

# Large Language Models Are Natural Video Popularity Predictors

Anonymous ACL submission

## Abstract

Predicting video popularity is often framed as a supervised learning task, relying heavily on meta-information and aggregated engagement data. However, video popularity is shaped by complex cultural and social factors that such approaches often overlook. We argue that Large Language Models (LLMs), with their deep contextual awareness, can better capture these nuances. To bridge the gap between pixel-based video data and token-based LLMs, we convert frame-level visuals into sequential text representations using Vision-Language Models. This enables LLMs to process multimodal content—titles, frame-based descriptions, and captions—capturing both engagement intensity (view count) and geographic spread (number of countries where a video trends). On 13,639 popular videos, a supervised neural network using content embeddings achieves 80% accuracy, while our LLM-based approach reaches 82% without fine-tuning. Combining the neural network’s predictions with the LLM further improves accuracy to 85.5%. Moreover, the LLM generates interpretable, theory-grounded explanations for its predictions. Manual validations confirm the quality of these hypotheses and address concerns about hallucinations in the video-to-text conversion process. Overall, our findings suggest that LLMs, equipped with text-based multimodal representations, offer a powerful, interpretable, and data-efficient solution for tasks requiring rich contextual insight, such as video popularity prediction.

## 1 Introduction

Video consumption now accounts for the majority of internet traffic and continues to grow rapidly, making video popularity prediction an important task for content creators, social media platforms, and advertisers (Cisco, 2021). Beyond its commercial importance, accurately identifying which videos will become popular offers insights for researchers studying information diffusion (Park

et al., 2017; Rezvanian et al., 2023), social influence (Park et al., 2016; Lin et al., 2023b), cultural dynamics (Park et al., 2017; Haldar et al., 2023), and misuse in online networks (Beni et al., 2023). Despite its significance, video popularity prediction remains challenging due to factors such as historical context (Ng and Taneja, 2023), cultural trends (Park et al., 2017), and emotional engagement (Guadagno et al., 2013; Park et al., 2016), compounded by the enormous diversity of online video content.

Despite the growing interest in this area, most research has used statistical or supervised learning approaches centered on meta-information and aggregated engagement metrics (e.g., uploader reputation, view/comment/like counts, external social network size) (Zhou et al., 2010; Shamma et al., 2011; Borghol et al., 2012; Park et al., 2016). While these signals are valuable, they generally reflect early user reactions rather than the deeper contextual and cultural nuances of the video content. These intrinsic qualities of video content may play a critical role in how a video resonates with both local and global audiences. However, traditional approaches struggle to process and leverage such rich information due to limited capacity in processing complex multimodal data.

In this paper, we propose a novel approach to video popularity prediction that shifts the focus to the intrinsic qualities of a video’s textual, verbal, and visual content, excluding after-the-fact user engagement data such as early view counts and social network signals. Our method leverages the power of Vision-Language Models (VLMs) and Large Language Models (LLMs) to extract and interpret these intrinsic qualities, complemented by conventional descriptors such as titles and descriptions. More specifically, to address the challenge of integrating pixel-based video data with the token-based architecture of LLMs, we use VLMs to transform frame-level visual data into sequential

textual representations. These representations are then combined with conventional video descriptors such as titles, descriptions, and captions (extracted from the video’s audio) to create a comprehensive multimodal textual representation. This transformation allows LLMs to perform the video popularity prediction task, effectively capturing both vertical (e.g., view counts) and horizontal (e.g., geographic spread) aspects through its contextual understanding of textual content.

Empirically, we introduce a novel prompting strategy that integrates supervised learning signals into LLM-based predictions, outperforming deep learning models based on content embeddings and yielding interpretable, attribute-based explanations. Additionally, we present the *Global Popular Video Dataset (GPVD)*, a large-scale dataset of 1.3M unique popular YouTube videos, enriched with titles, descriptions, and detailed metadata. This dataset uniquely includes three key popularity metrics—view counts, number of countries where a video trended, and number of days spent in trending—spanning 109 countries. The latter two metrics, which reflect the video’s global reach and sustained popularity, have been overlooked in prior research but are crucial for a comprehensive understanding of video virality. Our experiments use a balanced subset covering different popularity classes, varying engagement intensity and global reach, as detailed in the Methods section. Furthermore, we address concerns related to hallucinations in video-to-text conversion and validate the quality of attribute-based explanations through survey experiments.

The key contributions of this paper are as follows:

- We formulate a video popularity prediction task that considers both engagement intensity (e.g., view counts) and geographic spread (e.g., the number of countries where the video trends), addressing the limitation of prior research that primarily focuses on one-dimensional metrics. This nuanced formulation can enable a more comprehensive understanding of video success.
- We introduce the *GPVD*, a large-scale dataset of 1.3M representative YouTube videos from 109 countries, supplemented with metadata such as titles, descriptions, and detailed popularity metrics. This dataset can facilitate research on multimodal video analysis and

cross-cultural trends, bridging gaps in existing datasets.

- We develop an LLM-based prediction framework leveraging a VLM-LLM pipeline that transforms multimodal video content into textual representations, addressing modality gaps and enabling unified reasoning over visual and textual data.
- We systematically explore and evaluate various prompting strategies, including hypothesis generation, KNN-based example retrieval, and the innovative integration of supervised signals. These experiments provide insights into the relative performance of different strategies and their impact on prediction accuracy and interpretability.
- We conduct rigorous human evaluations to validate the accuracy of visual-to-text transformations and assess the quality and interpretability of the generated explanations. These evaluations confirm the robustness of the pipeline and the reliability of the predictions.

## 2 Related Work

**Video Popularity Prediction** Video popularity prediction has traditionally been approached as a supervised learning task, often aimed at estimating view counts using features such as title length, run-time, and early user engagement metrics (e.g., comments and likes) (Zhou et al., 2010; West, 2011; Borghol et al., 2012; Wang et al., 2012). Later work introduced time-series analysis and user behavior modeling to track changes in popularity over time (Broxton et al., 2013; Pinto et al., 2013; Vallet et al., 2015; Park et al., 2016; Jog et al., 2021). However, these methods mostly rely on aggregated platform metrics and struggle to capture the cultural and social dynamics that can significantly influence a video’s success.

Further, the majority of these studies treat popularity as a single-dimensional concept, typically focusing on view counts. Yet a video may accumulate high views in a limited region without achieving broader international reach. For instance, as Park et al. (2017) illustrate that region-specific cultural preferences can inflate view counts locally without translating into global appeal. To address this gap, we expand the notion of popularity to include *geographic spread* (global vs. local reach) alongside *engagement intensity* (view counts). We argue

that Large Language Models (LLMs), equipped with extensive contextual knowledge, are uniquely poised to capture these subtleties in a way that earlier approaches—tied to simpler, often unimodal features—could not.

**Multimodal Learning and Modality Gaps** Incorporating multimodal data, particularly for video understanding tasks, has long been challenging due to the modality gap between pixel-based imagery and token-based language models. Existing systems often adopt modular pipelines that use pre-trained visual encoders (e.g., ViT) for images and language models (e.g., BERT) for text, then concatenate their embeddings (Zeng et al., 2022; Alayrac et al., 2022; Li et al., 2023a,b; Lu et al., 2022). Although these methods capture some visual-linguistic interactions, they may overlook crucial temporal and contextual information (Chen et al., 2023b; Qin et al., 2023).

Recent advancements have leveraged VLMs to transform video frames into textual representations, enabling LLMs to reason over multimodal content (Bhattacharyya et al., 2023; Khandelwal et al., 2024; Chen et al., 2023a). These approaches utilize pre-trained models and modular pipelines to generate textual summaries of videos, which are then employed for tasks such as classification and user behavior modeling. Collectively, they demonstrate the effectiveness of integrating VLMs with LLMs for video understanding tasks, emphasizing the role of textual intermediaries to bridge the modality gap.

Building on these innovations, our approach introduces a frame-to-text transformation pipeline that converts video frames into sequential textual descriptions via VLMs. This enables LLMs to process visual data as richly contextualized text, facilitating unified reasoning across modalities. Compared to modular architectures that merely concatenate embeddings, our method achieves deeper integration of multimodal data, capturing nuanced information from visual and textual content. Moreover, our work extends this pipeline to an end prediction task using systematic prompt engineering, enabling comparisons across prompting techniques and offering insights into performance improvements.

**Natural Language Explanation and Prompt Engineering** Generating explanations alongside predictions has been shown to enhance both model understanding and performance in complex tasks.

Techniques like Chain-of-Thought prompting (Wei et al., 2022) and self-consistency sampling (Wang et al., 2022) have demonstrated how reasoning chains can improve accuracy while maintaining interpretability. Two-stage approaches, such as hypothesis generation followed by task solving (Wang et al., 2023), highlight the value of explanations in boosting performance, though they often add complexity to the prediction pipeline.

Building on Hanu et al. (2023), who demonstrated the effectiveness of textual descriptions for multimodal classification, we integrate hypothesis generation directly into the prediction process, streamlining the two-stage approach into an efficient one-step process. This approach not only reduces computational overhead but also allows the model to generate explanations alongside predictions, enhancing both usability and interpretability. Additionally, by incorporating supervised learning outputs as additional signals within the prompt, we further improve performance, enabling the model to combine external insights with its internal reasoning. This integration bridges the gap between data-driven learning and explainable AI, providing a unified framework for interpretable and high-performing predictions.

### 3 Methods

Our method predicts video popularity by converting key frames into textual descriptions and integrating them with video metadata for LLM-based inference. Figure 1 illustrates two main stages, following the task definition and dataset description:

1. **Transforming Video Content (§3.2):** Generate textual summaries from sampled frames using *VideoLLava* (Lin et al., 2023a) and align them with time-synced captions, titles, and descriptions.
2. **Popularity Prediction (§3.3):** Use structured prompting in an LLM, combining context, reasoning, and transfer sets with supervised signals, few-shot, and near-example guidance.

The subsections below detail the task formulation, data collection, video-to-text transformations, and how we predict local-hit vs. global-big-hit categories.

#### 3.1 Task Definition and Dataset

We redefine the video popularity prediction task by introducing two dimensions:

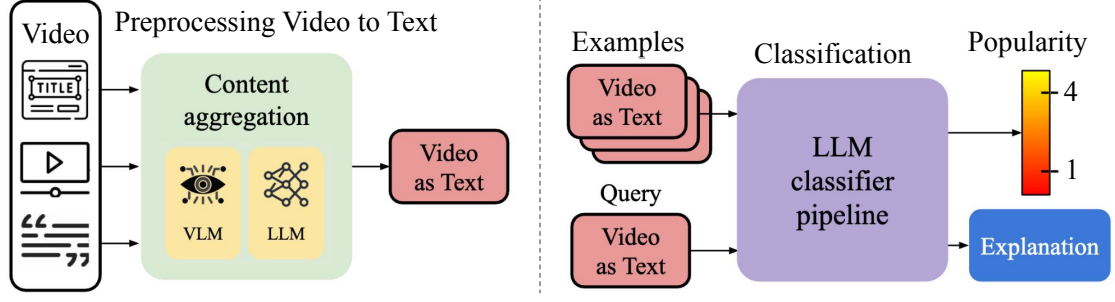


Figure 1: A training-free framework for video popularity prediction utilizing modality-aligned VLMs and LLMs. The **left** shows the video preprocessing and content aggregation stages, where video content is transformed into sequential text representations through VLMs. This transformation generates *Video as Text* summaries, combining visual and textual information. The **right** illustrates the classification and prediction stages, where the LLM processes the *Video as Text* summaries to predict a popularity score and provides an explanation based on the identified patterns.

- **Engagement Intensity:** The total view count, reflecting audience size and engagement.
- **Geographic Spread:** The number of countries where the video appears on trending lists, capturing global reach.

Formally, given a video  $v$  with multimodal features (frames, audio, captions), the goal is to predict  $p(v) \in [0, 1]$  indicating the likelihood of being classified as either a ‘local hit’ or a ‘global big hit.’

