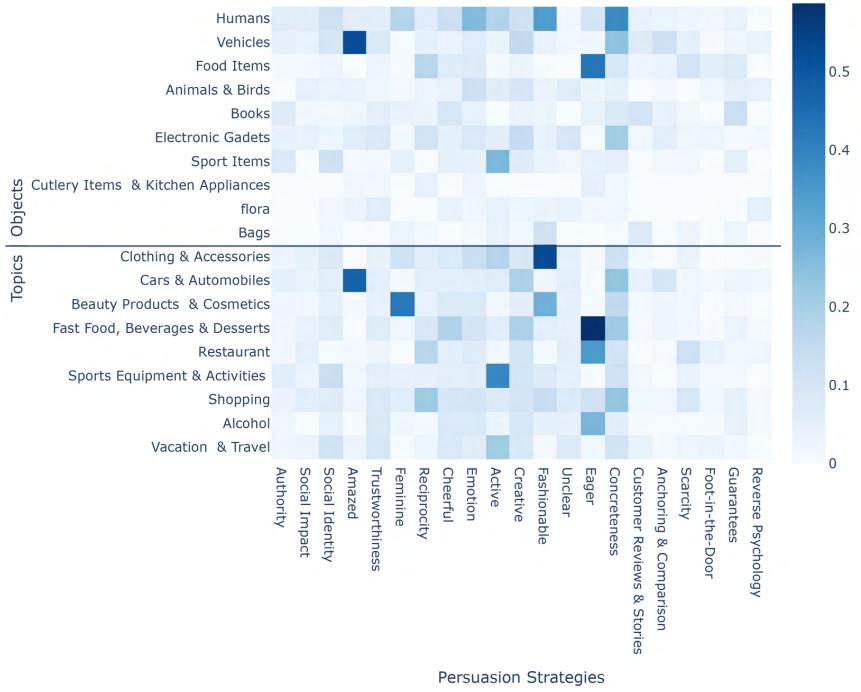


Figure 2.8: Dice correlation between topics and strategies. Topics are taken from the Pitts Ad dataset and further similar topics are combined to get these values.

2.3.2 Persuasion Strategy Dataset For Video Advertisements

For this task, we collected 2203 video advertisements from popular brands publicly available on the web. We use the following 12 strategies as our target persuasion strategy set: *Social Identity*, *Concreteness*, *Anchoring and Comparison*, *Overcoming Reactance*, *Reciprocity*, *Foot-in-the-Door*, *Authority*, *Social Impact*, *Anthropomorphism*, *Scarcity*, *Social Proof*, and *Unclear*. We use non-experts human annotators to label this dataset (as compared to expert humans for the image ads dataset). In order to make the class labels easier to understand for non-expert human annotators, we make a list of 15 yes/no type-questions containing questions like “*Was there any expert (person or company) (not celebrity) encouraging to use the product/brand?* *Was the company showcasing any awards (e.g., industrial or government)?* *Did the video show any customer reviews or testimonials?*” (complete list in Table 2.1).

Each human annotator watches 15 videos such that each video gets viewed by at least two annotators and answers these questions for each video. Based on all the responses for a video, we assign labels to that video. We remove videos with an inter-annotator score of less than 60%. After removing those, we get a dataset with 1002 videos, with an average length

of 33 secs and a distribution as shown in Fig. 2.9. This dataset is then used for the persuasion strategy identification task.

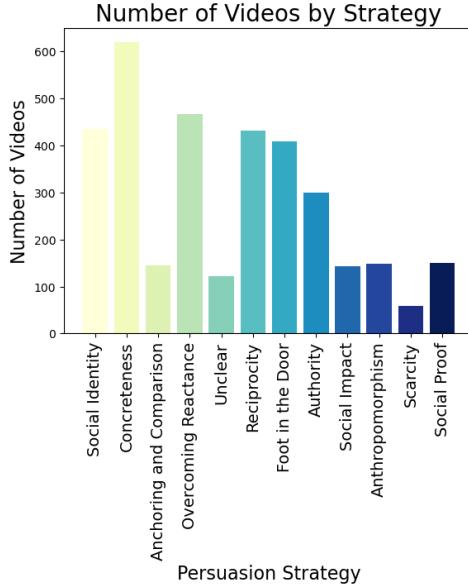


Figure 2.9: Distribution of persuasion strategies in our video persuasion strategy dataset

2.4 Modeling: Persuasion Strategy Prediction

2.4.1 Modelling Persuasion Strategy For Image Advertisements

The proposed Ads dataset \mathcal{D} annotated with the persuasion strategies comprises of samples where each sample advertisement a_i is annotated with a set of annotation strategies S_i such that $1 \leq |S_i| \leq 3$. The unique set of the proposed persuasion strategies \mathcal{P} is defined in Table 2.2. Given a_i , the task of the modeling is to predict the persuasion strategies present in the input ad. As we observe from Fig. 2.4, advertisements use various rhetoric devices to form their messaging strategy. The strategies thus are in the form of multi-modalities, including images, text and symbolism. To jointly model the modalities, we design an attention fusion multi-modal framework, which fuses multimodal features extracted from the ad, *e.g.*, the ad image, text present in the ad extracted through the OCR (Optical Character Recognition), regions of interest (ROIs) extracted using an object detector, and embeddings of captions obtained through an image captioning model (see Fig. 2.11). The information obtained through these modalities are firstly embedded independently through their modality specific encoders followed

| Question | Strategy | Question | Strategy |
|--|--|--|-------------------------------|
| Was there any expert (person or company) (not celebrity) encouraging to use the product/brand? | Authority | Did the video show any normal customers (non-expert, non-celebrity) using the product? | Social Identity |
| Did the video showcase any awards or long usage history of the product/brand? | Authority | Did the video show any customer reviews or testimonials? | Social Proof |
| Was the product/brand comparing itself with other competitors or existing solutions? | Anchoring and Comparison | Were any number/statistics mentioned? | Concreteness |
| Did the video talk about any specific features or provide information about the product/brand? | Concreteness | Were there any mention of any offers on the brand/product? | Reciprocity |
| Were the offers limited or available for a short period of time? | Scarcity | Was the product/brand told to be free or available on a discount? | Foot in the Door, Reciprocity |
| Was the brand/product described as simple, easy to use, or can start using with minimal resistance? | Overcoming Reactance, Foot in the Door | Was the brand/product talking about bigger societal impact? | Social Impact |
| Did the brand provide any guarantees that might help reduce the risk of people trying out the product? | Overcoming Reactance | Did the video provide any resources, tips, guides, or tools related to the product? | Reciprocity |
| Is the brand or product portrayed as human-like? | Anthropomorphism | | |

Table 2.1: The questions we asked to the non-expert annotators to help them identify persuasion strategy contained in the video advertisement.

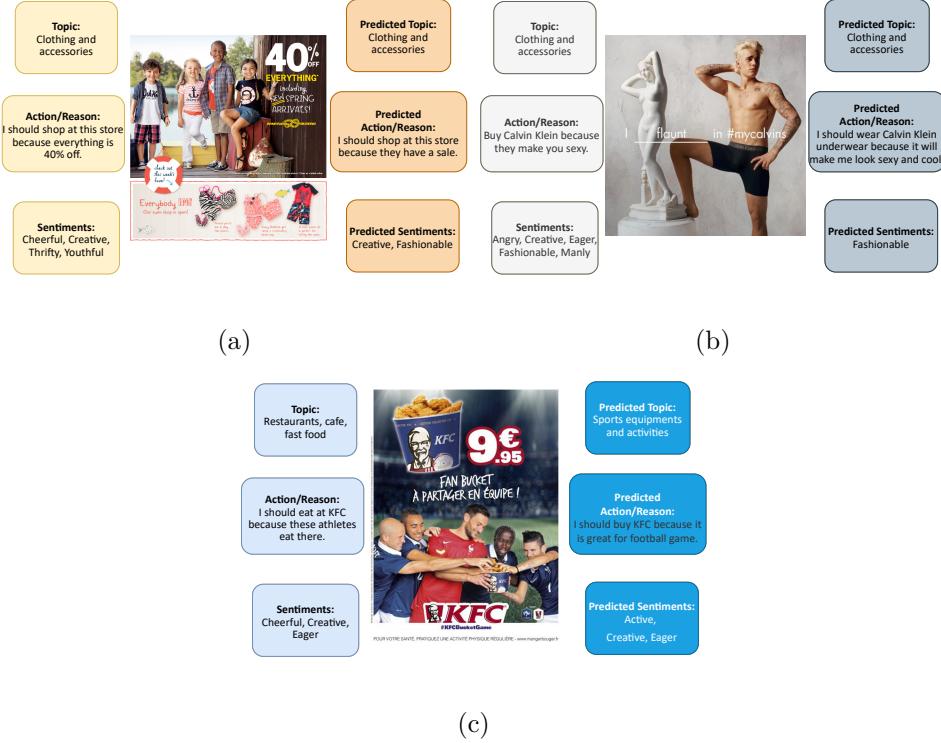


Figure 2.10: Some samples from the Pitts Ads dataset along with the ground truth and predicted action-reason statement, topic and sentiments.

by a transformer-based cross-attention module to fuse the extracted features from different modalities. The fused embeddings from the attention module are then used as input for a classifier that predicts a probability score for each strategy $p \in \mathcal{P}$. The overall architecture of the proposed model is illustrated in Fig.2.11. In the following, we describe each step in the prediction pipeline in detail.

2.4.1.1 Feature Extractors

In order to capture different rhetoric devices, we extract features from the image, text, and symbolism modalities.

Image Feature: We use the Vision Transformer [142] (ViT) model for extracting image features from the entire input image. The model resizes the input image to size 224×224 and divides it into patches of size 16×16 . The model used has been pre-trained on the ImageNet 21k dataset. We only use the first output embedding, which is the CLS token embedding, a 768 dimension tensor, as we only need a representation of the entire image. Then, a fully connected layer is used to reduce the size of the embedding, resulting in a tensor of dimension 256.

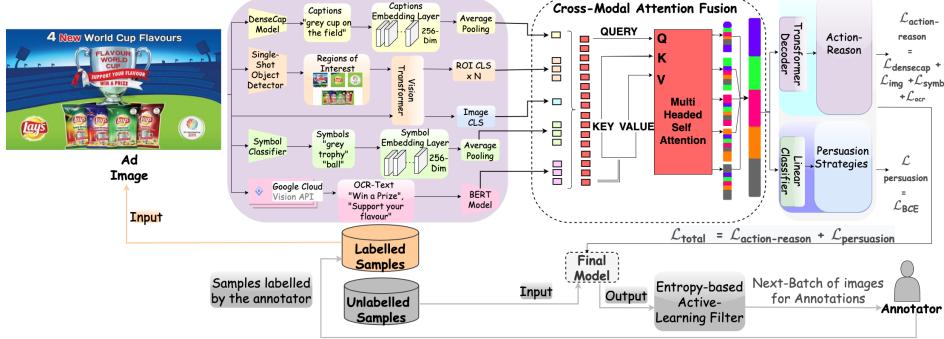


Figure 2.11: Architecture of the Persuasion Strategy Prediction model. To capture the different rhetoric devices, we extract features for the image, text, and symbolism modalities and then apply cross-modal attention fusion to leverage the interdependence of the different devices. Further, the model trains over two tasks: persuasion strategies and the reasoning task of action-reason prediction.

Regions of Interest (RoIs) from Detected Objects and Captions:

Ad images contain elements that the creator deliberately chooses to create *intentional impact* and deliver some *message* in addition to the ones that occur *naturally* in the environment. Therefore, it is important to identify the composing elements of an advertisement to understand the creator's intention and the ad's message to the viewer. We detect and extract objects as regions of interest (RoIs) from the advertisement images. We get the RoIs by training the single-shot object detector model [143] on the COCO dataset [144]. We compare it with the recent YOLOv5 model [145]. We also extract caption embeddings to detect the most important activity from the image using a caption generation mode. We compare DenseCap [146] and the more recent BLIP [147] for caption generation.

OCR Text: The text present in an ad presents valuable information about the brand, such as product details, statistics, reasons to buy the product, and creative information in the form of slogans and jingles that the company wants its customers to remember and thus making it helpful in decoding various persuasion strategies. Therefore, we extract the text from the ads and use it as a feature in our model. We use the Google Cloud Vision API for this purpose. All the extracted text is concatenated, and the size is restricted to 100 words. We pass the text through a BERT model and concatenate the embeddings for those 100 words. Similar to image embeddings, an FC layer is used to convert embeddings to 256 dimensions. The final embedding of the OCR is a tensor of dimension 100×256 .