To support this, we introduce the *Global Popular Video Dataset (GPVD)*, which includes the top 50 daily trending YouTube videos across 109 countries over 589 days (February 13, 2021, to March 17, 2023). The dataset contains approximately 5,450 observations per day, totaling 3,210,050 observations (1,302,698 unique videos), each with video IDs, trending countries, metadata, and popularity metrics (e.g., views, likes, dislikes).<sup>1</sup> Given that the majority of videos neither go viral nor achieve significant consumption levels, gathering a representative sample of globally popular videos is inherently challenging and valuable. Unlike prior studies with shorter collection periods (1-5 months) (Park et al., 2017; Ng and Taneja, 2023), our dataset spans over two years with broader geographic coverage. This enables both short-term trend analysis and long-term cross-cultural comparisons.

Videos are classified into 16 groups via  $4 \times 4$  quantiles of Engagement Intensity and Geographic Spread. Our experiments focus on two extreme classes, namely ‘global big hit’ (top 25% in both

dimensions) and ‘local hit’ (bottom 25% in both dimensions). We randomly select a balanced sample of 6,279 local-hit and 7,360 global-big-hit videos to evaluate our framework.

### 3.2 Transforming Video Content into Text

Our pipeline converts raw video content into coherent textual narratives:

1. **Frame Extraction:** Sample frames from the video at uniform intervals (10 frames per minute),  $\mathbf{V} = \{f_1, f_2, \dots, f_M\}$ , to capture visual highlights.
2. **Frames to Text Conversion:** Apply *VideoLLava* (Lin et al., 2023a) to each frame (or frame window), generating descriptive text  $S_{\text{VideoLLava}}(\mathbf{V}) = \{S_1, S_2, \dots, S_S\}$ , where each  $S_i$  is a textual description generated from a window of frames around  $f_i$  ( $S_i = \text{VideoLLava}(f_i^W)$ ).
3. **Caption Alignment:** Combine the generated frame-level summaries with time-synced captions,  $\mathbf{C} = \{C_1, C_2, \dots, C_C\}$ , aligning visual and audio cues for richer context.
4. **Summarization:** An LLM ( $\Phi_{\text{LLM}}$ ) ingests the combined text  $\mathcal{I}$  alongside  $T$  (title) and  $D$  (description), producing a concise final summary  $\mathcal{F}$ . This step ensures the text is polished, removing redundancies while maintaining narrative consistency.

This process consolidates disparate video elements (visual, textual, and spoken) into a cohesive textual narrative, facilitating effective LLM-based

<sup>1</sup>We will release the dataset, including video IDs, metadata, and Python code for video downloading, in a publicly accessible repository upon the paper’s acceptance, ensuring reproducibility and support for future research.



inference. Below is an excerpt of the prompt used for summarization:

```
You're an expert in YouTube videos with extensive
experience in analyzing video content and trends.
...
<output>
<scratchpad>
Step 1: Identify key elements -... Step 2: Analyze
Segment Transitions - ...
...
</scratchpad>
<complete_summary>
<segment_summaries>
<introduction> Provide an overview of the video's
theme and initial setting. </introduction>
...
</segment_summaries>
</output>
```

### 3.3 Popularity Prediction through Prompting

We employ a structured prompting strategy in three categories, progressively introducing reasoning steps, few-shot learning, near examples, hypothesis generation, and supervised signals:

1. **Context Set:** Establishes the task objective by specifying instructions, task definitions, and output format.
2. **Reasoning Set:** Encourages intermediate reasoning through prompts like “think before evaluating” and “hypothesis generation,” enhancing the model’s ability to process complex information and provide explanations for its prediction.
3. **Transfer Set:** Refines predictions using external knowledge via techniques such as few-shot learning and near-example selection. Also incorporates signals from a supervised learning output.

#### 3.3.1 Sequential Prompting Approach

We frame the task as a four-class classification ( $c \in \{1, 2, 3, 4\}$ ) representing increasing popularity, then aggregate classes  $\{1, 2\}$  (local hits) vs.  $\{3, 4\}$  (global big hits). This yields more robust boundaries and calibration than a strict binary setup.<sup>2</sup> The input consists of integrated summaries,

<sup>2</sup>During experiments, we observed that the LLM set a very high threshold for the Global ‘Big’ Hit class, resulting in a noticeable bias towards the Local Hit class. Introducing buffer classes (1 and 2 for local hits and 3 and 4 for global hits) addressed this issue by (1) allowing the model to express uncertainty through intermediate predictions, avoiding forced binary decisions; (2) creating a more granular classification system that better reflects real-world ambiguities between extreme categories; (3) enabling better-calibrated confidence levels by incorporating buffer zones between classes; and (4) establishing a smoother decision boundary between extremes, reducing the potential for overly rigid classifications. This

titles, and descriptions, denoted as  $\mathcal{F}(\mathcal{I}, T, D)$ . Prompts were added sequentially to refine performance, as described below:

**Vanilla LLM Prompt (Context Set)** A baseline prompt  $\mathcal{P}_{\text{vanilla}}$  included instructions, the task definition, and the required output format:

$$\mathcal{P}_{\text{vanilla}} = \mathcal{P}_{\text{instructions}} + \mathcal{P}_{\text{task}} + \mathcal{P}_{\text{output}}$$

**Thinking (Context + Reasoning Set)** We added  $\mathcal{P}_{\text{think}}$ , a Chain-of-Thought (Wei et al., 2022) component that encourages explicit intermediate reasoning, improving the LLM’s ability to interpret information within video content.

**Few-shot Learning (Context + Reasoning + Transfer Set)** Few-shot learning (Brown et al., 2020) was then introduced by providing labeled examples  $\mathcal{E} = \{(T_i, \mathcal{F}(\mathcal{I}_i, T_i, D_i), y_i)\}_{i=1}^N$ , where  $y_i \in \{1, 2, 3, 4\}$  denotes the popularity class. This enabled the model to generalize from analogous examples.

**Near Examples (Context + Reasoning + Transfer Set)** We then incorporated semantically similar  $k$  examples,  $\mathcal{E}_{\text{near}}$ , selected based on cosine similarity between title embeddings generated using the MPNet encoder (Song et al., 2020). These examples,  $\mathcal{E}_{\text{near}} \subseteq \mathcal{E}$ , were added to the prompt:  $\mathcal{P}_{\text{near}} = \mathcal{P}_{\text{vanilla}} + \sum_{(T_i, \mathcal{F}(\mathcal{I}_i, T_i, D_i), y_i) \in \mathcal{E}_{\text{near}}} \dots$  where  $T_i$ ,  $\mathcal{F}(\mathcal{I}_i, T_i, D_i)$ , and  $y_i$  represent title, full integrated description, and popularity, respectively. This provided the LLM with relevant context to enhance predictions.

**Hypothesis Generation (Full Context + Reasoning + Transfer Set)** Hypothesis generation prompted the LLM to create a set of hypotheses  $\mathcal{H} = \{h_j\}_{j=1}^M$  based on  $\mathcal{E}_{\text{near}}$ , using a hypothesis generation function  $\Phi_{\text{hypothesis}}$ , a prompt designed to produce hypotheses (see Figure A1 for the final prompt). While Wang et al. (2023) adopts a two-stage approach—first generating hypotheses and then solving the task—for inductive reasoning tasks, we streamline the process by integrating hypothesis generation directly into the ‘thinking’ process. This adjustment enables the LLM to process and synthesize information more effectively, leading to more accurate and interpretable predictions by embedding reasoning within the task-solving step.

adjustment significantly enhanced prediction performance, leading us to adopt the four-class structure in our experimentation pipeline. Details of the class setup can be found in Section A.4 of the Appendix, which includes the final prompt used in our experiments.

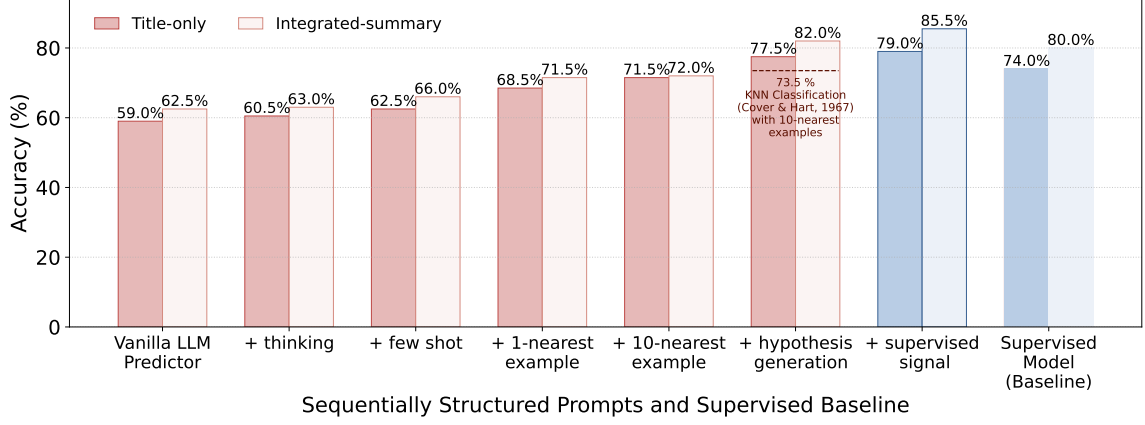


Figure 2: Comparison of accuracy for video-title-only and full-integrated-description-based predictions across models using different prompt sets. The plot shows performance improvements with each model enhancement, beginning with the baseline ‘Vanilla’ model and culminating in the final configuration incorporating supervised signals. The KNN (Cover and Hart, 1967) and supervised models are included as baselines.

**Supervised Signal (Final Prompt)** The final enhancement incorporated supervised signals from a baseline classifier  $\mathcal{F}_{\text{classifier}}$ , appending information like: “A supervised model ( $x\%$  accurate) predicts a popularity rating of {prediction} with {confidence}.” This signal, though noisy, provided the LLM with an external estimate of the video’s potential popularity, thereby encouraging the LLM to weigh external insights as well as to consider the inherent uncertainty in such predictions.

The final prompt,  $\mathcal{P}_{\text{Final}}$ , was constructed through a straightforward concatenation:

$$\mathcal{P}_{\text{Final}} = \mathcal{P}_{\text{vanilla}} + \mathcal{P}_{\text{think}} + \mathcal{P}_{\text{few-shot}} + \mathcal{P}_{\text{near}} + \mathcal{P}_{\text{hypothesis}} + \mathcal{P}_{\text{supervised}}$$

This comprehensive prompt guided the LLM to integrate reasoning, near examples, and supervised signals for informed predictions.

## 4 Experiments and Results

**Implementation Details** We conducted experiments using Claude Sonnet 3.5 (Anthropic AI, 2024) for reasoning tasks and LLaMa 3 (Touvron et al., 2023) for zero-shot and in-context learning. Additional implementation details, including specific configurations for LLaMa, are provided in Appendix A.3. The supervised baseline model was trained with the Adam optimizer (learning rate = 0.001) and early stopping (patience = 6). Additional details are included in Appendix A.8.

**Evaluation Metrics** We used accuracy as the primary metric, supplemented by precision and recall to evaluate classification quality.

**Ablation Study** We conducted extensive ablation studies to evaluate the robustness of our model. Specifically, we analyzed the effects of temperature settings, the number of near-examples, and different embedding type choices on performance. Results in Appendix A.6 reveal the model’s stability and performance across configurations. Further experiments with Gemini 1.5 Pro, a state-of-the-art Vision-Language Large Model (VLLM), demonstrated consistent results, underscoring the generalizability of our framework and prompting strategies (see Appendix A.9).

### 4.1 Supervised Approach (Baseline)

The supervised baseline integrates multimodal embeddings for video features. Textual features like titles ( $T$ ), descriptions ( $D$ ), and captions ( $C$ ) were encoded using the MPNet base v2 encoder (Song et al., 2020), yielding embeddings  $\mathbf{e}_T$ ,  $\mathbf{e}_D$ , and  $\mathbf{e}_C$ . Visual features, including video frames ( $V$ ) and thumbnails, were processed using CLIP (Radford et al., 2021) and its video counterpart (Mendele- vitch and Aguyname, 2023), generating frame- and video-level embeddings ( $\mathbf{e}_{I_i}$ ,  $\mathbf{e}_V$ ). These embeddings were concatenated into a unified representation ( $\mathbf{v}$ ), input to a deep neural network classifier for binary prediction ( $\mathcal{F}_{\text{classifier}} : \mathcal{Z} \rightarrow \{0, 1\}$ ), distinguishing ‘local hits’ from ‘global big hits.’ This model provided a robust benchmark for comparison.

### 4.2 Impact of Sequential Prompts

We evaluated the performance of LLM-based models using both title-only and full-integrated-

description-based predictions. As shown in Figure 2, incremental prompt enhancements led to significant accuracy gains.

Starting with the vanilla prompt ( $\mathcal{P}_{\text{vanilla}}$ ), the model achieved 59.0% accuracy for title-only predictions and 62.5% for integrated descriptions, indicating that the richer information provided in the description resulted in a noticeable improvement in accuracy. Adding thinking prompts ( $\mathcal{P}_{\text{think}}$ ) and few-shot examples ( $\mathcal{P}_{\text{few-shot}}$ ) progressively improved accuracy to 62.5% (title) and 66.0% (integrated-description). Incorporating 1-nearest-example and 10-nearest-example retrieval ( $\mathcal{P}_{\text{near}}$ ) further boosted accuracy to 68.5% and 71.5% (title) and 71.5% and 72.0% (integrated-description), respectively. Hypothesis generation ( $\mathcal{P}_{\text{hypothesis}}$ ) added another 6.0–10.0 percentage points, reaching 77.5% for title-based prediction and 82.0% for integrated-description-based. This underscores the value of hypothesis generation, allowing the model to generate and test multiple hypotheses, in enhancing both accuracy and explainability.

The final model, integrating supervised signals ( $\mathcal{P}_{\text{supervised}}$ ), achieved 79.0% accuracy for title-only predictions and 85.5% for integrated-description-based predictions. This marks improvements of 1.5 and 3.5 percentage points over the previous model and 5.0 and 5.5 percentage points over the baseline supervised model. The results demonstrate that combining LLM reasoning, near examples, hypothesis generation, and supervised signals significantly enhances performance, outperforming traditional baselines and common prompting strategies like thinking and few-shot learning. Additionally, the use of hypothesis generation provides attribute-based explanations, further improving the interpretability and transparency of the predictions.

### 4.3 Manual Validations of Hallucinations and Hypothesis Quality

To evaluate the video-to-text conversion process and the quality of the LLM-generated hypotheses, we conducted two surveys with human evaluators in the US, recruited through Mechanical Turk (MTurk). All participants held at least a Master’s degree and were compensated at an hourly rate equivalent to USD 15 for tasks taking approximately 10–15 minutes each. The survey instructions and questions are fully provided in Section A.2 of the Appendix.

For the evaluation, we selected 5 videos from each popularity category (i.e., 5 local hits and 5

global big hits; 10 videos in total), with each video reviewed by 30 independent evaluators. This setup resulted in a total of 300 evaluations for each task, i.e., assessing (1) video-to-text conversion and (2) hypothesis and analysis quality. Screening questions were implemented at the end of the survey to ensure high-quality feedback. These questions tested attention to video content, focusing on the video’s title, activity, and evaluation metrics. Participants answering all questions correctly were classified as having “passed.” Notably, the results showed consistency across both groups—those who passed and those who did not—demonstrating the robustness and reliability of our pipeline’s outputs.

#### 4.3.1 Validation 1: Video-to-Text Conversion Quality

Participants viewed short clips alongside model-generated text and rated the text on a 1–5 scale for accuracy, adherence, consistency, and alignment with the main topic. Mean scores exceeded 4.22 across all categories (Table 1), indicating that the text was generally well-aligned with the visual content and free of significant factual errors. Notably, even participants who did not pass the screening questions provided average ratings statistically comparable to those who passed, further supporting the robustness and reliability of our summarization approach.

Metric	All ( $N = 30$ )	Passed ( $N = 12$ )
Accuracy	$4.35 \pm 0.30$	$4.32 \pm 0.28$
Adherence	$4.28 \pm 0.40$	$4.22 \pm 0.10$
Consistency	$4.40 \pm 0.25$	$4.36 \pm 0.22$
Main Topic	$4.56 \pm 0.24$	$4.55 \pm 0.30$

Table 1: Survey results evaluating video-to-text conversion quality.

#### 4.3.2 Validation 2: Hypothesis Quality and LLM Analysis

The second survey measured the clarity and relevance of hypotheses explaining each prediction, along with the perceived quality of the LLM’s overall analysis (1–5 scale). Mean ratings exceeded 4.15 across both full samples and the “passed” subsamples (Table 2), demonstrating that annotators generally regarded the model’s explanations as coherent and insightful.

These evaluations confirm that our video-to-text pipeline produces accurate summaries with mini-

Metric	All ( $N = 30$ )	Passed ( $N = 13$ )
Hypothesis	$4.44 \pm 0.26$	$4.45 \pm 0.30$
LLM Analysis	$4.15 \pm 0.25$	$4.24 \pm 0.42$

Table 2: Survey results on the qualities of the generated hypothesis and LLM analysis.

mal hallucination and that the model’s hypotheses offer meaningful, interpretable insights into video popularity. Collectively, they underscore the effectiveness of our approach in both generating reliable content representations and delivering transparent, high-quality predictions.

## 5 Discussion and Future Work

The results of our experiments demonstrate the effectiveness of our approach to video popularity prediction, where LLMs are progressively enhanced through structured prompting techniques. In particular, the performance improvements—from the vanilla prompt to the final model incorporating hypothesis generation and supervised signals—underscore the potential of combining **LLM reasoning capabilities with supervised learning methods**. Notably, **hypothesis generation** not only improved performance as a major contributor to the gains but also enhanced explainability—validated through survey experiments—making predictions more transparent and providing insights into factors driving video popularity.

Our extensive ablation studies, presented in Sections A.6-A.10 of the Appendix, validate the robustness of our framework across various temperature settings, model architectures, and video languages. The results highlight its ability to maintain stable performance under varying conditions, emphasizing its generalizability and reliability even when model parameters fluctuate.

Overall, the results indicate that our approach provides a highly effective solution for multimodal prediction tasks, surpassing traditional supervised models and simpler methods that rely on example-based guidance, such as few-shot and near examples. Additionally, by incorporating two key dimensions of popularity—**engagement intensity** and **geographic spread**—our framework delivers a more nuanced and comprehensive understanding of the factors driving video success, moving beyond the one-dimensional view count focus prevalent in previous research.

Beyond video popularity prediction, the tech-

niques developed in this study have broader implications and applications across various domains. For instance, our VLM-to-LLM pipeline and hypothesis generation methods could generalize to social media analysis, enabling trend, sentiment, or engagement prediction on multimodal platforms. In healthcare, the interpretability of hypotheses generated by LLMs could enhance transparency in medical imaging and diagnosis, where explainable AI is critical for trust and adoption. Similarly, the approach could support education by offering explainable feedback for student assessments or personalized content. Finally, in computational social science, the high-quality hypothesis generation demonstrated by our framework could transform theoretical exploration by offering nuanced explanations, shifting the focus beyond simple statistical coefficients to a richer understanding of sociological and cultural phenomena. These broader applications underscore the versatility and transformative potential of our approach.

Future research could explore several promising directions. Refining the hypothesis generation process, particularly by incorporating advanced reinforcement learning techniques, holds potential for further accuracy improvements. Specifically, the LLM can act as an agent generating hypotheses about video popularity factors, with each hypothesis representing an action within the state space of possible predictions. Prediction accuracy serves as a reward signal, guiding the system to learn which hypotheses are most effective. Additionally, improving the model’s ability to account for cultural and emotional factors could enhance predictions for content like humor or emotionally resonant videos, where traditional metrics (e.g., view counts) may be insufficient. A deeper understanding of these factors could also enable the model to better align with user contexts, gauging video appeal based on situational factors and intent. Further integration of multimodal data, such as audio analysis or granular sentiment analysis of comments, could offer richer insights into the drivers of video popularity. Finally, future work could explore real-time prediction capabilities and adapt this framework to predict trends across platforms beyond YouTube.

## 6 Limitations

Although our framework demonstrates strong performance and generalizability, several limitations merit further attention:



**Reliance on Targeted Dimensions** Our approach effectively captures both *geographic spread* and *engagement intensity*, reflecting how global appeal and influencer-driven novelty can drive popularity (e.g., videos featuring prominent teams or unique gaming challenges). However, it remains unclear whether the LLM consistently relies on these specific dimensions throughout its reasoning process. While hypothesis generation provides a transparent lens into the model’s predictions, future research is needed to confirm whether these explanations truly align with the model’s internal decision-making. For illustrative examples and additional analysis, refer to Appendix A.11.

**Cultural and Emotional Nuances** The model struggles with content heavily influenced by cultural specificity or emotional resonance—such as comedic or highly local videos—where it may overestimate or underestimate global appeal. This suggests that subtle factors like humor, sentiment, or region-specific context are not yet fully captured. Potential improvements include incorporating cultural embeddings or advanced sentiment analysis to handle these nuanced elements more effectively.

**Scope of the Dataset** Although our dataset includes popular YouTube videos at varying levels of success, it focuses on content that has already gained noticeable traction on a single platform. Consequently, generalizing these findings to other platforms or less prominent videos may be limited. Incorporating data from more diverse sources, including niche or region-specific platforms, could provide a fuller picture of the factors driving video success.

**Video-to-Text Conversion** Because our method relies on converting frames and audio into textual summaries, errors in this step can propagate through the pipeline. While manual checks indicate generally high-quality text with minimal hallucinations, ensuring consistent performance across different video genres and formats may still be a challenge. Future work could explore additional multimodal consistency checks or robust alignment techniques to further reduce the risk of transformation errors.

## 7 Ethical Considerations

This study was deemed exempt by the IRB at [Institution Name], as the classification tasks did not

involve human subjects and no personally identifiable information was collected during the manual validation process. Crowdworkers were recruited via Amazon Mechanical Turk and compensated at an average rate of \$15/hour for tasks requiring approximately 10–15 minutes. Participants were provided with a clear description of the study’s purpose before opting in and were free to withdraw at any time.

The dataset used in this study consists of publicly available data collected via the YouTube API and excludes any personally identifiable information. While the dataset primarily features popular videos, reducing the likelihood of harmful or sensitive content, we implemented additional safeguards to further mitigate potential risks. Specifically, during the manual evaluation process, we randomly sampled a larger set of random videos (e.g., 100 videos) to screen for potentially problematic content, such as toxicity or inappropriate material, and found no such instances in the sampled data. This thorough screening ensured that annotators were not exposed to harmful content and confirmed that the dataset does not contain severely harmful material, as anticipated.

To prevent misuse, we plan to share the dataset under controlled access. Researchers will be required to agree to terms of use prohibiting malicious activities, such as targeted harassment or orchestrated misinformation campaigns. These measures reflect our commitment to ethical research practices and responsible data stewardship.