Symbolism: While the names of the detected objects convey the names or literal meaning of the objects, creative images often also use objects for their symbolic and figurative meanings. For example, an upward-going arrow represents growth or the north direction or movement towards the up-

ward direction depending on the context; similarly, a person with both hands pointing upward could mean danger (*e.g.*, when a gun is pointed) or joy (*e.g.*, during dancing). In Fig. 2.4, in the creative Microsoft ad, a symbol of a balloon is created by grouping multiple mice together. Therefore, we generate symbol embeddings to capture the symbolism behind the most prominent visual objects present in an ad. We use the symbol classifier by Hussain *et al.* [79] on ad images to find the distribution of the symbolic elements present and then convert this to a 256 dimension tensor.

2.4.1.2 Cross-Modal Attention

To capture the inter-dependency of multiple modalities for richer embeddings, we apply a cross-modal attention (CMA) layer [148] to the features extracted in the previous steps. Cross-modal attention is a fusion mechanism where the attention masks from one modality (*e.g.* text) are used to highlight the extracted features in another modality (*e.g.* symbolism). It helps to link and extract common features in two or more modalities since common elements exist across multiple modalities, which complete and reinforce the message conveyed in the ad. For example, the pictures of the silver cup, stadium, and ball, words like “Australian”, “Pakistani”, and “World Cup” present in the chips ad shown in Fig. 2.11 link the idea of buying *Lays* with supporting one’s country’s team in the World Cup. Cross attention can also generate effective representations in the case of missing or noisy data or annotations in one or more modalities [148]. This is helpful in our case since marketing data often uses implicit associations and relations to convey meaning.

The input to the cross-modal attention layer is constructed by concatenating the image, RoI, OCR, caption, and symbol embeddings. This results in a 114×256 dimension input to our attention layer. The cross-modal attention consists of two layers of transformer encoders with a hidden dimension size of 256. The output of the attention layer gives us the final combined embedding of our input ad. Given image embeddings E_i , RoI embeddings E_r , OCR embeddings E_o , caption embeddings E_c and symbol embeddings E_s , the output of the cross-attention layer E_{att} is formulated as:

$$\text{Enc}(X) = \text{CMA}([E_i(X), E_r(X), E_o(X), E_c(X), E_s(X)]),$$

where $[.\dots,.]$ is the concatenation operation. For the advertisement in Fig. 2.11, we observed that the caption “grey cup on the field” attends to OCR text (containing words like “win”) and ViT features of the RoI (of “cup” and “field”).

2.4.1.3 Persuasion Strategy Predictor

This module is a persuasion strategy predictor, which processes the set of feature embedding $\text{Enc}(X)$ obtained through cross-modality fusion. Specifically, $\text{Enc}(X)$ is passed through a self-attention layer as:

$$o_1 = \text{softmax}(\text{Enc}(X) \otimes W_{\text{self-attn}})^T \otimes \text{Enc}(X) \quad (2.1)$$

where $\text{Enc}(X)$ is of the dimension 114×256 , $W_{\text{self-attn}} \in \mathcal{R}^{256 \times 1}$, \otimes denote tensor multiplication and o_1 denotes the output of self attention layer, which is further processed through a linear layer to obtain $o_{|\mathcal{P}|}$ to represent the logits for each persuasion strategy. We apply sigmoid over each output logit such that the i^{th} index of the vector after applying sigmoid denotes p_i - the probability with which i^{th} persuasion strategy is present in the ad image. Our choice of using sigmoid over softmax is motivated by the fact that multiple persuasion strategies can be present simultaneously in an ad image. Consequently, the entire model is trained in an end-to-end manner using binary cross-entropy loss \mathcal{L}_s over logit for each strategy:

$$\mathcal{L}_s = [-y_i \log(p_i) - (1 - y_i) \log(1 - p_i)] \quad (2.2)$$

where, y_i is 1 if i^{th} persuasion strategy is present in the ad and 0 otherwise. It can be observed in Table 2.2 that our model achieves an accuracy of 59.2%, where a correct match is considered if the strategy predicted by the model is present in the set of annotated strategies for a given ad. Further, we perform several ablations where we exclude each modality while retaining all the other modalities. We note that for each modality, excluding the modality results in a noticeable decrease in accuracy, with significant decreases observed when excluding DenseCap ($\sim 3.6\%$) and OCR ($\sim 4.4\%$). Further, we observe that using DenseCap for obtaining caption embeddings, and SSD for object detection works better than BLIP and YOLOv5, respectively (see Table 2.3). We also explore using focal loss [149] in place of cross-entropy loss to handle class imbalance but observed that it led to degradation instead of improvements (top-1 acc.[§] of 56.4% vs 59.2% using cross-entropy). We also train the model of Hussain *et al.* [79] for strategy prediction through a similar configuration as ours (along with action-reason generation using an LSTM branch). We find that their top-1 and top-3 accuracy is 52.4% (vs. 59.2% ours) and 75.7% (vs. 84.8% ours), which is lesser compared to our model.

[§] *Top-1 Accuracy*: It is defined as the fraction of images, where the highest predicted strategy is present in the ground-truth strategies. *Top-3 Accuracy* : It is defined as the fraction of images, where any of the top-3 highest predicted strategies is present in the ground-truth strategies.

| Models | Top-1 Acc. | Top-3 Acc. |
|-------------------------------|-------------|-------------|
| Our Model | 59.2 | 84.8 |
| w/o DenseCap | 55.6 | 80.8 |
| w/o Symbol | 58.8 | 81.6 |
| w/o DenseCap & Symbol | 55.2 | 80.8 |
| w/o OCR | 54.8 | 82 |
| w/o Symbol, OCR & DenseCap | 58 | 78.8 |
| w/o Action-Reason Task | 56.4 | 80.4 |
| Random Guess | 6.25 | 18.75 |

Table 2.2: Effect of different Modalities and Tasks on the accuracy and performance of the strategy prediction task.

2.4.1.4 Multi Task Learning

One of the key opportunities for our persuasion strategies data labeling and modeling task was the presence of additional labels already given in the base Pitts Ads dataset. In that, authors had given labels about the reasoning task. For the reasoning task, the annotators were asked to provide answers in the form “I should [Action] because [Reason].” for each ad. In other words, they asked the annotators to describe *what the viewer should do and why*, according to the ad. Similar to the reasoning task, persuasion strategies provide various cognitive, behavioral, and affective reasons to try to elicit the motivation of the ad viewers towards their products or services. Therefore, we hypothesize that these natural language descriptions of *why the viewers should follow* the ad will be informative in inferring the ad’s persuasion strategy.

We formulate obtaining action-reason statement as a sequence generation task where the model learns to generate a sentence $Y^g = (y_1^g, \dots, y_T^g)$ of length T conditioned on advertisement X by generating the sequence of tokens present in the action-reason statement. To achieve this, we use a transformer decoder module that attends on the features $\text{Enc}(X)$ as shown in Fig. 2.11. The annotated action-reason statement is used to train the transformer decoder as an auxiliary task to strategy prediction through the standard teacher forcing technique used in Seq2Seq framework. Please refer to the Supplementary for more architectural details about the action-reason generation branch. As shown in Table 2.2, generating action-reason as an auxiliary task improves the strategy prediction accuracy by 2.8%. We

| Model Used | Top-1 Accuracy | Top-3 Accuracy | Recall |
|---------------------------|----------------|----------------|--------|
| Model with DenseCap & SSD | 59.2 | 84.8 | 74.59 |
| Model with BLIP & YOLOv5 | 58.4 | 83.8 | 71.58 |

Table 2.3: Comparison of caption and object detection models. We noticed that BLIP while being more recent and trained on a larger dataset, generates more informative captions for background objects which DenseCap successfully ignores.

evaluate the performance on action-reason generation on following metrics: BLEU-1, BLEU-2, BLEU-3, BLEU-4, METEOR, ROUGE, CIDEER, SPICE and observed a score of 53.6, 42.0, 33.1, 25.7, 26.3, 48.4, 42.8, 8.9 respectively.

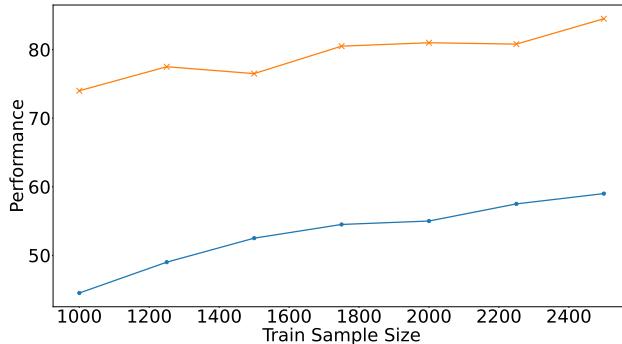


Figure 2.12: Incremental effect of introducing new data through active learning; Results for prediction of persuasion strategies on the test set

2.4.1.5 Active Learning

We use an active learning method to ease the large-scale label dependence when constructing the dataset. As in every active learning setting, our goal is to develop a learner that selects samples from unlabeled sets to be annotated by an oracle. Similar to traditional active learners [150, 151], we use uncertainty sampling to perform the sample selection. In doing so, such function learns to score the unlabeled samples based on the expected performance gain they are likely to produce and used to update the current version of the localization model being trained. To evaluate each learner, we measure the performance improvements, assessed on a labeled test set at different training dataset sizes.

| | #ads with 1 strategy | #ads with 2 strategies | #ads with 3 strategies | Avg. strategies | Std. Dev. |
|-----------|-------------------------|---------------------------|---------------------------|--------------------|--------------|
| Train-Set | 1440 | 905 | 155 | 1.486 | 0.612 |
| Val-Set | 132 | 98 | 20 | 1.552 | 0.639 |
| Test-Set | 147 | 93 | 10 | 1.452 | 0.574 |
| Total | 1719 | 1096 | 185 | 1.49 | 0.592 |

Table 2.4: Distribution of test, train, validation, and the total dataset

At every learning step t , a set of labeled samples L_t is first used to train a model f_t . Then, from an unlabeled pool $U_t = D - L_t$, an image instance a is chosen by a selection function g . Afterwards, an oracle provides temporal ground-truth for the selected instance, and the labeled set L_t is augmented with this new annotation. This process repeats until the desired performance is reached or the set U_t is empty.

In our implementation, we instantiate the active learning selection function as the entropy of the probability distribution predicted by the model over the set of persuasion strategies for a given ad image instance a . Formally, $g = -\sum_{i=1}^{|P|} p_i^n * \log(p_i^n)$, where p_i^n denotes the normalized probability with which i^{th} persuasion strategy is present in a as per the model prediction. The normalized probability p_i^n is estimated as $p_i / \sum_{j=1}^{|P|} p_j$. Intuitively, ad samples with high entropy selection values indicate that the model trained on limited data has a higher degree of confusion while predicting the persuasion strategy since it is not decisively confident about predicting few strategies. Hence, we rank the unlabeled ad images in the decreasing order of difficulty according to the corresponding values of the entropy selection function and select the top-k ads in the subsequent batch for annotation followed by training. As shown in Fig. 2.12, we set k to be 250 and analyze the effect of incrementally introducing new samples selected through active learning. It can be seen that both top-1 and top-3 accuracy increases with the addition of new training data. We stop at the point when 2500 training samples are used since the model performs reasonably well with a top-1 and top-3 strategy prediction accuracy of 59.2% and 84.8% (see Fig. 2.12).

2.4.2 Modelling Persuasion Strategy For Video Advertisements

Large Language Models (LLMs) have been demonstrated to perform well for downstream classification tasks in the text domain. This powerful ability has been widely verified on natural language tasks, including text classification, semantic parsing, mathematical reasoning, *etc.* Inspired by these

advances of LLMs, we aim to explore whether they could tackle reasoning tasks on multimodal data (*i.e.* videos). Therefore, we propose a storytelling framework, which leverages the power of LLMs to verbalize videos in terms of a text-based story and then performs downstream video understanding tasks on the generated story instead of the original video. Our pipeline can be used to verbalize videos and understand videos to perform complex downstream tasks such as emotion, topic, and persuasion strategy detection.

We show the performance of our framework on fifteen distinct tasks across five datasets. Firstly, we employ a video story dataset to evaluate the story generation task. Secondly, we utilize a video advertisements dataset to assess topic and emotion classification, as well as action and reason generation. Then, the persuasion strategy dataset to evaluate the task of understanding persuasion strategies within stories, and finally, HVU and LVU for concept, user engagement, and attribute prediction. These diverse datasets allow us to evaluate the performance and capabilities of our framework thoroughly.

1. The Video story dataset [152] contains 105 videos, from four types of common and complex events (*i.e.* birthday, camping, Christmas, and wedding) and corresponding stories written by annotators. It has longer videos (average length 12.4 mins) and longer descriptions (162.6 words on average). Moreover, the sentences in the dataset are more sparsely distributed across the video (55.77 sec per sentence). *Metrics:* Following [152], we use several NLP metrics, *viz.*, BLEU-N, ROUGE-L, METEOR and CIDEr to measure the similarity between the story generated by the model and ground truth.

2. The Image and Video Advertisements [79] contains 3,477 video advertisements and the corresponding annotations for emotion and topic tags and action-reason statements for each video. There are a total of 38 topics and 30 unique emotion tags per video. Further, we have 5 action-reason statements for each video for the action-reason generation task. For our experiment, we use 1785 videos, due to other videos being unavailable/privated from Youtube.

Metrics: Following [79], for the topic and emotion classification task, we evaluate our pipeline using top-1 accuracy as the evaluation metric. Further, since [79] did not use any fixed set of vocabulary for annotations, rather they relied on annotator-provided labels, the labels are often very close (like cheerful, excited, and happy). Therefore, based on nearness in Plutchik’s [153] wheel of emotions, we club nearby emotions and use these seven main categories: joy, trust, fear, anger, disgust, anticipation, and unclear. For the action-reason task, following [79], we evaluate our accuracy on the action and reason retrieval tasks where 29 random options along with 1 ground truth are provided to the model to find which one is the ground truth. Further, we also generate action and reason statements and evaluate the generation’s faithfulness with the ground truth using metrics like ROUGE,

BLEU, CIDEr, and METEOR.

3. Persuasion strategy dataset: This is the dataset we contribute for understanding persuasion strategies.

Metrics: We evaluate the performance using top-1 accuracy metric. Videos have a varied number of strategies, therefore, we consider a response to be correct if the predicted strategy is present among the list of ground-truth strategies.

4. Long-Form Video Understanding (LVU): We *et al.* [154] released a benchmark comprising of 9 diverse tasks for long video understanding and consisting of over 1000 hours of video. The various tasks consist of content understanding ('relationship', 'speaking style', 'scene/place'), user engagement prediction ('YouTube like ratio', 'YouTube popularity'), and movie metadata prediction ('director', 'genre', 'writer', 'movie release year'). We *et al.* [154] use top-1 classification accuracy for content understanding and metadata prediction tasks and MSE for user engagement prediction tasks.

5. Holistic Video Understanding (HVU): HVU [155] is the largest long video understanding dataset consisting of 476k, 31k, and 65k samples in train, val, and test sets, respectively. A comprehensive spectrum includes the identification of various semantic elements within videos, consisting of classifications of scenes, objects, actions, events, attributes, and concepts. To measure performance on HVU tasks, similar to the original paper, we use the mean average precision (mAP) metric on the validation set.

Next, we explain our pipeline to solve these tasks.

2.4.2.1 Video Verbalization

To obtain a verbal representation of a video, we employ a series of modules that extract unimodal information from the multimodal video. This information is then used to prompt a generative language model (such as GPT-3.5 [59] and Flan-t5 [156]) to generate a coherent narrative from the video. The overall pipeline is depicted in Fig. 2.13. In the following, we delve into each component of the framework in details.

1. Video Metadata: Understanding the context of a story is crucial, and we achieve this by gathering information about the communicator (brand). We leverage the publicly available video title and channel name from the web. Additionally, we utilize Wikidata [157], a collaborative knowledge base that provides comprehensive data for Wikipedia, to obtain further details such as the company name, product line, and description. This information helps us comprehend the story elements and establish connections with the brand's business context. For non-advertisement videos, we skip this step and retrieve only the video title.

2. Text Representation of Video Frames: We extract two types of textual information from video frames. Firstly, we capture the literal text

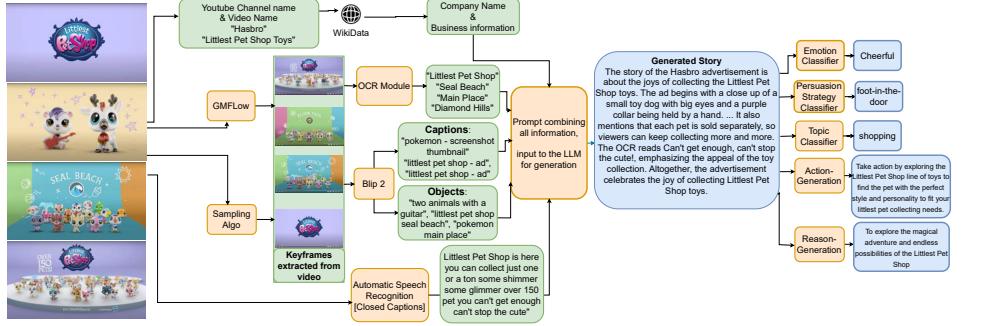


Figure 2.13: The overview of our framework to generate a story from a video and perform downstream video-understanding tasks. First, we sample keyframes from the video which are verbalized using BLIP-2. We also extract OCR from all the frames. Next, using the channel name and ID, we query Wikidata to get company and product information. Next, we obtain automatically generated captions from YouTube videos using the YouTube API. All of these are concatenated as a single prompt and given as input to an LLM and ask it to generate the story of the advertisement. Using the generated story, we then perform the downstream tasks of emotion and topic classification and persuasion strategy identification. This video can be watched at <https://youtu.be/ZBLkTALi1CI>.

present on the frames. Secondly, we analyze the scene depicted in each frame to gain a deeper understanding. In the upcoming sections, we will elaborate on both of these aspects.

| Method | Frame Extraction | METEOR | CIDEr | Rouge-L | BLEU-1 | BLEU-2 | BLEU-3 | BLEU-4 |
|---------|------------------|--------|-------|---------|--------|--------|--------|--------|
| GPT-3.5 | Uniform Sampling | 24.8 | 102.4 | 24.3 | 63.8 | 56.4 | 47.2 | 38.6 |
| GPT-3.5 | Pyscenedetect | 24.17 | 67.8 | 21.17 | 54.59 | 49.05 | 41.54 | 33.88 |

Table 2.5: Comparison of Pyscenedetect [158] with uniform sampling of choosing video frames. Based on downstream performance, we can see that uniform sampling works better than Pyscenedetect

a. *Visual and Scenic Elements in Frames:* For videos with a duration shorter than 120 seconds, we employ an optical flow-based heuristic using the GMFlow model [159] to extract keyframes. In shorter advertisement videos, scene changes often indicate transitions in the story, resulting in keyframes with higher optical flow values. The GMFlow model effectively captures these story transitions. We select frames with an optical flow greater than 50 and prioritize frames with maximum pixel velocity. However, for longer videos, this approach yields a large number of frames that are difficult to accommodate within a limited context. To address this, we sample frames at a uniform rate based on the native frames-per-second (fps) of the video (see Table 2.5 for a comparison between uniform sampling and Pyscenedetect).



Transcript : well hey Zuma that's the wrong pinata how old are you my birthday cake is its zoo paw patrol you got as a present ... Thank You Marshall... are you ready to blow out the candle.... Marshall let's push chase and the pool ... Wow all the surprises we got.... Go first okay put on the blindfold ready for more fun click the red subscribe button

Generated Story : A family had the perfect idea for a Paw Patrol-themed birthday party. The kids had a blast blowing bubbles in the backyard and playing in the pool with a mascot in a Paw Patrol costume. The birthday boy got to play with a pinata in the backyard and a hot air balloon flew over the house. The kids played with a toy bike and a Paw Patrol playhouse, and the birthday girl hugged a mascot. They also set up a LEGO train set and a DC Comics mascot in the pool. The birthday party ended with a game of catch and a pin the badge on the mascot game. Everyone had a great time, and the birthday cake was a Paw Patrol-themed delight. This vlog shows that family fun can be had at any birthday party with a little imagination.

| | | |
|--|---|---|
| Topic : Celebrations | Emotion : Cheerful | Persuasion Strategy: Social Identity, Social Proof |
| Action : Consider using your imagination to create a fun and memorable birthday party for your family or friends. | Reason : The reason is to create a uniquely fun-filled and memorable birthday party for your family or friends | |

Figure 2.14: An example of a story generated by the proposed pipeline along with the predicted outputs of the video-understanding tasks on the generated story. The generated story captures information across scenes, characters, event sequences, dialogues, emotions, and the environment. This helps the downstream models to get adequate information about the video to reason about it correctly. The original video can be watched at https://youtu.be/_amwPjAcoC8.

Additionally, we discard frames that are completely dark or white, as they may have high optical flow but lack informative content.

Using either of these methods, we obtain a set of frames that represent the events in the video. These frames are then processed by a pretrained BLIP-2 model [160]. The BLIP model facilitates scene understanding and verbalizes the scene by capturing its most salient aspects. We utilize two different prompts to extract salient information from the frames. The first prompt, “*Caption this image*”, is used to generate a caption that describes what is happening in the image, providing an understanding of the scene. The second prompt, “*Can you tell the objects that are present in the image?*”, helps identify and gather information about the objects depicted in each frame.

b. *Textual elements in frames:* We also extract the textual information present in the frames, as text often reinforces the message present in a scene and can also inform viewers on what to expect next [161]. For the OCR module, we sample every 10th frame extracted at the native frames-per-second of the video, and these frames are sent to PP-OCR [157]. We filter the OCR text and use only the unique words for further processing.

3. Text Representation of Audio: The next modality we utilize from

the video is the audio content extracted from it. We employ an Automatic Speech Recognition (ASR) module to extract transcripts from the audio. Since the datasets we worked with involved YouTube videos, we utilized the YouTube API to extract the closed caption transcripts associated with those videos.

4. Prompting: We employ the aforementioned modules to extract textual representations of various modalities present in a video. This ensures that we capture the audio, visual, text, and outside knowledge aspects of the video. Once the raw text is collected and processed, we utilize it to prompt a generative language model in order to generate a coherent story that represents the video. To optimize the prompting process and enable the generation of more detailed stories, we remove similar frame captions and optical character recognition (OCR) outputs, thereby reducing the overall prompt size.

The prompt template is given in Section 2.4.2.2. Through experimentation, we discovered that using concise, succinct instructions and appending the text input signals (such as frame captions, OCR, and automatic speech recognition) at the end significantly enhances the quality of video story generation. For shorter videos (up to 120 seconds), we utilize all available information to prompt the LLM for story generation. However, for longer videos, we limit the prompts to closed captions and sampled frame captions. The entire prompting pipeline is zero-shot and relies on pre-trained LLMs. In our story generation experiments, we employ GPT-3.5 [59], Flan-t5 [156], and Vicuna [162]. A temperature of 0.75 is used for LLM generation. The average length of the generated stories is 231.67 words. Subsequently, these generated stories are utilized for performing video understanding tasks.

2.4.2.2 Prompt format

For verbalization, a template prompt format has been used, including all the data components as objects, captions, asr, ocr, meta-data.

“Please write a coherent story based on the following video advertisement. Use only the information provided and make sure the story feels like a continuous narrative and at the end include one sentence about what product the advertisement was about. Do not include any details not mentioned in the prompt. Use the elements given below to create a coherent narrative, but don’t use them as it is. The advertisement for the company {company_name} The video is titled {title}, with captions that include {caption}, voice-over : {transcripts}, and object recognition descriptions : {ocr}. The following objects are present in the advertisement and should be used to help create the story: {objects} Please exclude any empty or stop words from the final text.”

For downstream tasks, a template prompt format with an instruction about the specific task, the previous generated verbalization and vocabulary

for the downstream task is prompted to the LLM. Here is the example for the topic detection task, for other tasks context and vocab were changed accordingly.