## References

- Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. 2022. Flamingo: a Visual Language Model for Few-Shot Learning. In *Advances in Neural Information Processing Systems*, volume 35, pages 29463–29478.
- Anthropic AI. 2024. The claude 3 model family: Opus, sonnet, haiku. *Claude-3 Model Card*. Accessed: 2024-09-05.
- Nima Ghorbani Beni, Ahlem Bouyer, and Alireza Aghajani Rezaei. 2023. [An improved independent cascade model for influence maximization in social networks](#). *IEEE Transactions on Network Science and Engineering*. Early Access.
- Aanisha Bhattacharyya, Yaman K Singla, Balaji Krishnamurthy, Rajiv Ratn Shah, and Changyou Chen. 2023. [A video is worth 4096 tokens: Verbalize](#)

781	videos to understand them in zero shot. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 9822–9839, Singapore. Association for Computational Linguistics.	
786	Yassine Borghol, Sacha Ardon, Niklas Carlsson, Derek Eager, and Anirban Mahanti. 2012. Content-agnostic factors affecting YouTube video popularity. In <i>Proceedings of the 18th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining</i> , pages 1186–1194.	
792	Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D. Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. <i>Language Models are Few-Shot Learners</i> . <i>arXiv preprint arXiv:2005.14165</i> .	
797	Tim Broxton, Yannet Interian, Jonathan Vaver, and Markus Wattenhofer. 2013. Catching a viral video. In <i>Proceedings of the 2013 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining</i> , pages 296–303.	
802	Guo Chen, Yin-Dong Zheng, Jiahao Wang, Jilan Xu, Yifei Huang, Junting Pan, Yi Wang, Yali Wang, Yu Qiao, Tong Lu, et al. 2023a. Videollm: Modeling video sequence with large language models. <i>arXiv preprint arXiv:2305.13292</i> .	
807	Yiwei Chen, Yunjin Choi, Suk-Jae Hwang, and Jungseock Kim. 2023b. Enhancing social media post popularity prediction with visual content. <i>arXiv preprint arXiv:2405.02367</i> .	
811	Cisco. 2021. Cisco Visual Networking Index: Forecast and Trends, 2017–2022. <a href="https://cloud.report/Resources/Whitepapers/eea79d9b-9fe3-4018-86c6-3d1df813d3b8_white-paper-c11-741490.pdf">https://cloud.report/Resources/Whitepapers/eea79d9b-9fe3-4018-86c6-3d1df813d3b8_white-paper-c11-741490.pdf</a> .	
816	Thomas Cover and Peter Hart. 1967. Nearest neighbor pattern classification. <i>IEEE transactions on information theory</i> , 13(1):21–27.	
819	Rosanna E Guadagno, Daniel M Rempala, Shannon Murphy, and Bradley M Okdie. 2013. What makes a video go viral? an analysis of emotional contagion and internet memes. <i>Computers in human behavior</i> , 29(6):2312–2319.	
824	Rajarshi Haldar, Don Towsley, and Tauhid Zaman. 2023. A temporal cascade model for understanding information spreading dynamics in evolving networks. <i>IEEE Transactions on Network Science and Engineering</i> .	
828	Laura Hanu, Anita L Verő, and James Thewlis. 2023. Language as the medium: Multimodal video classification through text only. <i>arXiv preprint arXiv:2309.10783</i> .	
832	Shivam Jog, Bhargav Siras, Akash Fender, Parikshit Mandurkar, Yash Nikalje, and Sharda Chhabria. 2021. Video popularity prediction using machine learning. <i>International Research Journal of Modernization in Engineering Technology and Science</i> , 3(5).	
	Ashmit Khandelwal, Aditya Agrawal, Aanisha Bhattacharyya, Yaman Kumar, Somesh Singh, Uttaran Bhattacharya, Ishita Dasgupta, Stefano Petrangeli, Rajiv Ratn Shah, Changyou Chen, and Balaji Krishnamurthy. 2024. <i>Large content and behavior models to understand, simulate, and optimize content and behavior</i> . In <i>The Twelfth International Conference on Learning Representations</i> .	837 838 839 840 841 842 843 844
	Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023a. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. <i>arXiv preprint arXiv:2301.12597</i> .	845 846 847 848
	KunChang Li, Yinan He, Yi Wang, Yizhuo Li, Wenhai Wang, Ping Luo, Yali Wang, Limin Wang, and Yu Qiao. 2023b. Videochat: Chat-centric video understanding. <i>arXiv preprint arXiv:2305.06355</i> .	849 850 851 852
	Bin Lin, Bin Zhu, Yang Ye, Munan Ning, Peng Jin, and Li Yuan. 2023a. Video-llava: Learning united visual representation by alignment before projection. <i>arXiv preprint arXiv:2311.10122</i> .	853 854 855 856
	Yiming Lin, Lixin Gao, Yike Guo, and Yingying Zhu. 2023b. A tensor-based independent cascade model for influence maximization in social networks. <i>IEEE Transactions on Knowledge and Data Engineering</i> .	857 858 859 860
	Jiasen Lu, Christopher Clark, Rowan Zellers, Roozbeh Mottaghi, and Aniruddha Kembhavi. 2022. Unified-io: A unified model for vision, language, and multimodal tasks. In <i>The Eleventh International Conference on Learning Representations</i> .	861 862 863 864 865
	Daniel Mendelevitch and Rich Aguinamed. 2023. <i>iejMac/clip-video-encode</i> . github.	866 867
	Yee Man Margaret Ng and Harsh Taneja. 2023. Web use remains highly regional even in the age of global platform monopolies. <i>PloS one</i> , 18(1):e0278594.	868 869 870
	Minsu Park, Mor Naaman, and Jonah Berger. 2016. A data-driven study of view duration on youtube. In <i>Proceedings of the international AAAI conference on web and social media</i> , volume 10, pages 651–654.	871 872 873 874
	Minsu Park, Jaram Park, Young Min Baek, and Michael Macy. 2017. Cultural values and cross-cultural video consumption on youtube. <i>PLoS one</i> , 12(5):e0177865.	875 876 877 878
	Henrique Pinto, Jussara M Almeida, and Marcos A Gonçalves. 2013. A novel time series based approach to predict popularity of online content. <i>Journal of Information and Data Management</i> , 4(3):347.	879 880 881 882
	Jie Qin, Yifan Ding, Yiwei Chen, Jiawei Jiang, and Yong Xu. 2023. Predicting video popularity based on video covers and titles using a multimodal large-scale model and pipeline parallelism. <i>ResearchGate</i> .	883 884 885 886
	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from	887 888 889 890

891	natural language supervision. In <i>International conference on machine learning</i> , pages 8748–8763. PMLR.	
892		
893	Alireza Rezvanian, S Mehdi Vahidipour, and Mohammad Reza Meybodi. 2023. A new stochastic diffusion model for influence maximization in social networks. <i>Scientific Reports</i> , 13(1):6122.	
894		
895		
896		
897	David Shamma, Jude Yew, Lyndon Kennedy, and Elizabeth Churchill. 2011. Viral actions: Predicting video view counts using synchronous sharing behaviors. In <i>Proceedings of the International AAAI Conference on Web and Social Media</i> , volume 5, pages 618–621.	
898		
899		
900		
901		
902	Kaitao Song, Xu Tan, Tao Qin, Jianfeng Lu, and Tie-Yan Liu. 2020. MpNet: Masked and permuted pre-training for language understanding. <i>arXiv preprint arXiv:2004.09297</i> .	
903		
904		
905		
906	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. <i>arXiv preprint arXiv:2302.13971</i> .	
907		
908		
909		
910		
911		
912	David Vallet, Miriam Fernandez, and Pablo Castells. 2015. Social media and information retrieval: A survey. In <i>Proceedings of the 2015 International Conference on Social Media and Society</i> , pages 1–10.	
913		
914		
915		
916		
917	Ruocheng Wang, Eric Zelikman, Gabriel Poesia, Yewen Pu, Nick Haber, and Noah D Goodman. 2023. Hypothesis search: Inductive reasoning with language models. <i>arXiv preprint arXiv:2309.05660</i> .	
918		
919		
920		
921	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2022. Self-consistency improves chain of thought reasoning in language models. <i>arXiv preprint arXiv:2203.11171</i> .	
922		
923		
924		
925		
926	Zhi Wang, Shaojie Ye, Bernardo A Huberman, Chuang Li, and Dongmei Sun. 2012. Characterizing and modeling the dynamics of online popularity. In <i>Proceedings of the 21st international conference on World Wide Web</i> , pages 623–632.	
927		
928		
929		
930		
931	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. <i>Advances in neural information processing systems</i> , 35:24824–24837.	
932		
933		
934		
935		
936	Tyler West. 2011. Going viral: Factors that lead videos to become internet phenomena. <i>Public Relations Review</i> , 37(1):101–104.	
937		
938		
939	Andy Zeng, Maria Attarian, Brian Ichter, Krzysztof Choromanski, Adrian Wong, Stefan Welker, Federico Tombari, Aavek Purohit, Michael Ryoo, Vikas Sindhwani, Johnny Lee, Vincent Vanhoucke, and Pete Florence. 2022. Socratic models: Composing zero-shot multimodal reasoning with language. <i>arXiv</i> .	
940		
941		
942		
943		
944		
	Ruoming Zhou, Sartaj Khemmarat, and Lixin Gao. 2010. The impact of youtube recommendation system on video views. In <i>Proceedings of the 10th ACM SIGCOMM conference on Internet measurement</i> , pages 404–410. ACM.	945 946 947 948 949

**A Appendix**

The appendix provides an in-depth exploration of our video popularity prediction framework, offering detailed analyses and insights to supplement the main text. Its primary objectives are to demonstrate the robustness of our approach, provide a comprehensive understanding of the dataset, and highlight the performance improvements achieved over baseline models. Additionally, we present key ablation studies to examine the impact of various hyperparameters on our model’s performance.

**A.1 Overview of Pipeline**

Figure 1 presents an overview of our approach. This high-level view depicts our training-free framework for video popularity prediction, which leverages modality-aligned Vision-Language Models (VLMs) and Large Language Models (LLMs) to generate *Video as Text* summaries. In the preprocessing stage (left), video content is transformed into sequential text representations using VLMs. During the content aggregation stage, visual and textual information is aligned and combined. The LLM then processes these text summaries to predict a video’s popularity score (ranging from 1 to 4) and generates explanations based on identified patterns (right). These explanations take the form of hypotheses grounded in theoretically sound attributes. For example, if the video is about a national football game organized by FIFA, the model may highlight its global appeal due to the prominence of the organization and the attention drawn by specific teams, such as Brazil.

**A.2 Surveys for Human Evaluation**

To evaluate the video-to-text conversion process and the quality of LLM-generated hypotheses, we conducted two separate surveys. The first focused on assessing video-to-text conversion quality, while the second evaluated the quality of hypotheses and LLM analysis. Each survey included clear instructions, detailed evaluation criteria, and screening questions to ensure participant attention and understanding. Participants were recruited through Amazon Mechanical Turk (MTurk) and compensated at an hourly rate equivalent to USD 15, reflecting fair pay for tasks requiring approximately 10–15 minutes each.

**A.2.1 Survey 1: Video-to-Text Conversion Quality**

**Video-to-Text Conversion Quality: Initial Instructions**

Welcome to our research study on video-to-text transcription quality. We are academic researchers from \*\*\*\*, investigating the accuracy of automated transcription systems.

**Survey Overview**

In this survey, you will:

- Watch short video clips
- Read the corresponding automated transcriptions
- Answer questions about the accuracy and quality of the transcriptions

The survey should take approximately 8-12 minutes to complete. You will receive:

- ☐ Base compensation: USD 1.25
- ☐ Potential bonus: Up to USD 5 total for high-quality responses

**Your Role as an Evaluator**

Describing video content accurately is an incredibly complex task for AI. It requires understanding context, nuance, and implied information—skills that come naturally to humans but are extremely challenging for machines. Your task will involve:

- Watching diverse video clips



- Reviewing AI-generated content descriptions
- Providing detailed feedback on accuracy and quality

### **Key Evaluation Areas**

When assessing the AI-generated descriptions, please consider:

- Overall accuracy in capturing key themes and concepts
- AI's ability to understand context and implied information
- Areas where the AI shows particular understanding
- Opportunities for improvement

### **Important Considerations**

Please note:

- The AI system provides a comprehensive overview, not word-for-word transcription
- Focus on overall meaning and key points rather than exact phrasing
- The AI may make contextual inferences

### **Ready to Begin?**

- ☐ I understand all the above instructions thoroughly

976

## **Video-to-Text Conversion Quality: Evaluation Criteria**

### **1. Overall Accuracy**

How accurately does the text description match the content of the video?