“Given {topics} identify the most relevant topic from the dictionary keys from topic_vocab related to the story of the video advertisement given below. Consider the definitions given with topics in the topic_vocab dictionary, to identify which topic is most relevant, don’t add any extra topics that are not given in dictionary keys and answer with just the most relevant topic. Story : {verbalization}”

| Training | Model | Topic | Emotion | | Persuasion | Action | Reason |
|----------------------|--|-------|------------|---------|------------|--------|--------|
| | | | All labels | Clubbed | | | |
| Random | Random | 2.63 | 3.37 | 14.3 | 8.37 | 3.34 | 3.34 |
| Finetuned | VideoMAE [163] | 24.72 | 29.72 | 85.55 | 11.17 | - | - |
| | Hussain <i>et al.</i> [79] | 35.1 | 32.8 | - | - | - | 48.45 |
| | Intern-Video [164] | 57.47 | 36.08 | 86.59 | 5.47 | 6.8 | 7.1 |
| Zero-shot | VideoChat [165] | 9.07 | 3.09 | 5.1 | 10.28 | - | - |
| Our Framework | GPT-3.5 Generated Story + GPT-3.5 Classifier | 51.6 | 11.68 | 79.69 | 35.02 | 66.27 | 59.59 |
| | GPT-3.5 Generated Story + Flan-t5-xxl Classifier | 60.5 | 10.8 | 79.10 | 33.41 | 79.22 | 81.72 |
| | GPT-3.5 Generated Story + Vicuna Classifier | 22.92 | 10.8 | 67.35 | 29.6 | 21.39 | 20.89 |
| | Vicuna Generated Story + GPT-3.5 Classifier | 46.7 | 5.9 | 80.33 | 27.54 | 61.88 | 55.44 |
| | Vicuna Generated Story + Flan-t5-xxl Classifier | 57.38 | 9.8 | 76.60 | 30.11 | 77.38 | 80.66 |
| | Vicuna Generated Story + Vicuna Classifier | 11.75 | 10.5 | 68.13 | 26.59 | 20.72 | 21.00 |
| | Generated Story + Roberta Classifier | 71.3 | 33.02 | 84.20 | 64.67 | 42.96* | 39.09* |

Table 2.6: Comparison of all the models across topic, emotion, and persuasion strategy detection tasks. We see that our framework, despite being zero-shot, outperforms finetuned video-based models on the topic classification, persuasion strategy detection and action and reason classification tasks and comes close on the emotion classification task. Further, the Roberta classifier trained on generated stories outperforms both finetuned and zero-shot models on most tasks. Best models are denoted in green and runner-ups in blue .

2.4.2.3 Results

Video Storytelling: The performance comparison between our pipeline and existing methods is presented in Table 2.7. We evaluate multiple generative and retrieval-based approaches and find that our pipeline achieves state-of-the-art results. It is important to note that as our method is entirely generative, the ROUGE-L score is lower compared to retrieval-based methods due to less overlap with ground truth reference video stories. However, overall metrics indicate that our generated stories exhibit a higher level of

| | Method | Model Type | METEOR | CIDEr | ROUGE-L | BLEU-1 | BLEU-2 | BLEU-3 | BLEU-4 |
|---------------|--------------------------------------|------------|--------|-------|---------|--------|--------|--------|--------|
| Random | Random | Retrieval | 13.1 | 30.2 | 21.4 | 43.1 | 23.1 | 10.0 | 4.8 |
| Finetuned | Narrator [152] | Retrieval | 19.6 | 98.4 | 29.5 | 69.1 | 43.0 | 25.3 | 15.0 |
| | EMB [152] | Retrieval | 19.1 | 88.8 | 28.9 | 64.5 | 39.3 | 22.7 | 13.4 |
| | BRNN [152] | Retrieval | 18.1 | 81.0 | 28.3 | 61.4 | 36.6 | 20.3 | 11.3 |
| | ResBRNN [152] | Retrieval | 19.6 | 94.3 | 29.7 | 66.0 | 41.7 | 24.3 | 14.7 |
| | Pseudo-GT+ | Retrieval | 20.1 | 103.6 | 29.9 | 69.1 | 43.5 | 26.1 | 15.6 |
| | ResBRNN-kNN [152] GVMF [166] | Retrieval | 20.7 | 107.7 | 30.8 | 70.5 | 44.3 | 26.9 | 15.9 |
| Zero-shot | VideoChat [165] | Generative | 15.49 | 42.9 | 17.88 | 50.00 | 43.30 | 34.76 | 27.21 |
| Zero-shot | GPT-3.5 | Generative | 24.8 | 102.4 | 24.3 | 63.8 | 56.4 | 47.2 | 38.6 |
| Our Framework | Vicuna | Generative | 17.4 | 73.9 | 20.9 | 70.49 | 60.0 | 48.25 | 38.20 |
| | Flan-t5-xxl | Generative | 4.8 | 34.6 | 10.58 | 7.9 | 6.8 | 5.4 | 4.3 |
| | Uniformly Sampled BLIP-2 Captions | Generative | 21.7 | 108.9 | 24.04 | 55.19 | 48.5 | 40.7 | 33.76 |

Table 2.7: Comparison on story generation task on the video-story dataset. We see that our framework despite being zero-shot outperforms all the finetuned generative prior art on all metrics. Further, it also outperforms finetuned retrieval models, which choose from a fixed set of frame descriptions on most metrics. Best models are denoted in green and runner-ups in blue.

| Task | Model | METEOR | CIDEr | ROUGE-L | BLEU-1 | BLEU-2 | BLEU-3 | BLEU-4 |
|---------------------|-------------|--------|-------|---------|--------|--------|--------|--------|
| Action | GPT-3.5 | 20.46 | 41.7 | 9.5 | 18.7 | 14.8 | 11.8 | 9.4 |
| Action | Flan-t5-xxl | 15.75 | 61.5 | 13.6 | 50.0 | 34.8 | 26.9 | 21.8 |
| Action | Vicuna | 21.20 | 42.6 | 7.6 | 16.8 | 13.08 | 10.08 | 7.7 |
| Reason | GPT-3.5 | 13.34 | 16.7 | 7.8 | 27.1 | 20.8 | 14.7 | 10.4 |
| Reason | Flan-t5-xxl | 8.35 | 24.9 | 5.9 | 39.4 | 24.7 | 16.7 | 12.0 |
| Reason | Vicuna | 15.82 | 27.9 | 7.75 | 24.6 | 19.3 | 14.1 | 10.3 |
| Reason given action | GPT-3.5 | 13.77 | 29.4 | 8.7 | 33.5 | 24.9 | 17.9 | 13.2 |
| Reason given action | Flan-t5-xxl | 4.29 | 19.0 | 7.6 | 23.2 | 15.0 | 10.2 | 7.5 |
| Reason given action | Vicuna | 13.62 | 24.4 | 7.61 | 22.6 | 17.7 | 12.8 | 9.2 |

Table 2.8: Comparison of the different zero-shot models on the action and reason generation tasks. Note that there are no fine-tuned generative models in the literature for this task and the number of annotated videos is too small to train a generative model. Best models are denoted in green.

| Training | Model | relationship | way_speaking | scene | like_ratio | view_count | director | genre | writer | year |
|-----------|----------------------------|--------------|---------------|-------|-------------|-------------|----------|-------|--------|-------|
| Trained | R101-slowfast+NL [154] | 52.4 | 35.8 | 54.7 | 0.386 | 3.77 | 44.9 | 53.0 | 36.3 | 52.5 |
| Trained | VideoBert [167] | 52.8 | 37.9 | 54.9 | 0.320 | 4.46 | 47.3 | 51.9 | 38.5 | 36.1 |
| Trained | Xiao <i>et al.</i> [168] | 50.95 | 34.07 | 44.19 | 0.353 | 4.886 | 40.19 | 48.11 | 31.43 | 29.65 |
| Trained | Qian <i>et al.</i> [169] | 50.95 | 32.86 | 32.56 | 0.444 | 4.600 | 37.76 | 48.17 | 27.26 | 25.31 |
| Trained | Object Transformers [154] | 53.1 | 39.4 | 56.9 | 0.230 | 3.55 | 51.2 | 54.6 | 34.5 | 39.1 |
| Zero-shot | GPT-3.5 generated | 64.1 | 39.07 | 60.2 | 0.061 | 12.84 | 69.9 | 58.1 | 52.4 | 75.6 |
| (Ours) | story + Flan-t5-xxl | | | | | | | | | |
| Zero-shot | GPT-3.5 generated | 68.42 | 32.95 | 54.54 | 0.031 | 12.69 | 75.26 | 50.84 | 32.16 | 75.96 |
| (Ours) | story + GPT-3.5 classifier | | | | | | | | | |
| Trained | GPT-3.5 generated | 62.16 | 38.41 | 68.65 | 0.054 | 11.84 | 45.34 | 39.27 | 35.93 | 7.826 |
| (Ours) | story + Roberta | | | | | | | | | |

Table 2.9: Comparison of various models on the LVU benchmark. We see that our framework, despite being zero-shot, outperforms fine-tuned video-based models on 8/9 tasks. Best models are denoted in green and runner-ups in blue.

| Training | Model | Scene | Object | Action | Event | Attribute | Concept | Overall |
|---------------------|---|-------|--------|--------|-------|-----------|---------|---------|
| Trained | 3D-Resnet | 50.6 | 28.6 | 48.2 | 35.9 | 29 | 22.5 | 35.8 |
| Trained | 3D-STCNet | 51.9 | 30.1 | 50.3 | 35.8 | 29.9 | 22.7 | 36.7 |
| Trained | HATNet | 55.8 | 34.2 | 51.8 | 38.5 | 33.6 | 26.1 | 40 |
| Trained | 3D-Resnet (Multitask) | 51.7 | 29.6 | 48.9 | 36.6 | 31.1 | 24.1 | 37 |
| Trained | HATNet (Multitask) | 57.2 | 35.1 | 53.5 | 39.8 | 34.9 | 27.3 | 41.3 |
| Zero-shot (Ours) | GPT-3.5 generated story + Flan-t5-xxl classifier | 59.66 | 98.89 | 98.96 | 38.42 | 67.76 | 86.99 | 75.12 |
| Zero-shot (Ours) | GPT-3.5 generated story + GPT-3.5 classifier | 60.2 | 99.16 | 98.72 | 40.79 | 67.17 | 88.6 | 75.77 |

Table 2.10: Comparison of various models on the HVU benchmark [155]. The models scores are as reported in [155]. We see that our framework, despite being zero-shot, outperforms fine-tuned video-based models on all the tasks. Best models are denoted in green and runner-ups in blue.

similarity to the reference stories and effectively capture the meaning of the source video.

Video Understanding: The performance comparison between our pipeline and other existing methods across six tasks (topic, emotion, and persuasion strategy classification, as well as action and reason retrieval and generation) is presented in Tables 2.6 and 2.8. Notably, our zero-shot model outperforms finetuned video-based baselines in all tasks except emotion classification. Further, our text-based finetuned model outperforms all other baselines on most of the tasks.

Unlike the story generation task, there are limited baselines available for video understanding tasks. Moreover, insufficient samples hinder training models from scratch. To address this, we utilize state-of-the-art video understanding models, VideoMAE and InternVideo. InternVideo shows strong performance on many downstream tasks. Analyzing the results, we observe that while GPT-3.5 and Vicuna perform similarly for story generation (Table 2.7), GPT-3.5 and Flan-t5 excel in downstream tasks (Table 2.6). Interestingly, although GPT-3.5 and Vicuna-generated stories yield comparable results, GPT-3.5 exhibits higher performance across most tasks. Vicuna-generated stories closely follow GPT-3.5 in terms of downstream task performance.

Next, we compare the best models (as in Table 2.6) on the LVU and HVU benchmarks with respect to the state-of-the-art models reported in the literature. Tables 2.9 and 2.10 report the results for the comparisons. As can be noted, the zero-shot models outperform most other baselines. For LVU, the zero-shot models work better than the trained Roberta-based classifier model. For HVU, we convert the classification task to a retrieval task, where in a zero-shot way, we input the verbalization of a video along with 30 randomly chosen tags containing an equal number of tags for each category (scene, object, action, event, attribute, and concept). The model is then prompted to pick the top 5 tags that seem most relevant to the video.

| Model | Topic | Emotion | | Persuasion | Action | Reason |
|------------------------------------|-------|------------|---------|------------|--------|--------|
| | | All labels | Clubbed | | | |
| BLIP-2 Captions + Flant-t5-xxl | 32.2 | 7.4 | 43.11 | 32.1 | 52.98 | 76.26 |
| BLIP-2 Captions + GPT-3.5 | 32.7 | 7.9 | 76.69 | 30.1 | 49.91 | 58.71 |
| Audio Transcription + Flant-t5-xxl | 49.37 | 10.1 | 63.56 | 21.9 | 66.17 | 79.68 |
| Audio Transcription + GPT-3.5 | 32.88 | 6.4 | 75.97 | 32.25 | 64.98 | 61.78 |

Table 2.11: Ablation study of using only visual (caption) or audio (transcripts) and LLMs for downstream tasks. It can be noted that the overall model does not perform as well (compared to Table 2.6) when using only audio or scene description without generating story.

These tags are mapped back to the main category tags, which are treated as the predicted labels.

Furthermore, as a comparative and ablation study of our approach, we evaluate the performance using only the BLIP-2 captions and audio transcriptions (Table 2.11). Our findings highlight that generated stories leveraging both audio and visual signals outperform those using vision or audio inputs alone. This emphasizes the significance of verbalizing a video in enhancing video understanding.

2.4.2.4 Ablation

Among the different components of information input present in the prompt, the LLM utilizes them differently while constructing the verbalization for the videos. For this experiment we use a subset of [79] dataset, considering videos that have spoken audio present.

We use ROUGE-l to get the longest common subsequence (LCS) between the generated verbalization and the individual components, which captures the overlapping content, providing an indication of their semantic similarity.

As generated verbalizations are abstractive as compared to extractive, we also use cosine similarity between the Roberta embeddings of the generated verbalization and the individual components.

We find that despite the order of the components in the prompt, the LLMs tend to utilize the audio components in the videos, in an extractive way.

2.4.2.5 A few examples of the stories generated using our method

1. <https://www.youtube.com/watch?v=1PdD8NvVfw0>: Kathy Ames had always wanted to pursue a doctoral degree but was unsure about the time commitment. When she discovered Grand Canyon University, she knew she had found the perfect fit. Grand Canyon University offered a flexible schedule that would allow her to balance her personal and family life with her studies. She - along with other students - gathered in the classroom,

| Model | Top-5 Accuracy | mAP |
|-----------------------------------|----------------|-------|
| VideoMAE | 25.57 | 24.79 |
| InternVideo | 7.477 | 15.62 |
| GPT-3.5 Generated Story + GPT-3.5 | 34.2 | 27.53 |
| Vicuna Generated Story + GPT-3.5 | 31.54 | 27.24 |
| GPT-3.5 Generated Story + Flant5 | 37 | 27.96 |
| Vicuna Generated Story + Flant5 | 31.13 | 27.32 |

Table 2.12: Top-5 accuracy, and mAP for persuasion strategy detection task

excitedly listening to their coach, Scott Saunders, explain the program. Afterward, Kathy made her way to the library and settled into a chair with her laptop.

She studied diligently, surrounded by her peers and classmates. In the evenings, she met with her peers around the table to discuss the topics of the day. Everyone was always eager to help and support each other. After a long day, Kathy made her way back to her living room where she relaxed on the couch with a glass of water and a lamp providing a soothing light.

Kathy was grateful for the opportunity to pursue her dream at Grand Canyon University. She was able to learn from experienced faculty and gain real-world experience that would prepare her for success after graduation.

The advertisement for Grand Canyon University was about offering a private, Christian education at an affordable price.

2. https://www.youtube.com/watch?v=f_6QQ6IVa6E: The woman holding the book stepped onto the patio and looked up to the sky. She was ready to take on the day. Taking out her phone, she opened the furniture catalog app, scrolling through the various designs. She quickly decided on the perfect pieces to brighten up her home. Next, she headed to The Home Depot for the supplies she needed. As she entered the store, the woman was delighted to find all the tools and materials she needed, from the Ryobi Cordless Vacuum to the Leaf Blower. She was even more excited when she spotted the Splatter an object recognition tool that allowed her to easily find the perfect paint color for her project. With her shopping done, the woman made her way to the checkout line with a cup of coffee in hand. She couldn't wait to get to work and make her home more beautiful. She knew that with the help of The Home Depot, Today was the Day for Doing. This advertisement was for The Home Depot - the one-stop-shop for all your home improvement needs.

3. <https://www.youtube.com/watch?v=PJ1HiQJBDMw>: The advertisement for the company Sherwin-Williams opens on a kitchen table strewn with shells and wicker baskets, with two glasses of iced water beside them. A vase with a blue pattern sits in the foreground, and a person holds up a

phone with the Sherwin-Williams logo on the screen. A girl appears from behind a white sheet, peeking out of a white tent as if to signify the timelessness and neutrality of this color. The voice-over begins, as the camera pans to a living room with a staircase, and then to a dining room with a white table, chairs, and a white vase. The words “Color of the Month: Shell White, Sherwin-Williams” appear on the screen, as the camera zooms in on the vase. The words are followed by Our app makes it a snap, referring to Color Snap, the company’s new way of painting a home. The advertisement ends with the Sherwin-Williams logo, emphasizing the company’s commitment to excellence in home painting. This advertisement was promoting the company’s color of the month, Shell White.

4. <https://www.youtube.com/watch?v=CDjBIt70fp4>: The story began with a green light glowing in the dark, symbolizing the presence of a powerful technology that can change the way we work. This technology was an advanced graphics card, the NVIDIA Quadro FX 1700. It was compared side-by-side with its successor, the Quadro 2000, and it was clear that the Quadro 2000 was far more powerful. The Quadro FX 1700 had a qt - x700 vs qt - x700 capacity, while the Quadro 2000 had a green light that shone brighter and further than before. As the comparison was being made, the results were clear: the Quadro 2000 was the superior product. This advertisement for the company NVIDIA showcased the power of the Quadro 2000, and the improved performance it could bring to an organization. The advertisement concluded with a green light, signaling that NVIDIA had the answer to improving workflows. The product the advertisement was about was the NVIDIA Quadro 2000.”

5. <https://www.youtube.com/watch?v=SXTSJC5NrLs>: Mary Kom was training hard to make her mark as a boxer, and she had nothing but her focus and determination to help her succeed. She was a force to be reckoned with in the boxing ring and her hard work was soon paying off. But she knew that she had to take her training to the next level if she wanted to win the boxing championship. That’s when she decided to take her own advice and step out of her comfort zone. She headed out to the woods, determined to make the most of her time in nature. The woods were a perfect place for Mary to challenge herself physically. She was soon standing on the rocks of a river, pushing her body to its limits with every movement. She was standing tall, her black top glinting in the sunlight, and her energy was palpable. The sun was setting, and it was the perfect time to relax and unwind. Finally, she ended her day with a yoga routine, her body twisting and turning in a beautiful display of strength and grace. It was a day of hard work and fun, and it left Mary feeling inspired and ready to take on anything. The video was an advertisement for the movie Mary Kom, released by Viacom18 Studios.”

2.4.2.6 Hallucinations Present In the Automatically Generated Stories

The current large language models , do hallucinate or produce imaginary outputs. However, they are able to capture the gist of the videos, helping improve on downstream tasks. Some examples are shown below:

1. https://youtu.be/_Usyxm9S_Jg:
 - (a) Story: The woman in the advertisement is standing in front of a red background, smiling brightly. She and a man sitting at a table look content, as if they've just come from a conversation. The Netflix logo is visible, indicating that the advertisement is for the company and its products. The voice-over for the advertisement introduces Selena Gomez as the star of the trailer for the popular Twilight Saga. Her popularity has made her a perfect choice to promote Netflix's products. The voice-over then informs viewers that the company's streaming service is not available. The man and woman in the advertisement seem to have a strong connection, and we can infer that they are friends. Perhaps they share an understanding of each other's struggles, since they both seem to be comfortable and happy with each other. The advertisement ends with the Netflix logo, leaving viewers with the knowledge that they can find solace and understanding in the company's products, no matter what they are going through. The advertisement was for Netflix."
 - (b) Hallucinations: The story falsely attributes "Selena Gomez" to "Twilight Saga". The LLMs often confuse named entity linkages.
2. <https://youtu.be/c6ihySYvKYM>:
 - (a) Story: The man in glasses was pointing at the Skytech Gaming computer, as he and the man in the headset prepared to start a game. The man with the headset was lost in the world of gaming, his fingers flying across the keyboard as he tried to achieve domination. Meanwhile, the man with glasses watched intently as the RGB TV lit up the room. The woman was playing a game of PC gaming, while the man with glasses and headphones was engrossed in a game of Esports. A man in a headset played a game of PS4, while another man held up a computer screen. A girl was playing a game of PC gaming, while a man was looking at a computer case with colorful lights. The man wearing glasses and headphones was sitting in front of a computer, playing a game of PC gaming. The woman wearing headphones was playing a game on a computer, and the man with glasses was smiling as he played a game of PS4. A person's hand was holding a gaming mouse,

while a woman was playing a game of Esports. Finally, a computer case with green lights on it and an Acer Aspire R7 MK2 PC case were in the mix, as the video ended with a man playing a guitar with a lamp and a purple light in the dark. The advertisement was for Skytech Gaming, a company offering gaming PCs and laptops that are more than just a pastime - they are a passion and a battle of the heart. One sentence summarizing the video would be: Skytech Gaming shows that their gaming PCs and laptops can bring out the best of gamers, allowing them to take on every battle with anticipation, adrenaline, and camaraderie.

- (b) Hallucination: The entity “Acer Aspire R7 MK2” and “PS4” are not referenced in the prompt and are hallucinated by the LLM

3. <https://youtu.be/SZwWT51nHf8>

- (a) Story: A young girl is walking through her bedroom when she notices a man running up the stairs in his sneakers. She is intrigued and follows him, curious to see where he is going. Arriving at the top of the stairs, she finds a young man standing with his backpack. Across the room, she sees a girl in a dress standing in front of a mirror. The girl smiles and looks content. At that moment, the voice-over of the Academy Sports + Outdoors ad can be heard: “Academy Sports and Outdoors is the place to get all the best brands in thousands of combinations for back to school. Right stuff, low price, every day.” Suddenly, a black cat can be seen flying in the sky and the girl in the tank top is now standing in front of the mirror. The man is wearing blue pants and sneakers, and the shoes, rug, and sneakers of the apartment are visible. The girl smiles and looks content. The scene shifts to the bedroom where the bed, dresser, and mirror are present, along with a sign for Academy Sportswear. This advertisement for Academy Sports + Outdoors was about the company’s back-to-school collection of products. It showcased how Academy Sports + Outdoors provides the right stuff, at a low price, every day, for students to get ready for school.
- (b) The scene with a black cat is hallucinated though not mentioned in the prompt

4. <https://youtu.be/B63ubKGNqos>

- (a) Story: The advertisement by Netflix, titled ”Many Worlds, One Home”, is a stunning vision of many different universes. We see a man standing on a red moon, a desert landscape, and a live performance by Adam Sandler. Purple backgrounds with the word ”worlds” convey the idea of many possibilities, while a boy playing with a ball in front of a statue shows a more peaceful and homely side to this. As the night falls, a man is standing on a car, and the Ultraman. There are also glimpses of Tse

Seeds, Devilman Crybaby, Saint Seiya: Knights of the Zodiac, Eden Vs. Eden, Carole and Tuesday, Aggretsuko, Seis Manos, Godzilla, Castlevania, B: The Beginning, Baki, Cannon Busters, Rilakkuma and Kaoru, Beastars, and Ultramarine Magmell. We see a teddy bear, knights of the zodiac, a character walking in the water, a man in a snowy area, a pair of shoes, a man running away from the fire, the moon, a man wearing a suit with a green light, eyes being the only thing that can be seen, birds flying in the sky, a red and black background with a red light, bats and castles, the eyes, a man holding a gun, a man flying in the air, a face, a girl with long hair, and two hands. All of this culminates to signify the many different worlds that one can find "home", and the advertisement concludes with the Netflix logo. This advertisement for Netflix captures the idea of many possibilities and of being at home in a single place.

- (b) Hallucination: Blip2 hallucinates entities such as "Adam Sandler" who is not associated with the advertisement. But it causes the generated story to have the same association.