- ☐ 1: Completely inaccurate
- ☐ 2: Mostly inaccurate
- ☐ 3: Somewhat accurate
- ☐ 4: Mostly accurate
- ☐ 5: Highly accurate

### **2. Content Accuracy**

How closely does the text description stick to the content presented in the video?

- ☐ 1: Mostly unrelated to video content
- ☐ 2: Significant deviations from video content
- ☐ 3: Moderate adherence to video content
- ☐ 4: Close adherence to video content
- ☐ 5: Perfectly matches video content

977

### 3. Main Topic Capture

How well does the text description capture the main topic(s) discussed in the video?

- ☐ 1: Misses all main topics
- ☐ 2: Captures few main topics
- ☐ 3: Captures some main topics
- ☐ 4: Captures most main topics
- ☐ 5: Accurately captures all main topics

### 4. Key Point Coverage

To what extent are key points from the video included in the text description?

- ☐ 1: Misses all key points
- ☐ 2: Includes few key points
- ☐ 3: Includes some key points
- ☐ 4: Covers most key points
- ☐ 5: Covers all key points

## Video-to-Text Conversion Quality: Task: Video Transcription Evaluation

Watch this Video, you will be given the task of evaluating a short transcript about this video next:  
Video title: Cristiano Ronaldo Hat-Trick! | Manchester United 3-2 Norwich | Highlights

[Youtube Video]

### Important Note

Your careful attention to this video is essential for accurately understanding the issues related to AI behavior and the quality of AI-generated video transcriptions. The more accurately you understand and remember the video's content, the more accurate and valuable your evaluation will be. We encourage you to watch the entire video attentively, as your insights will directly impact the assessment of AI performance!

Participants who demonstrate a thorough understanding of the video content will be eligible for bonus compensation. Thank you for your dedication to this task!

## Transcript

### Introduction

This video is a recording of a football match between two teams, featuring commentary and analysis throughout. The initial setting is a stadium with players from both teams on the field.

### Segment Details

- **Segment 1:** The first segment shows a man walking on the field while another man walks in the background. The commentator describes the scene, mentioning the ball and a goal scored by Adrian Luna.
- **Segment 2:** In this segment, a man is seen walking towards the camera, wearing a yellow shirt, while another man sits on the ground, wearing a red shirt. The commentators discuss the game, highlighting great deliveries and goals scored.
- **Segment 3:** This segment shows more gameplay, with the commentators analyzing the players' moves and discussing the score.

### Overall

The overall impact of this video is an immersive and engaging experience for football fans. The combination of exciting commentary, intense gameplay, and skilled players creates a thrilling narrative that will likely appeal to viewers who enjoy sports content.

### Evaluation

	1 Not at all	2 Slightly	3 Somewhat	4 Mostly	5 Completely
How accurately does the text description match the content of the video?	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
How closely does the text description stick to the content presented in the video?	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
How consistent is the information in the text description with the facts presented in the video?	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
How well does the text description capture the main topic(s) discussed in the video?	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>

980

### Video-to-Text Conversion Quality: Screening Questions

**Q1.** Which videos did you evaluate in this survey?

- ☐ SUPAHOTFIRE vs BLUEFACE
- ☐ Cristiano Ronaldo Hat-Trick
- ☐ Kerala Blasters FC vs Jamshedpur FC Highlights
- ☐ Where I'm Travelling Next- Solo Trip?

**Q2.** In the videos you evaluated, which of the following activities was *not* mentioned?

- ☐ Playing rock-paper-scissors
- ☐ Eating cake
- ☐ A football match
- ☐ Skydiving

**Q3.** In the evaluation process, what were you asked to rate about the video transcriptions?

- ☐ Overall Accuracy
- ☐ Content Accuracy
- ☐ Main Topic Capture
- ☐ Video Production Quality
- ☐ Key Point Coverage

**Q4.** Thank you for participating in our study. Your insights will help us improve our AI model for predicting video popularity. Do you have any additional comments or feedback about the model's predictions or this survey?

981

A.2.2 Survey 2: Hypothesis Quality and LLM Analysis

Hypothesis Quality and LLM Analysis: Initial Instructions

Welcome

Welcome to our research study on video popularity prediction using AI models. We are academic researchers from ... , investigating how Large Language Models (LLMs) can predict video popularity.

Survey Overview

In this survey, you will:

- Read video descriptions and LLM-generated hypotheses about video popularity
- Rate the accuracy and relevance of these hypotheses
- Assess the LLM’s ability on critical analysis and judgment

The survey should take approximately 8-12 minutes to complete. You will receive:

- ☐ Base compensation: USD 1.25
- ☐ Potential bonus: Up to USD 5 total for high-quality responses

Study Overview

We are evaluating a Language Learning Model (LLM) designed to predict video popularity based on content analysis. The LLM analyzes videos from YouTube’s trending page and generates hypotheses about what makes videos popular, as well as providing a detailed analysis of each video’s content. Your role is to:

- Rate the hypotheses generated by the LLM
- Assess the LLM’s critical analysis and judgment for TWO separate videos

Video Popularity Rating Scale

The LLM rates videos on a 4-point scale:

- **Popular:** Likely to have general appeal and be popular for a short while
- **Moderately Popular:** Has several appealing elements for more than basic popularity
- **Highly Popular:** Likely to be popular among a broad audience but may not reach ultra popularity
- **Ultra Popular:** Strong potential to become ultra popular, featuring unique, engaging, and broadly appealing content

*Note: While the LLM uses this 4-point scale to rate video popularity, your task will be to rate your agreement with the LLM’s hypotheses and analysis using a different 4-point scale.*

Hypothesis Quality and LLM Analysis: Instructions - Part 2

Your Tasks

Task 1: Rate the LLM’s Hypotheses

You will be presented with video descriptions and the LLM’s hypotheses about what makes them popular. You will rate your agreement with 4-5 specific hypotheses generated by the LLM about what makes the video popular. For example, a hypothesis looks like ‘*Sport highlights, especially from important matches, tend to be ultra popular*’, to which you can Strongly agree, Agree, Disagree or Strongly Disagree. Your job is to rate how much you generally agree or disagree with given hypothesis based on the same information provided to the LLM.

Task 2: Assess the LLM’s Critical Analysis

You will evaluate the LLM’s critical analysis of the video, including its assessment of factors influencing popularity and its final popularity prediction.



### Rating Scale for Your Responses

You will use a 4-point scale to rate both the LLM's hypotheses and its critical analysis. This scale is designed to encourage you to form a definitive opinion based on your knowledge and the information provided. For both tasks, use the following scale:

- **1 - Strongly Disagree:** The hypothesis or analysis is clearly incorrect or irrelevant
- **2 - Disagree:** The hypothesis or analysis has major flaws or inaccuracies
- **3 - Agree:** The hypothesis or analysis is mostly accurate and relevant
- **4 - Strongly Agree:** The hypothesis or analysis is highly accurate and insightful

### Tips for Completing the Tasks

- Read each video description and LLM hypothesis carefully before rating
- Consider each hypothesis and analysis point carefully. Draw on your own knowledge of popular online content, but focus primarily on the information provided in the video description
- For instance, if you think the LLM's hypothesis about sports highlights is accurate based on the video description and your knowledge, you might select 'Agree' or 'Strongly Agree'
- Try to be consistent in your ratings across similar types of content

Your thoughtful evaluations will help us improve the LLM's ability to predict video popularity, ultimately contributing to a better understanding of content trends on platforms like YouTube.

985

### Hypothesis Quality and LLM Analysis: Video: Minecraft but there's Cartoon Hearts

#### Video summary

Introduction: The video opens with a mysterious green figure walking around a dark room, setting the tone for a fantastical and humorous adventure. Segment Details - Segment 1: Introduces the green figure, referencing Shrek and showcasing magic powers . . . Conclusion: The video's impact is significant, as it showcases the creators imagination and ability to blend disparate elements into a cohesive narrative. The humor, entertainment value, and references to popular franchises will likely appeal to viewers who enjoy fantasy, scifi, and comedy.

In the next two pages, you will be shown 2 tasks:

- Task 1: Rate the hypotheses generated by the LLM
- Task 2: Assess the LLM's critical analysis and judgment

986

## Task 1: Rate the LLM's Hypotheses

### Video Information

#### Important Note

Your careful attention to this video description is essential for accurately understanding the quality of AI-generated hypothesis. The more accurately you understand the video's content, the more accurate and valuable your evaluation will be. We encourage you to read the entire video description attentively, as your insights will directly impact the assessment of AI performance!

Participants who demonstrate a thorough understanding of the video content will be eligible for bonus compensation. Thank you for your dedication to this task!

### Model's Hypotheses

	Strongly Disagree	Disagree	Agree	Strongly Agree
H1: Videos with unique Minecraft concepts tend to be ultra-popular	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
H2: Content that blends multiple franchises or pop culture elements has broader appeal	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
H3: Videos with humorous and imaginative content encourage sharing and discussion	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
H4: Fast-paced content with diverse visual elements keeps viewers more engaged	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>

If you Disagree/Strongly Disagree, do you have any better hypotheses or suggestions for improving the provided hypotheses?

## Task 2: Assess the LLM's Critical Analysis

### Factors in the given video that could influence popularity:

- F1: Creative concept: "Minecraft but there's Cartoon Hearts" (very positive)
- F2: Blending of multiple franchises (Shrek, Teen Titans, Star Wars, Scooby-Doo) (positive)
- F3: Humorous and playful tone (positive)
- F4: Imaginative scenarios (magic powers, cyber crystals, outer space) (positive)
- F5: Alignment with geek culture trends (positive)
- F6: Potential for viewer engagement and discussion (positive)

**LLM's Final Analysis:** Considering all factors, especially the similarity to other ultra popular videos and the supervised model prediction, this video is likely to be Ultra Popular. It has all the elements of highly engaging content that tends to perform exceptionally well, particularly in the gaming and geek culture niches.

	Strongly Disagree	Disagree	Agree	Strongly Agree
F1: Creative concept	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
F2: Blending of multiple franchises	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
F3: Humorous and playful tone	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
F4: Imaginative scenarios	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
F5: Alignment with geek culture trends	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
F6: Potential for viewer engagement	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>

## Hypothesis Quality and LLM Analysis: Screening Questions

**Q1.** What was the primary task you were asked to perform in this survey?

- ☐ Predict which videos would become viral
- ☐ Assess the LLM's hypotheses and analysis about video popularity
- ☐ Provide your own theories about what makes videos popular
- ☐ Compare different AI models' performance in analyzing videos

These video-specific screening questions were randomly assigned based on the viewed video:

**Q2.** Which of the following elements were present in the video you analyzed? (select all that apply)

- ☐ Commentary and analysis
- ☐ Great deliveries and goals scored
- ☐ Interviews with team managers
- ☐ Slow-motion replays
- ☐ Adrian Luna scoring a goal
- ☐ Penalty shootout

**Q3.** Which of the following elements were present in the video you analyzed? (select all that apply)

- ☐ References to Shrek
- ☐ Mining of cyber crystals
- ☐ Outer space scenes
- ☐ Pokémon battles
- ☐ Teen Titans-inspired content
- ☐ Underwater exploration

**Q4.** Which of the following elements were present in the video you analyzed? (select all that apply)

- ☐ Internal struggle of the protagonist
- ☐ Car chase scenes
- ☐ Self-harm depicted
- ☐ Comedic dialogue
- ☐ Apology and plea for forgiveness
- ☐ Transformation sequences

**Q5.** Which of the following elements were present in the video you analyzed? (select all that apply)

- ☐ Two men playing video games
- ☐ Reaction to a music video
- ☐ Wearing headphones
- ☐ Dancing performances
- ☐ Occasional singing into microphones
- ☐ Cooking demonstrations

### A.3 Implementation Details

The video popularity prediction pipeline was implemented using PyTorch 2.1.0 and the transformers 4.35.0 library. The content aggregation was performed using a custom module that combined textual information from titles, descriptions, and generated captions. We employed Claude 3.5 Sonnet (accessed via Anthropic’s API) (Anthropic AI, 2024) as our primary LLM for classification and hypothesis generation, while also using LLaMa 3 70B (Touvron et al., 2023) Instruct offline on two NVIDIA RTX A4000 GPUs using tensor parallelism and quantization (2.8 bits per weight) for efficient video-to-text generations. The VLM component (of VideoLLava) used a fine-tuned CLIP model to align visual and textual features. The entire pipeline was orchestrated using a custom Python script that handled data flow between components, with batching implemented to optimize throughput. For the offline LLaMa 3 70B setup, we achieved approximately 2 tokens per second inference speed. All experiments maintained consistent hyperparameters (Temperature: 0.5, Max Tokens: 4096, Top-p: 0.95) to ensure reproducibility.

### A.4 Final prompt

Figure A1 illustrates the final prompt structure for video popularity prediction using LLMs. The prompt incorporates step-by-step instructions, contextual elements, and a structured output format designed to guide the LLM in analyzing and predicting video popularity.

Key features of the prompt include:

- Instructions: Clear definitions of the task and the four popularity classes, ranging from “Locally Moderately Popular” to “Globally Ultra Popular,” to ensure the model understands the classification criteria.
- Step-by-Step Reasoning: A scratchpad mechanism encourages the model to reason through intermediate steps, such as comparing the given video to similar examples, considering supervised model predictions, and generating hypotheses to explain observed patterns.
- Structured Output: The prompt specifies a coherent format for outputs, including an evaluation rating and explanatory reasoning, ensuring interpretability and consistency.

### A.5 Dataset Analysis

Our dataset analysis uncovers several key insights that help contextualize the dynamics of video popularity and inform our prediction task design. We highlight the importance of considering both *geographical reach* and *view count* as critical factors in assessing a video’s popularity. Our analysis shows a positive correlation between a video’s international presence and its view count, with notable clusters of videos distributed across different quadrants of this relationship.

Additionally, a three-dimensional analysis that includes the duration a video remains on trending lists reveals a more nuanced relationship between trending duration, geographical reach, and view counts. We observe an optimal range of 100-200 units of trending duration and a reach of 15-35 countries as indicators of peak performance. However, outliers in this analysis suggest that content quality and other unquantified factors play important roles in determining a video’s success.

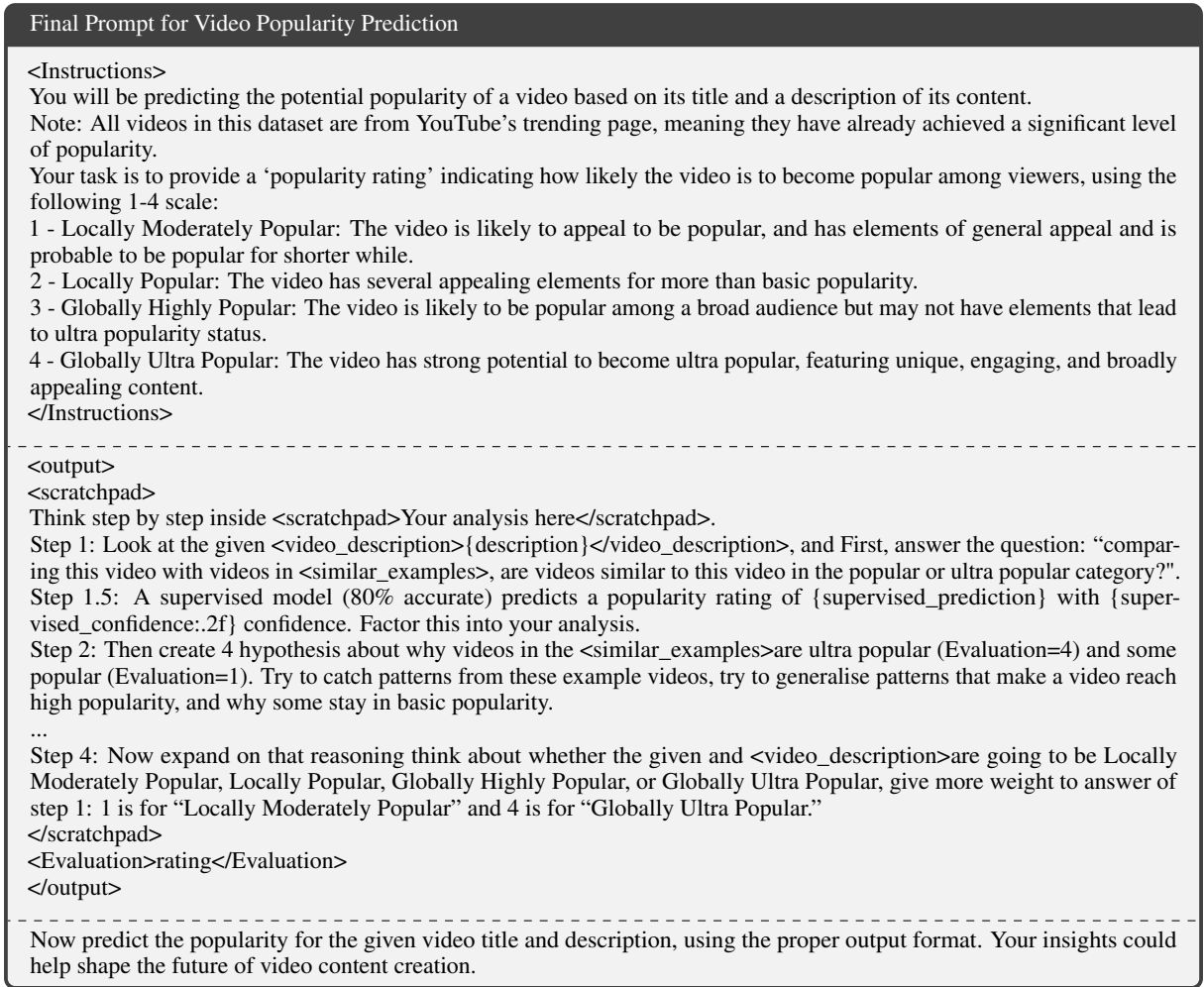
#### A.5.1 Dataset features and video categorization

We present the features used for each video (Figure A2, left) and a heatmap categorizing videos based on engagement intensity and geographical reach (Figure A2, right). This information provides crucial context for understanding the nature of our dataset and how we distinguish globally viral videos from those with more localized popularity.

#### A.5.2 Geographical reach vs. view count

To better understand how a video’s international presence relates to its popularity, we analyzed our dataset, focusing on the connection between the number of countries a video reaches and its total views. Figure A3 visualizes this relationship.





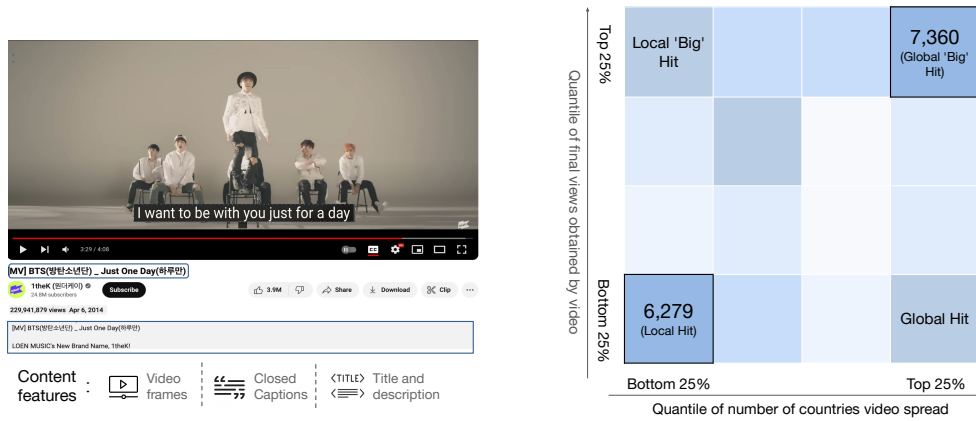


Figure A2: **Left:** The list of features for a youtube video. **Right:** Heatmap categorization of YouTube videos into 16 quantiles based on two key dimensions: the number of views and the number of countries in which the video trended. Videos are classified into ‘Global Big Hit’ (top 25% in both dimensions) and ‘Local Hit’ (bottom 25% in both dimensions), with cell colors indicating the relative density of each class.

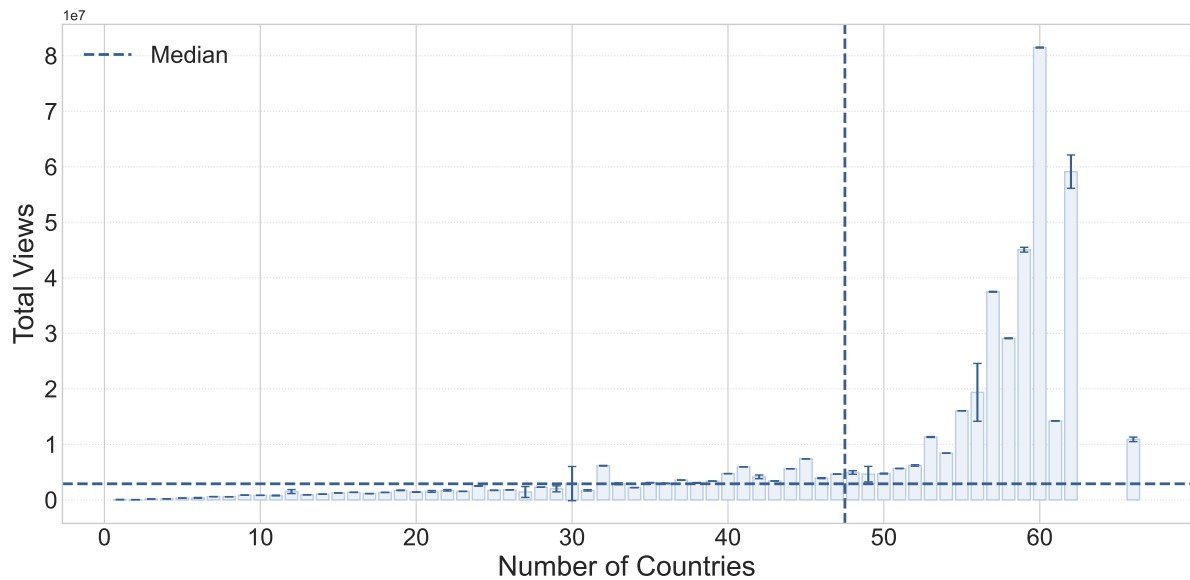


Figure A3: The plot depicts the relationship between the number of countries a video reaches and its total view count. Blue lines represent the median values for each dimension, dividing the plot into quadrants.

increases from 0 to about 100 units. Beyond this point, the relationship becomes more complex, with videos between 100-200 units performing particularly well, especially when they reach a moderate to high number of countries. Interestingly, there’s a slight decline in viewership for extremely long sequences (200+), suggesting an optimal range for sequence length.

The impact of country reach on viewership is evident, with videos reaching more countries generally receiving more views. However, this relationship varies across different sequence lengths. For shorter sequences (0-50), the impact of reaching more countries is less pronounced, while it becomes more significant for medium to long sequences.

Several notable patterns emerge from this analysis:

- A clear region of low viewership is visible in the bottom-left corner, corresponding to short sequences with minimal international reach.
- A ‘hot zone’ appears in the middle-right area of the heatmap, encompassing sequence lengths of 100–175 and country counts of 15–35, where viewership is consistently high.