5. <https://youtu.be/WCIEmrN8ot4>

- (a) Story: The sun shone brightly in the savanna as the cheetah cub frolicked in the grass. Suddenly, a cheetah emerged from the shadows. The cheetah and its cub plan to chase the gazelles. Despite the presence of death, life in the wild continued on, just as it had for centuries. This story is featured in the advertisement for Walt Disney Studios, which captures the beauty and unpredictable nature of life in the African savanna. This advertisement ultimately serves to promote the film African Cats, which tells the story of a cheetah family as they strive to protect their respective prides against a variety of obstacles in the wild.
- (b) Hallucination: The generated caption "a fox cub is playing in the grass cheetah eats in the wild" was insufficient to capture the overall information of the video, however, due to the presence of ASR and world knowledge of the LLM. The LLM is able to generate a coherent story.

6. <https://youtu.be/Lh5n4HcLvCU>

- (a) Story: The woman in the Ulta Beauty advertisement skillfully applied makeup to her face, creating a beautiful pastel tint. With each brush stroke, her look became brighter and more dreamy. She beamed with confidence and joy, her face glowing with beauty. As the advertisement continued, she highlighted her eyes and cheeks with subtle pastel shades, adding to the stunning effect. With each addition, her face lit up with a vibrant luminosity. When she was finished, her face was a work of art, the perfect combination of bright and dreamy. The Ulta Beauty advertise-

ment showcased the power of a perfect pastel tint, and how it can enhance any look. This advertisement was for Ulta Beauty.

- (b) Hallucination: The video contains only music. However, the generated verbalization can pick the advertisement's gist.

7. <https://youtu.be/VdhRmVnbQo>

- (a) Story: Jackie had been wanting diamond hoop earrings for what seemed like forever. She had been dreaming of the perfect pair that she could wear on special occasions and that would bring out the best in her. One day Jackie happened to come across an advertisement for Costco, which featured a pair of Round Brilliant Diamond Hoop Earrings in 14kt White Gold. She immediately knew this was the perfect pair for her. The advertisement had a voice-over that said, "These earrings sparkle with glamour and sophistication." Jackie was mesmerized by the sparkle of the diamonds, and she was sure anyone who saw them would also be taken in by their beauty. Jackie quickly decided to purchase the earrings. She was sure they would make the perfect accessory for any special occasion. From date nights to family gatherings, she knew these earrings would make her look simply stunning. The advertisement for the Costco Round Brilliant Diamond Hoop Earrings in 14kt White Gold proved to be the perfect choice for Jackie. They not only added a touch of glamour and sophistication to her wardrobe, but made her feel extra special when she wore them. The advertisement for the company Costco was for a pair of Round Brilliant Diamond Hoop Earrings in 14kt White Gold.
- (b) Hallucination: The video has no audio transcription, and minimal visual information. However the LLM hallucinates a story, maintaining the essence of the advertisement. It also attributes a false voice-over which is not present in the video.

2.5 Conclusion

What does an advertisement say that makes people change their beliefs and actions? With limited works, the computational study of rhetoric of this all-pervasive form of marketing communication is still in its infancy. In this chapter, based on the well-developed social psychology and marketing literature, we develop and release the largest vocabulary of persuasion strategies and labeled dataset. We develop several models for predicting persuasion strategies for video and image based ads. Further, we show the performance of these models on several other advertisement-understanding related tasks, including topic, emotion, and question-answering.

Chapter 3

Content and Behavior Models

Chapter 4

Generating Content Leading to Optimal Behavior

Publications

1. Khurana, V., Kumar, Y., Hollenstein, N., Kumar, R., & Krishnamurthy, B. (2023). Synthesizing Human Gaze Feedback for Improved NLP Performance. In Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics (pp. 1895–1908).
2. Kumar, Y., Jha, R., Gupta, A., Aggarwal, M., Garg, A., Malyan, T., Bhardwaj, A., Ratn Shah, R., Krishnamurthy, B., & Chen, C. (2023). Persuasion Strategies in Advertisements. Proceedings of the AAAI Conference on Artificial Intelligence, 37(1), 57-66. <https://doi.org/10.1609/aaai.v37i1.25076>
3. Bhattacharya, A., Singla, Y. K., Krishnamurthy, B., Shah, R. R., & Chen, C. (2023). A Video Is Worth 4096 Tokens: Verbalize Videos To Understand Them In Zero Shot. In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, pages 9822–9839, Singapore. Association for Computational Linguistics. (Nominated for the best paper award!)
4. Khandelwal, A., Agrawal, A., Bhattacharyya, A., Singla, Y.K., Singh, S., Bhattacharya, U., Dasgupta, I., Petrangeli, S., Shah, R.R., Chen, C. and Krishnamurthy, B., 2023. Large Content And Behavior Models To Understand, Simulate, And Optimize Content And Behavior. arXiv preprint arXiv:2309.00359. (Under review).
5. S I, H., Singh, S., K Singla, Y., Krishnamurthy, B., Chen, C., Baths V., & Ratn Shah, R. (2023). Sharingan: How Much Will Your Customers Remember Your Brands After Seeing Your Ads?. arxiv preprint (Under review).
6. Khurana, V., Singla, Y.K., Subramanian, J., Shah, R.R., Chen, C., Xu, Z. and Krishnamurthy, B., 2023. Behavior Optimized Image Generation. arXiv preprint arXiv:2311.10995. (Under review)

Bibliography

- [1] Claude E. Shannon and Warren Weaver. *The mathematical theory of communication*. The mathematical theory of communication. University of Illinois Press, Champaign, IL, US, 1949.
- [2] Leo Breiman. Statistical modeling: The two cultures (with comments and a rejoinder by the author). *Statistical science*, 16(3):199–231, 2001.
- [3] Jake M Hofman, Amit Sharma, and Duncan J Watts. Prediction and explanation in social systems. *Science*, 355(6324):486–488, 2017.
- [4] Galit Shmueli. To explain or to predict? *Statistical Science*, 2010.
- [5] Stanley Milgram and Christian Gudehus. Obedience to authority, 1978.
- [6] Solomon E Asch. The doctrine of suggestion, prestige and imitation in social psychology. *Psychological review*, 55(5):250, 1948.
- [7] Open Science Collaboration. Estimating the reproducibility of psychological science. *Science*, 349(6251):aac4716, 2015.
- [8] Philip E Tetlock. Expert political judgment. In *Expert Political Judgment*. Princeton University Press, 2017.
- [9] The Forecasting Collaborative. Insights into the accuracy of social scientists’ forecasts of societal change. *Nature human behaviour*, 7(4):484–501, 2023.
- [10] Chenhao Tan, Lillian Lee, and Bo Pang. The effect of wording on message propagation: Topic-and author-controlled natural experiments on twitter. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 175–185, 2014.
- [11] Phillip Isola, Jianxiong Xiao, Devi Parikh, Antonio Torralba, and Aude Oliva. What makes a photograph memorable? *IEEE transactions on pattern analysis and machine intelligence*, 36(7):1469–1482, 2013.

- [12] Ivan Habernal and Iryna Gurevych. What makes a convincing argument? empirical analysis and detecting attributes of convincingness in web argumentation. In *EMNLP*, pages 1214–1223, 2016.
- [13] Yaman Kumar, Rajat Jha, Arunim Gupta, Milan Aggarwal, Aditya Garg, Tushar Malyan, Ayush Bhardwaj, Rajiv Ratn Shah, Balaji Krishnamurthy, and Changyou Chen. Persuasion strategies in advertisements. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37, pages 57–66, 2023.
- [14] Kelvin Luu, Chenhao Tan, and Noah A Smith. Measuring online debaters’ persuasive skill from text over time. *Transactions of the Association for Computational Linguistics*, 7:537–550, 2019.
- [15] Aanisha Bhattacharyya, Yaman K Singla, Balaji Krishnamurthy, Rajiv Ratn Shah, and Changyou Chen. A video is worth 4096 tokens: Verbalize videos to understand them in zero shot. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 9822–9839, Singapore, December 2023. Association for Computational Linguistics.
- [16] Justin Cheng, Lada Adamic, P Alex Dow, Jon Michael Kleinberg, and Jure Leskovec. Can cascades be predicted? In *Proceedings of the 23rd international conference on World wide web*, pages 925–936, 2014.
- [17] Travis Martin, Jake M Hofman, Amit Sharma, Ashton Anderson, and Duncan J Watts. Exploring limits to prediction in complex social systems. In *Proceedings of the 25th international conference on world wide web*, pages 683–694, 2016.
- [18] Hyunyoung Choi and Hal Varian. Predicting the present with google trends. *Economic record*, 88:2–9, 2012.
- [19] Chaoming Song, Zehui Qu, Nicholas Blumm, and Albert-László Barabási. Limits of predictability in human mobility. *Science*, 327(5968):1018–1021, 2010.
- [20] Matthew J Salganik. *Bit by bit: Social research in the digital age*. Princeton University Press, 2019.
- [21] Yaman Kumar Singla, Jui Shah, Changyou Chen, and Rajiv Ratn Shah. What do audio transformers hear? probing their representations for language delivery & structure. In *2022 IEEE International Conference on Data Mining Workshops (ICDMW)*, pages 910–925. IEEE, 2022.

- [22] Shobhana Chandra, Sanjeev Verma, Weng Marc Lim, Satish Kumar, and Naveen Donthu. Personalization in personalized marketing: Trends and ways forward. *Psychology & Marketing*, 39(8):1529–1562, 2022.
- [23] Jonathan L Herlocker, Joseph A Konstan, Loren G Terveen, and John T Riedl. Evaluating collaborative filtering recommender systems. *ACM Transactions on Information Systems (TOIS)*, 22(1):5–53, 2004.
- [24] H Brendan McMahan, Gary Holt, David Sculley, Michael Young, Dietmar Ebner, Julian Grady, Lan Nie, Todd Phillips, Eugene Davydov, Daniel Golovin, et al. Ad click prediction: a view from the trenches. In *Proceedings of the 19th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 1222–1230, 2013.
- [25] Reid Pryzant, Youngjoo Chung, and Dan Jurafsky. Predicting sales from the language of product descriptions. *eCOM@ SIGIR*, 2311, 2017.
- [26] Phillip Isola, Jianxiong Xiao, Antonio Torralba, and Aude Oliva. What makes an image memorable? In *CVPR 2011*, pages 145–152. IEEE, 2011.
- [27] Aditya Khosla, Akhil S Raju, Antonio Torralba, and Aude Oliva. Understanding and predicting image memorability at a large scale. In *Proceedings of the IEEE international conference on computer vision*, pages 2390–2398, 2015.
- [28] Harini SI, Somesh Singh, Yaman K Singla, Aanisha Bhattacharyya, Veeky Baths, Changyou Chen, Rajiv Ratn Shah, and Balaji Krishnamurthy. Long-term ad memorability: Understanding & generating memorable ads. *arXiv preprint arXiv:2309.00378*, 2023.
- [29] Kate Newstead and Jenni Romaniuk. Cost per second: The relative effectiveness of 15-and 30-second television advertisements. *Journal of Advertising Research*, 50(1):68–76, 2010.
- [30] Barbara Plank. Keystroke dynamics as signal for shallow syntactic parsing. In *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers*, pages 609–619, 2016.
- [31] Sigrid Klerke, Yoav Goldberg, and Anders Søgaard. Improving sentence compression by learning to predict gaze. In *NAACL: Human Language Technologies*, San Diego, California, June 2016. Association for Computational Linguistics.