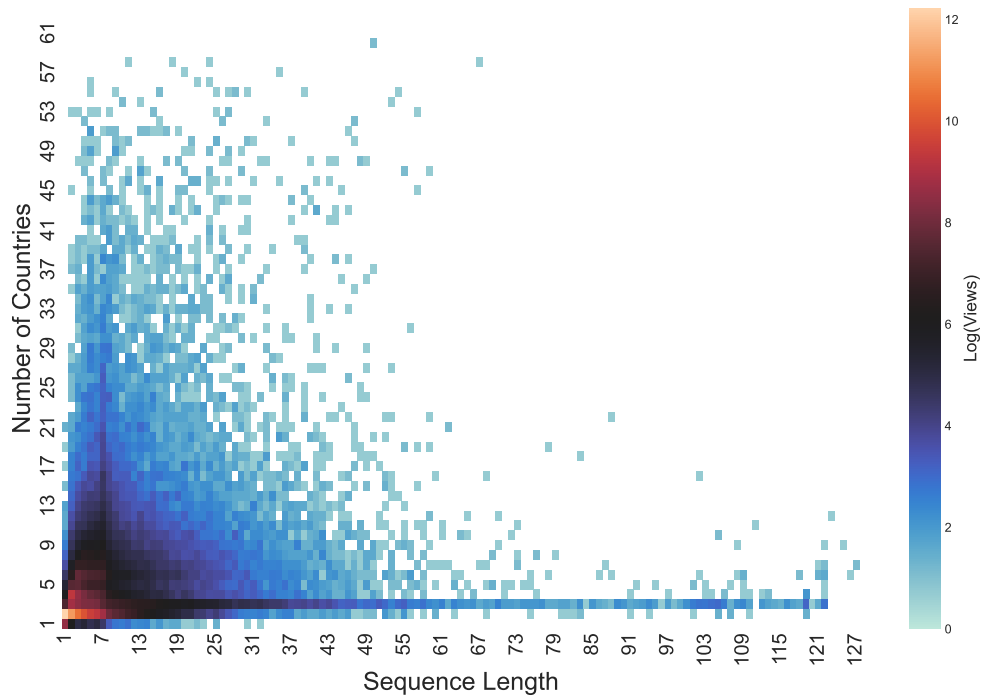


Figure A4: Heatmap depicting the relationship between trending duration (x-axis), number of countries reached (y-axis), and the logarithm of total views (color intensity) (say that image is truncated to 127 sequence length).

- Evidence of potential diminishing returns is observed for sequences exceeding 200 in length and country counts above 40. 1069 1070
- Scattered ‘hot spots’ are present across the heatmap, highlighting outlier videos with unexpectedly high viewership. 1071 1072

#### A.5.4 Video categories 1073

YouTube’s content ecosystem is diverse, encompassing a wide range of video categories. Our dataset provides a unique opportunity to analyze popularity trends across these categories, offering insights that are typically challenging to obtain. In this section, we present an analysis based on 11 heatmaps, each representing a distinct YouTube category (Figure A5). These heat maps visualize the complex interaction between the length of the sequence, the number of countries reached, and the total views for each category. This multi-dimensional analysis reveals both overarching trends and category-specific patterns in video popularity. 1074 1075 1076 1077 1078 1079 1080

Our analysis reveals several consistent patterns across categories. The maximum number of views ranges from 182 million to 1.48 billion views. The average number of views for most categories falls between 63,000 and 251,000 views. Interestingly, for almost all categories, maximum views occur at very short sequence lengths (mostly 1) and low number of countries (2), suggesting that brief, targeted content can achieve high viewership. 1081 1082 1083 1084 1085

Each category exhibits unique characteristics. The Games category demonstrates the highest maximum views (about 1.48 billion views) and one of the highest average views, indicating high engagement. In contrast, the News category contains the largest number of videos but shows a lower average view count, suggesting a high volume of content with more moderate individual performance. The Music category, despite having the fewest videos, maintains a competitive average view count, indicating that music videos tend to perform well relative to their number. The NonCategory exhibits the lowest average views, which might be expected for content that doesn’t fit into standard categories. 1086 1087 1088 1089 1090 1091 1092

The relationship between sequence length, country reach, and views varies across categories. Across most categories, shorter sequence lengths (0-20) tend to have higher average view counts. The Games category shows particularly high performance for short sequences (about 31,000 views on average). 1093 1094 1095

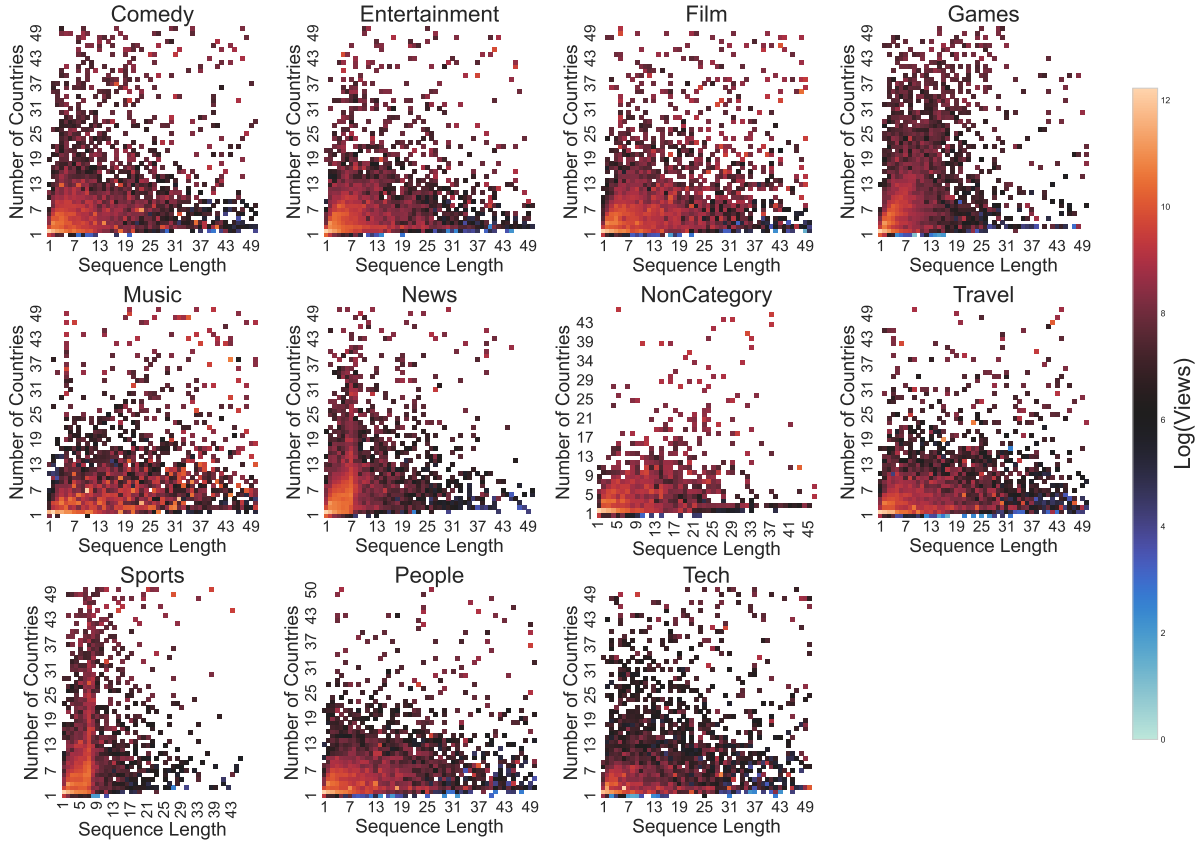


Figure A5: Heatmaps depicting the relationship between sequence length, number of countries reached, and logarithm of total views for 11 YouTube categories. Each heatmap represents a different category, with color intensity indicating  $\log_{10}(\text{Views})$ .

Regarding country reach, videos reaching 6-10 countries often have the highest average views across categories. There's a consistent decline in average views as the number of countries increases beyond 15, suggesting that very broad international appeal is rare.

**Studying Popularity and Reach** We looked at how different video categories do across popularity and international reach using a grid of pie charts (Figure A6). Each pie chart represents a different quantile combination of video popularity (views) and international reach (number of countries), providing a comprehensive view of category distribution across various levels of success.

Our analysis uncovers the universal appeal of certain video categories across different levels of popularity and geographical spread. Notably, "People," "News," and "Sports" stand out, appearing in the majority of quantiles, indicating their widespread popularity. "Comedy" and "Film" also show strong presence, suggesting their content resonates across various levels of success and international reach.

## A.6 Ablation Study

A detailed ablation study was conducted to examine the effects of key model hyperparameters on prediction accuracy. This investigation focused on the interaction between the embedding type used for near example retrieval, the count of these examples, and the temperature parameter of the Large Language Model (LLM). By methodically adjusting these parameters, the study aimed to reveal how variations in these elements influence the framework's predictive capabilities. The analysis closely examines the impact of embedding types, whether sourced from video descriptions or titles, the strategic selection of near example quantities, and the temperature settings within the LLM, providing a thorough examination of their combined effects on performance. This study is crucial as it illuminates the framework's operational nuances and informs potential adjustments, enhancing its effectiveness in the intricate task of video content analysis and prediction. Through this rigorous analysis, we aim to explore the framework's responsiveness to different

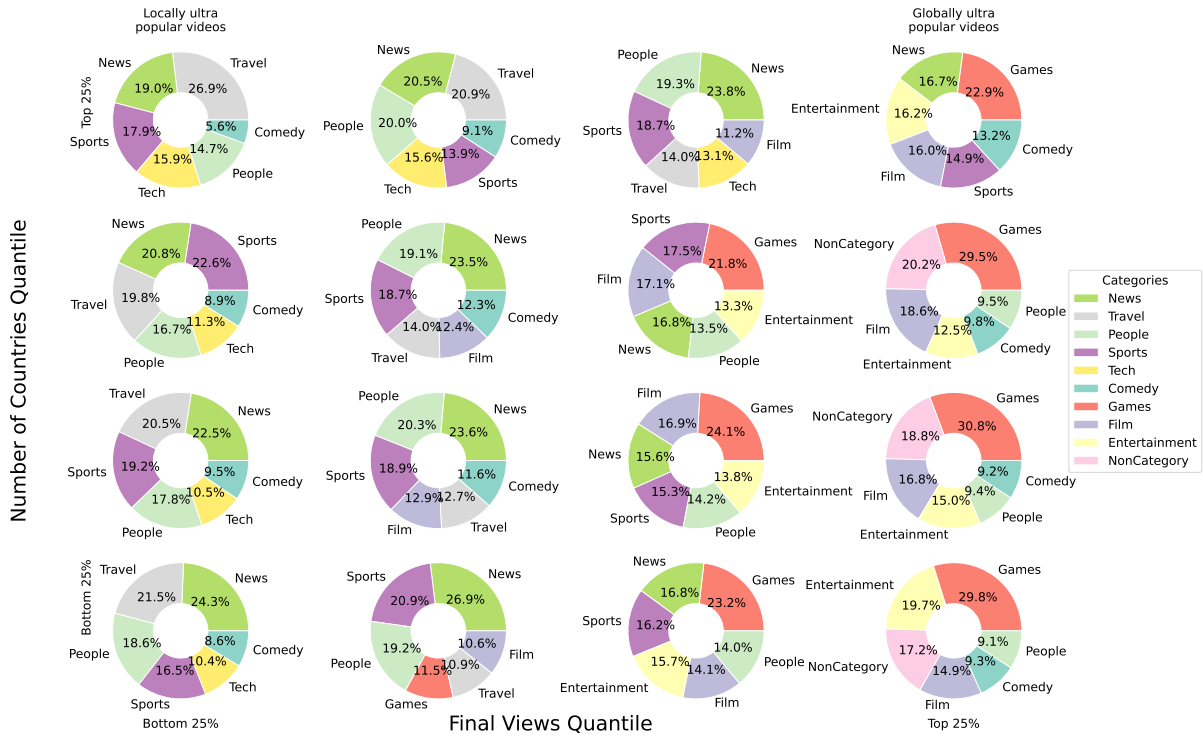


Figure A6: 4x4 grid of pie charts showing the distribution of video categories across different quantiles of popularity (views) and international reach (number of countries).

hyperparameters, contributing to a nuanced understanding of its predictive mechanisms. This endeavor is not about optimizing the model per se but about uncovering how the framework behaves under varied conditions, offering valuable insights into its structure and function.