- [32] Varun Khurana, Yaman Kumar, Nora Hollenstein, Rajesh Kumar, and Balaji Krishnamurthy. Synthesizing human gaze feedback for improved NLP performance. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 1895–1908, Dubrovnik, Croatia, May 2023. Association for Computational Linguistics.
- [33] Julia Kreutzer, Shahram Khadivi, Evgeny Matusov, and Stefan Riezler. Can neural machine translation be improved with user feedback? *arXiv preprint arXiv:1804.05958*, 2018.
- [34] Nisan Stiennon, Long Ouyang, Jeffrey Wu, Daniel Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul F Christiano. Learning to summarize with human feedback. *Advances in Neural Information Processing Systems*, 33:3008–3021, 2020.
- [35] Daniel M Ziegler, Nisan Stiennon, Jeffrey Wu, Tom B Brown, Alec Radford, Dario Amodei, Paul Christiano, and Geoffrey Irving. Fine-tuning language models from human preferences. *arXiv preprint arXiv:1909.08593*, 2019.
- [36] Reiichiro Nakano, Jacob Hilton, Suchir Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christopher Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, et al. Webgpt: Browser-assisted question-answering with human feedback. *arXiv preprint arXiv:2112.09332*, 2021.
- [37] Kimin Lee, Hao Liu, Moonkyung Ryu, Olivia Watkins, Yuqing Du, Craig Boutilier, Pieter Abbeel, Mohammad Ghavamzadeh, and Shixiang Shane Gu. Aligning text-to-image models using human feedback. *arXiv preprint arXiv:2302.12192*, 2023.
- [38] John David Pressman, Katherine Crowson, and Simulacra Captions Contributors. Simulacra aesthetic captions, 2023.
- [39] Xiaoshi Wu, Keqiang Sun, Feng Zhu, Rui Zhao, and Hongsheng Li. Better aligning text-to-image models with human preference. *arXiv preprint arXiv:2303.14420*, 2023.
- [40] Varun Khurana, Yaman K Singla, Jayakumar Subramanian, Rajaiv Ratn Shah, Changyou Chen, Zhiqiang Xu, and Balaji Krishnamurthy. Behavior optimized image generation. *arXiv preprint arXiv:2311.10995*, 2023.
- [41] Denis McQuail. *Mass communication theory: An introduction*. Sage Publications, Inc, 1987.
- [42] Klaus Krippendorff. *Content analysis: An introduction to its methodology*. Sage publications, 2018.

- [43] Harold D Lasswell. The structure and function of communication in society. *The communication of ideas*, 37(1):136–139, 1948.
- [44] Harold D Lasswell. *Propaganda technique in world war I*. MIT press, 1971.
- [45] Richard E Petty, John T Cacioppo, and Martin Heesacker. Effects of rhetorical questions on persuasion: A cognitive response analysis. *Journal of personality and social psychology*, 40(3):432, 1981.
- [46] Shelly Chaiken. Heuristic versus systematic information processing and the use of source versus message cues in persuasion. *Journal of personality and social psychology*, 39(5), 1980.
- [47] Andranik Tumasjan, Timm Sprenger, Philipp Sandner, and Isabell Welpe. Predicting elections with twitter: What 140 characters reveal about political sentiment. In *Proceedings of the international AAAI conference on web and social media*, volume 4, pages 178–185, 2010.
- [48] Sitaram Asur and Bernardo A Huberman. Predicting the future with social media. In *2010 IEEE/WIC/ACM international conference on web intelligence and intelligent agent technology*, volume 1, pages 492–499. IEEE, 2010.
- [49] Seungbae Kim, Jyun-Yu Jiang, Masaki Nakada, Jinyoung Han, and Wei Wang. Multimodal post attentive profiling for influencer marketing. In *Proceedings of The Web Conference 2020*, pages 2878–2884, 2020.
- [50] Yan Carrière-Swallow and Felipe Labbé. Nowcasting with google trends in an emerging market. *Journal of Forecasting*, 32(4):289–298, 2013.
- [51] Nicholas Generous, Geoffrey Fairchild, Alina Deshpande, Sara Y Del Valle, and Reid Priedhorsky. Global disease monitoring and forecasting with wikipedia. *PLoS computational biology*, 10(11):e1003892, 2014.
- [52] Giovanni De Toni, Cristian Consonni, and Alberto Montresor. A general method for estimating the prevalence of influenza-like-symptoms with wikipedia data. *Plos one*, 16(8):e0256858, 2021.
- [53] Márton Mestyán, Taha Yasseri, and János Kertész. Early prediction of movie box office success based on wikipedia activity big data. *PloS one*, 8(8):e71226, 2013.
- [54] Coco Krumme, Alejandro Llorente, Manuel Cebrian, Alex Pentland, and Esteban Moro. The predictability of consumer visitation patterns. *Scientific reports*, 3(1):1645, 2013.

- [55] Yves-Alexandre De Montjoye, Laura Radaelli, Vivek Kumar Singh, and Alex “Sandy” Pentland. Unique in the shopping mall: On the reidentifiability of credit card metadata. *Science*, 347(6221):536–539, 2015.
- [56] Dirk Brockmann, Lars Hufnagel, and Theo Geisel. The scaling laws of human travel. *Nature*, 439(7075):462–465, 2006.
- [57] Giovanna Miritello, Rubén Lara, Manuel Cebrian, and Esteban Moro. Limited communication capacity unveils strategies for human interaction. *Scientific reports*, 3(1):1950, 2013.
- [58] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *NAACL: Human Language Technologies, Volume 1 (Long and Short Papers)*. ACL, 2019.
- [59] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners, 2020.
- [60] Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language understanding by generative pre-training. 2018.
- [61] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research*, 21(1):5485–5551, 2020.
- [62] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- [63] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*, 2013.
- [64] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman

- Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- [65] Ashmit Khandelwal, Aditya Agrawal, Aanisha Bhattacharyya, Yamna K Singla, Somesh Singh, Uttaran Bhattacharya, Ishita Dasgupta, Stefano Petrangeli, Rajiv Ratn Shah, Changyou Chen, et al. Large content and behavior models to understand, simulate, and optimize content and behavior. *arXiv preprint arXiv:2309.00359*, 2023.
- [66] OpenAI. Gpt-4 technical report, 2023.
- [67] Stella Biderman, Kieran Bicheno, and Leo Gao. Datasheet for the pile. *arXiv preprint arXiv:2201.07311*, 2022.
- [68] Guilherme Penedo, Quentin Malartic, Daniel Hesslow, Ruxandra Cojocaru, Alessandro Cappelli, Hamza AlObaidli, Baptiste Pannier, Ebtesam Almazrouei, and Julien Launay. The refinedweb dataset for falcon llm: outperforming curated corpora with web data, and web data only. *arXiv preprint arXiv:2306.01116*, 2023.
- [69] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *arXiv preprint arXiv:2304.08485*, 2023.
- [70] Christof Rapp. Aristotle’s rhetoric. 2002.
- [71] Joan Meyers-Levy and Prashant Malaviya. Consumers’ processing of persuasive advertisements: An integrative framework of persuasion theories. *Journal of marketing*, 63(4_suppl1):45–60, 1999.
- [72] Punam Anand Keller, Isaac M Lipkus, and Barbara K Rimer. Affect, framing, and persuasion. *Journal of Marketing Research*, 40(1):54–64, 2003.
- [73] Carl Iver Hovland, Irving Lester Janis, and Harold H Kelley. Communication and persuasion. 1953.
- [74] Richard E Petty and John T Cacioppo. The elaboration likelihood model of persuasion. In *Communication and persuasion*, pages 1–24. Springer, 1986.
- [75] Chongyang Bai, Haipeng Chen, Srijan Kumar, Jure Leskovec, and VS Subrahmanian. M2p2: Multimodal persuasion prediction using adaptive fusion. *IEEE Transactions on Multimedia*, 2021.
- [76] Jungseock Joo, Weixin Li, Francis F Steen, and Song-Chun Zhu. Visual persuasion: Inferring communicative intents of images. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 216–223, 2014.

- [77] Jiaao Chen and Diyi Yang. Weakly-supervised hierarchical models for predicting persuasive strategies in good-faith textual requests. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(14):12648–12656, May 2021.
- [78] Paul AM Van Lange, E Tory Higgins, and Arie W Kruglanski. Handbook of theories of social psychology. *Handbook of Theories of Social Psychology*, pages 1–568, 2011.
- [79] Zaeem Hussain, Mingda Zhang, Xiaozhong Zhang, Keren Ye, Christopher Thomas, Zuha Agha, Nathan Ong, and Adriana Kovashka. Automatic understanding of image and video advertisements. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1705–1715, 2017.
- [80] Pranav Anand, Joseph King, Jordan Boyd-Graber, Earl Wagner, Craig Martell, Doug Oard, and Philip Resnik. Believe me—we can do this! annotating persuasive acts in blog text. In *Workshops at the Twenty-Fifth AAAI Conference on Artificial Intelligence*, 2011.
- [81] Demetrios Vakratsas and Tim Ambler. How advertising works: what do we really know? *Journal of marketing*, 63(1):26–43, 1999.
- [82] Chenhao Tan, Vlad Niculae, Cristian Danescu-Niculescu-Mizil, and Lillian Lee. Winning arguments: Interaction dynamics and persuasion strategies in good-faith online discussions. In *Proceedings of the 25th international conference on world wide web*, pages 613–624, 2016.
- [83] Jiaao Chen and Diyi Yang. Weakly-supervised hierarchical models for predicting persuasive strategies in good-faith textual requests. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 12648–12656, 2021.
- [84] Omar Shaikh, Jiaao Chen, Jon Saad-Falcon, Polo Chau, and Diyi Yang. Examining the ordering of rhetorical strategies in persuasive requests. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1299–1306, Online, November 2020. Association for Computational Linguistics.
- [85] Jialu Li, Esin Durmus, and Claire Cardie. Exploring the role of argument structure in online debate persuasion. In *EMNLP*, pages 8905–8912, Online, November 2020. Association for Computational Linguistics.
- [86] Stephanie Lukin, Pranav Anand, Marilyn Walker, and Steve Whitaker. Argument strength is in the eye of the beholder: Audience effects in persuasion. In *Proceedings of the 15th Conference of the*

European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers, pages 742–753, Valencia, Spain, April 2017. Association for Computational Linguistics.

- [87] Roxanne El Baff, Henning Wachsmuth, Khalid Al Khatib, and Benno Stein. Analyzing the Persuasive Effect of Style in News Editorial Argumentation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 3154–3160, Online, July 2020. Association for Computational Linguistics.
- [88] Tim Althoff, Cristian Danescu-Niculescu-Mizil, and Dan Jurafsky. How to ask for a favor: A case study on the success of altruistic requests. In *Proceedings of the International AAAI Conference on Web and Social Media*, volume 8, pages 12–21, 2014.
- [89] Zhongyu Wei, Yang Liu, and Yi Li. Is this post persuasive? ranking argumentative comments in online forum. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 195–200, 2016.
- [90] Ivan Donadello, Mauro Dragoni, and Claudio Eccher. Explaining reasoning algorithms with persuasiveness: a case study for a behavioural change system. In *Proceedings of the 35th Annual ACM Symposium on Applied Computing*, pages 646–653, 2020.
- [91] Emmanuel Hadoux, Anthony Hunter, and Sylwia Polberg. Strategic argumentation dialogues for persuasion: Framework and experiments based on modelling the beliefs and concerns of the persuadee. *arXiv preprint arXiv:2101.11870*, 2021.
- [92] Ivan Donadello, Anthony Hunter, Stefano Teso, and Mauro Dragoni. Machine learning for utility prediction in argument-based computational persuasion. In *proceedings of the AAAI conference on artificial intelligence*, volume 36, pages 5592–5599, 2022.
- [93] Jungseock Joo, Francis F Steen, and Song-Chun Zhu. Automated facial trait judgment and election outcome prediction: Social dimensions of face. In *Proceedings of the IEEE international conference on computer vision*, pages 3712–3720, 2015.
- [94] Behjat Siddiquie, Dave Chisholm, and Ajay Divakaran. Exploiting multimodal affect and semantics to identify politically persuasive web videos. In *Proceedings of the 2015 ACM on International Conference on Multimodal Interaction*, pages 203–210, 2015.
- [95] Elliot Aronson, Judith A Turner, and J Merrill Carlsmith. Communicator credibility and communication discrepancy as determinants

- of opinion change. *The Journal of Abnormal and Social Psychology*, 67(1):31, 1963.
- [96] Robert B Cialdini and Robert B Cialdini. *Influence: The psychology of persuasion*, volume 55. Collins New York, 2007.
- [97] Stanley Milgram. Behavioral study of obedience. *The Journal of abnormal and social psychology*, 67(4):371, 1963.
- [98] Elliott McGinnies and Charles D Ward. Better liked than right: Trustworthiness and expertise as factors in credibility. *Personality and Social Psychology Bulletin*, 6(3):467–472, 1980.
- [99] Kim Giffin. The contribution of studies of source credibility to a theory of interpersonal trust in the communication process. *Psychological bulletin*, 68(2):104, 1967.
- [100] Rahul Radhakrishnan Iyer and Katia Sycara. An unsupervised domain-independent framework for automated detection of persuasion tactics in text. *arXiv preprint arXiv:1912.06745*, 2019.
- [101] Henning Wachsmuth, Nona Naderi, Yufang Hou, Yonatan Bilu, Vinodkumar Prabhakaran, Tim Alberdingk Thijm, Graeme Hirst, and Benno Stein. Computational argumentation quality assessment in natural language. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers*, pages 176–187, 2017.
- [102] Esin Durmus and Claire Cardie. Exploring the role of prior beliefs for argument persuasion. In *NAACL: Human Language Technologies, Volume 1 (Long Papers)*, pages 1035–1045, New Orleans, Louisiana, June 2018. ACL.
- [103] Morton Deutsch and Harold B Gerard. A study of normative and informational social influences upon individual judgment. *The journal of abnormal and social psychology*, 51(3):629, 1955.
- [104] Richard E Petty, Duane T Wegener, and Leandre R Fabrigar. Attitudes and attitude change. *Annual review of psychology*, 48(1):609–647, 1997.
- [105] Wendy Wood. Attitude change: Persuasion and social influence. *Annual review of psychology*, 51(1):539–570, 2000.
- [106] Robert B Cialdini and Noah J Goldstein. Social influence: Compliance and conformity. *Annual review of psychology*, 55(1):591–621, 2004.

- [107] Nataly Levesque and Frank Pons. The human brand: A systematic literature review and research agenda. *Journal of Customer Behaviour*, 19(2):143–174, 2020.
- [108] Sara Rosenthal and Kathleen McKeown. Detecting influencers in multiple online genres. *ACM Transactions on Internet Technology (TOIT)*, 17(2):1–22, 2017.
- [109] Diyi Yang, Jiaao Chen, Zichao Yang, Dan Jurafsky, and Eduard Hovy. Let’s make your request more persuasive: Modeling persuasive strategies via semi-supervised neural nets on crowdfunding platforms. In *NAACL: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 3620–3630, 2019.
- [110] Fan Zhang, Diane Litman, and Kate Forbes-Riley. Inferring discourse relations from pdtb-style discourse labels for argumentative revision classification. In *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers*, pages 2615–2624, 2016.
- [111] Christian Stab and Iryna Gurevych. Parsing argumentation structures in persuasive essays. *Computational Linguistics*, 43(3):619–659, 2017.
- [112] Christopher Hidey, Elena Musi, Alyssa Hwang, Smaranda Muresan, and Kathleen McKeown. Analyzing the semantic types of claims and premises in an online persuasive forum. In *Proceedings of the 4th Workshop on Argument Mining*, pages 11–21, 2017.
- [113] Dennis T Regan. Effects of a favor and liking on compliance. *Journal of experimental social psychology*, 7(6):627–639, 1971.
- [114] Margaret S Clark. Record keeping in two types of relationships. *Journal of personality and social psychology*, 47(3), 1984.
- [115] Margaret S Clark and Judson Mills. Interpersonal attraction in exchange and communal relationships. *Journal of personality and social psychology*, 37(1), 1979.
- [116] Margaret S Clark, Judson Mills, and Martha C Powell. Keeping track of needs in communal and exchange relationships. *Journal of personality and social psychology*, 51(2):333, 1986.
- [117] Jonathan L Freedman and Scott C Fraser. Compliance without pressure: the foot-in-the-door technique. *Journal of personality and social psychology*, 4(2):195, 1966.
- [118] Jerry M Burger. The foot-in-the-door compliance procedure: A multiple-process analysis and review. *Personality and social psychology review*, 3(4):303–325, 1999.

- [119] Xuewei Wang, Weiyang Shi, Richard Kim, Yoojung Oh, Sijia Yang, Jingwen Zhang, and Zhou Yu. Persuasion for good: Towards a personalized persuasive dialogue system for social good. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5635–5649, Florence, Italy, July 2019. Association for Computational Linguistics.
- [120] John Paul Vargheese, Matthew Collinson, and Judith Masthoff. Exploring susceptibility measures to persuasion. In *Persuasive Technology. Designing for Future Change: 15th International Conference on Persuasive Technology, PERSUASIVE 2020, Aalborg, Denmark, April 20–23, 2020, Proceedings 15*, pages 16–29. Springer, 2020.
- [121] William J McGuire and Demetrios Papageorgis. The relative efficacy of various types of prior belief-defense in producing immunity against persuasion. *Journal of Abnormal Psychology*, 62(2), 1961.
- [122] Eric S Knowles and Jay A Linn. *Resistance and persuasion*. Psychology Press, 2004.
- [123] William J McGuire. Inducing resistance to persuasion. some contemporary approaches. *CC Haaland and WO Kaelber (Eds.), Self and Society. An Anthology of Readings, Lexington, Mass.(Ginn Custom Publishing) 1981, pp. 192-230.*, 1964.
- [124] Angela Y Lee, Punam Anand Keller, and Brian Sternthal. Value from regulatory construal fit: The persuasive impact of fit between consumer goals and message concreteness. *Journal of Consumer Research*, 36(5):735–747, 2010.
- [125] Adrian Furnham and Hua Chu Boo. A literature review of the anchoring effect. *The journal of socio-economics*, 40(1):35–42, 2011.
- [126] Duane T Wegener, Richard E Petty, Brian T Detweiler-Bedell, and W Blair G Jarvis. Implications of attitude change theories for numerical anchoring: Anchor plausibility and the limits of anchor effectiveness. *Journal of Experimental Social Psychology*, 37(1):62–69, 2001.
- [127] Amos Tversky and Daniel Kahneman. Judgment under uncertainty: Heuristics and biases. *science*, 185(4157):1124–1131, 1974.
- [128] Fritz Strack and Thomas Mussweiler. Explaining the enigmatic anchoring effect: Mechanisms of selective accessibility. *Journal of personality and social psychology*, 73(3):437, 1997.

- [129] Chitrabhan B Bhattacharya and Sankar Sen. Consumer–company identification: A framework for understanding consumers’ relationships with companies. *Journal of marketing*, 67(2):76–88, 2003.
- [130] Amy Zhang, Bryan Culbertson, and Praveen Paritosh. Characterizing online discussion using coarse discourse sequences. In *Proceedings of the International AAAI Conference on Web and Social Media*, volume 11, pages 357–366, 2017.
- [131] Liane Longpre, Esin Durmus, and Claire Cardie. Persuasion of the undecided: Language vs. the listener. In *Proceedings of the 6th Workshop on Argument Mining*, 2019.
- [132] Jack W Brehm. A theory of psychological reactance. 1966.
- [133] Michael Lynn. Scarcity effects on value: A quantitative review of the commodity theory literature. *Psychology & Marketing*, 8(1):43–57, 1991.
- [134] Alexander J Rothman, Steven C Martino, Brian T Bedell, Jerusha B Detweiler, and Peter Salovey. The systematic influence of gain-and loss-framed messages on interest in and use of different types of health behavior. *Personality and Social Psychology Bulletin*, 25(11):1355–1369, 1999.
- [135] Amos Tversky and Daniel Kahneman. The framing of decisions and the psychology of choice. In *Behavioral decision making*, pages 25–41. Springer, 1985.
- [136] Susan Fournier. Consumers and their brands: Developing relationship theory in consumer research. *Journal of consumer research*, 24(4):343–373, 1998.
- [137] Nicholas Epley, Adam Waytz, and John T Cacioppo. On seeing human: a three-factor theory of anthropomorphism. *Psychological review*, 114(4):864, 2007.
- [138] Sally Hibbert, Andrew Smith, Andrea Davies, and Fiona Ireland. Guilt appeals: Persuasion knowledge and charitable giving. *Psychology & Marketing*, 24(8):723–742, 2007.
- [139] Richard E Petty, John T Cacioppo, and David Schumann. Central and peripheral routes to advertising effectiveness: The moderating role of involvement. *Journal of consumer research*, 10(2):135–146, 1983.
- [140] He He, Derek Chen, Anusha Balakrishnan, and Percy Liang. Decoupling strategy and generation in negotiation dialogues. *arXiv preprint arXiv:1808.09637*, 2018.

- [141] Jacob Cohen. A coefficient of agreement for nominal scales. *Educational and psychological measurement*, 20(1):37–46, 1960.
- [142] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.
- [143] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C Berg. Ssd: Single shot multibox detector. In *ECCV*, pages 21–37. Springer, 2016.
- [144] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *ECCV*, pages 740–755. Springer, 2014.
- [145] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 779–788, 2016.
- [146] Linjie Yang, Kevin Tang, Jianchao Yang, and Li-Jia Li. Dense captioning with joint inference and visual context. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2193–2202, 2017.
- [147] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *ICML*, 2022.
- [148] Stella Frank, Emanuele Bugliarello, and Desmond Elliott. Vision-and-language or vision-for-language? on cross-modal influence in multimodal transformers. *arXiv preprint arXiv:2109.04448*, 2021.
- [149] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, pages 2980–2988, 2017.
- [150] Ran Gilad-Bachrach, Amir Navot, and Naftali Tishby. Query by committee made real. *Advances in neural information processing systems*, 18, 2005.
- [151] David D Lewis and Jason Catlett. Heterogeneous uncertainty sampling for supervised learning. In *Machine learning proceedings 1994*, pages 148–156. Elsevier, 1994.

- [152] Junnan Li, Yongkang Wong, Qi Zhao, and Mohan S. Kankanhalli. Video storytelling: Textual summaries for events. *IEEE Transactions on Multimedia*, 22(2):554–565, feb 2020.
- [153] ROBERT PLUTCHIK. Chapter 1 - a general psychoevolutionary theory of emotion. In Robert Plutchik and Henry Kellerman, editors, *Theories of Emotion*, pages 3–33. Academic Press, 1980.
- [154] Chao-Yuan Wu and Philipp Krahenbuhl. Towards long-form video understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1884–1894, 2021.
- [155] Ali Diba, Mohsen Fayyaz, Vivek Sharma, Manohar Paluri, Jürgen Gall, Rainer Stiefelhagen, and Luc Van Gool. Large scale holistic video understanding. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part V 16*, pages 593–610. Springer, 2020.
- [156] Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Alex Castro-Ros, Marie Pellat, Kevin Robinson, Dasha Valter, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Zhao, Yanping Huang, Andrew Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. Scaling instruction-finetuned language models, 2022.
- [157] Denny Vrandečić and Markus Krötzsch. Wikidata: A free collaborative knowledgebase. *Commun. ACM*, 57(10):78–85, sep 2014.
- [158] Breakthrough. Pyscenedetect: Video scene cut detection and analysis tool, 2023.
- [159] Haofei Xu, Jing Zhang, Jianfei Cai, Hamid Rezatofighi, and Dacheng Tao. Gmflow: Learning optical flow via global matching, 2022.
- [160] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models, 2023.
- [161] Jing Wang, Jinhui Tang, Mingkun Yang, Xiang Bai, and Jiebo Luo. Improving ocr-based image captioning by incorporating geometrical relationship. In *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1306–1315, 2021.

- [162] Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality, March 2023.
- [163] Zhan Tong, Yibing Song, Jue Wang, and Limin Wang. Videomae: Masked autoencoders are data-efficient learners for self-supervised video pre-training, 2022.
- [164] Yi Wang, Kunchang Li, Yizhuo Li, Yinan He, Bingkun Huang, Zhiyu Zhao, Hongjie Zhang, Jilan Xu, Yi Liu, Zun Wang, Sen Xing, Guo Chen, Junting Pan, Jiashuo Yu, Yali Wang, Limin Wang, and Yu Qiao. Internvideo: General video foundation models via generative and discriminative learning. *arXiv preprint arXiv:2212.03191*, 2022.
- [165] KunChang Li, Yinan He, Yi Wang, Yizhuo Li, Wenhai Wang, Ping Luo, Yali Wang, Limin Wang, and Yu Qiao. Videochat: Chat-centric video understanding. *arXiv preprint arXiv:2305.06355*, 2023.
- [166] Youwei Lu and Xiaoyu Wu. Video storytelling based on gated video memorability filtering. *Electronics Letters*, 58(15):576–578, 2022.
- [167] Chen Sun, Austin Myers, Carl Vondrick, Kevin Murphy, and Cordelia Schmid. Videobert: A joint model for video and language representation learning. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 7464–7473, 2019.
- [168] Fanyi Xiao, Kaustav Kundu, Joseph Tighe, and Davide Modolo. Hierarchical self-supervised representation learning for movie understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9727–9736, 2022.
- [169] Rui Qian, Tianjian Meng, Boqing Gong, Ming-Hsuan Yang, Huisheng Wang, Serge Belongie, and Yin Cui. Spatiotemporal contrastive video representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6964–6974, 2021.