### A.6.1 Impact of embedding type on model performance

In our study of our proposed framework, we focused on understanding how the choice of embedding type, whether video descriptions or titles, used to find similar videos affects the accuracy of predicting video popularity. It's important to note that while both methods use full video descriptions, the key difference lies in how these similar videos are identified. This study looks into how choosing between video descriptions or titles to find similar videos affects prediction accuracy, showing that using titles to retrieve examples, surprisingly make our pipeline's predictions better.

We report the findings in Figure A7, which show a nuanced result of different kinds of retrieved examples. Video-to-text achieved a mean accuracy of 81.75%, showcasing a consistent prediction capability with a standard deviation of 0.35%. This suggests that descriptions provide a reliable basis for similarity matching, albeit with a marginally lower accuracy compared to titles. Conversely, title embeddings yielded a higher mean accuracy of 85.5%, indicating their effectiveness in accurately identifying highly popular videos, albeit with increased variability, as evidenced by a standard deviation of 1.77%. This discrepancy may stem from the complexity and length of generated video descriptions, which could introduce extraneous information, diluting the core elements necessary for precise similarity matching. Titles, being more succinct, appear to offer a more focused approach for example retrieval, likely due to their ability to encapsulate the video's essence more directly. In contrast, the KNN model exhibited a mean accuracy of 79% with video descriptions and 73.5% with titles, highlighting a different pattern of performance that underscores the importance of model choice in leveraging embedding types effectively. This means that the way we select similar videos is key, and using titles, which are shorter, might predict popularity better by focusing on the main points of the video. The study shows that how we choose similar videos can make a big difference in how well our pipeline works, suggesting that we should think carefully about how we find these videos to improve predictions.



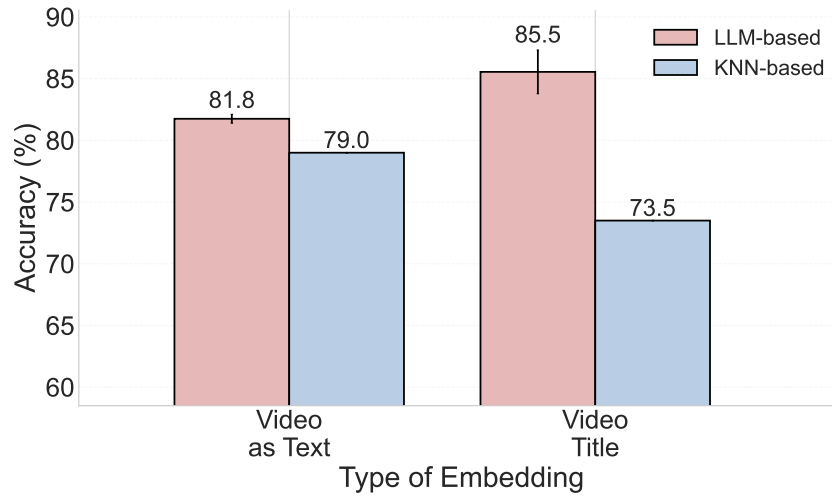


Figure A7: This graph compares how well our pipelining and near-examples models predict video popularity using 10 examples each, showing that using titles to find similar videos works better, even though both methods use examples containing full video descriptions as input. The only difference lies in how we find these examples: using titles or descriptions.

#### A.6.2 Impact of number of near examples on model performance

In this part of our study, we looked at how changing the number of similar videos (near examples) used by the our pipeline affects its ability to predict video popularity. Near examples are similar videos retrieved from the database that the model uses to identify patterns and make predictions. This analysis aims to determine the relationship of near examples, prediction accuracy and computational efficiency. We kept everything else the same and only changed the number of these examples to see how it impacts accuracy and how efficiently the model works.

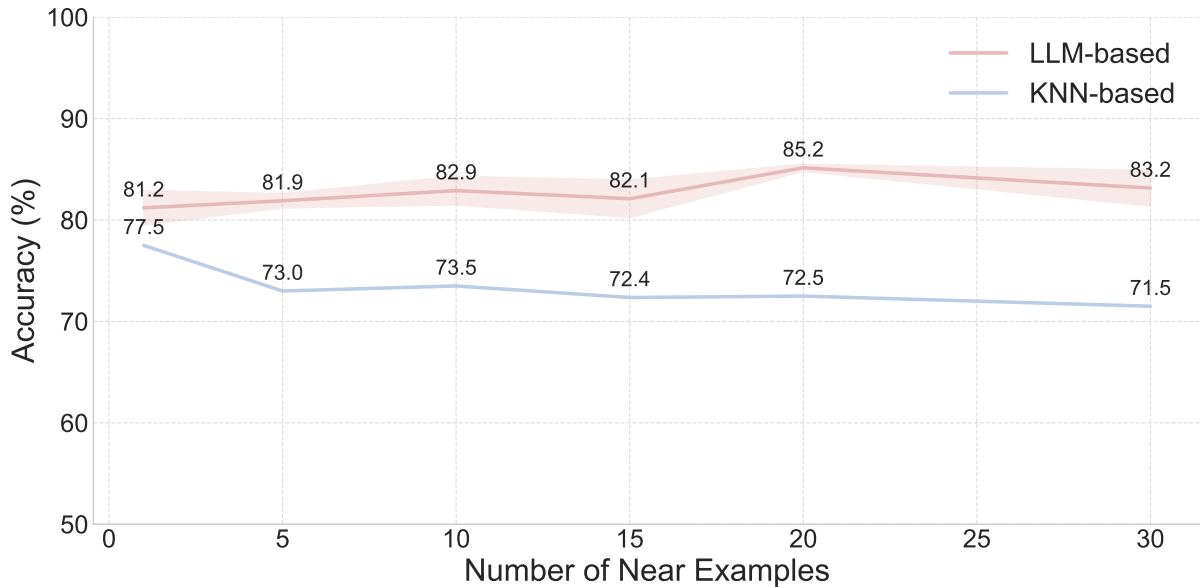


Figure A8: Impact of the number of near examples on accuracy. The plot shows mean accuracy and standard deviation for 1, 5, 10, 15, 20, and 30 near examples.

We tested using different numbers of near-examples, from 1 to 30, and found that more examples can actually help the model predict better, even though they might not seem as accurate when used in simpler models. Figure A8 shows this. The model's accuracy changes as we use more examples, with the best accuracy at 20 examples, suggesting this number might be just right for our model.

Interestingly, as we add more examples, the model’s predictions become a bit less consistent, but it gets better at predicting overall. This means that even if more examples don’t always lead to better results in simpler models, the LLM can use them to understand videos better and make more accurate predictions. However, using too many examples can actually make predictions a bit worse, showing there’s a sweet spot at 20 examples where the model is both accurate and consistent. This finding is important because it shows that the LLM can use more information to improve, but there’s a point where adding more doesn’t help as much. It also reminds us that while more examples can help, we need to consider how much work the model has to do. Looking at how different types of videos respond to more examples could help us fine-tune the model even more, making it better at predicting video popularity.

### A.6.3 Impact of temperature on model performance

In language models, temperature is a hyperparameter that controls the randomness of the model’s output. A lower temperature makes the model more deterministic in its predictions, while a higher temperature increases randomness and creativity. In the context of LLM based video popularity prediction, temperature plays a crucial role in balancing between making consistent, safe predictions and exploring more diverse, potentially insightful outcomes. Finding the optimal temperature is essential for maximizing the model’s predictive accuracy while maintaining its ability to generalize across various video types and popularity patterns.

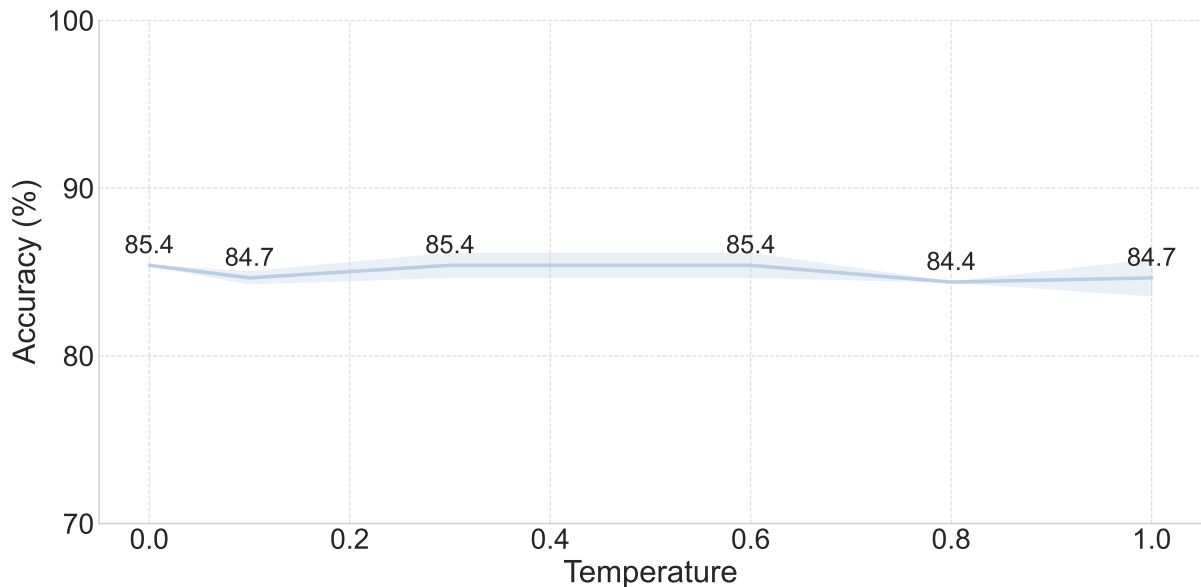


Figure A9: Impact of temperature settings. The plot shows mean accuracy for temperatures ranging from 0.0 to 1.0.

We looked at how different temperature settings affect the model’s accuracy to find the best balance. Figure A9 shows that the model works well across a range of temperatures, with 0.3 and 0.6 being optimal, both hitting 85.4% accuracy with little variation. This means the model can handle different temperatures well, but there’s a slight dip at 0.8 that needs more study. Future work could look closer at temperatures between 0.6 and 1.0 to understand why and improve predictions.

## A.7 Additional Ablation Study - Supervised Model

This section presents a detailed ablation study of our supervised model for video popularity prediction. By systematically analyzing the performance of different feature combinations, we aim to identify the most effective features and understand the trade-offs between model complexity and prediction accuracy. To investigate the effectiveness of different feature combinations in predicting video popularity, we conducted a series of experiments using various feature sets. We analyzed the performance of individual features, pairwise combinations, triple combinations, and complex feature sets. This comprehensive analysis aims to identify the most effective feature combinations and understand the trade-offs between model

complexity and prediction accuracy.

A.8 Additional Ablation Study - Baseline Multimodal Model

Feature	Embedding model	Embedding Size
Text (Title, Description, Captions)	MPNet	768
Image (Thumbnail)	CLIP	512
Video (Key Frame Aggregation)	VideoCLIP	512

Table A1: Baseline Embeddings

Component	Details
<b>Preprocessing Layer</b>	Uses a linear transformation to map input data to a fixed dimensionality (processed_dim).
<b>Main Layers</b>	Composed of fully connected layers with batch normalization and ReLU activation. Includes dropout for regularization. Layers are sequentially connected with increasing reduction in dimensionality (1024 → 512 → 256 → 128).
<b>Dimension Matching Layers</b>	Linear layers that adjust dimensions to enable addition of residual connections at each main layer stage.
<b>Attention Layer</b>	Consists of a linear transformation, a tanh activation, and a softmax output to produce attention weights.
<b>Final Classification Layer</b>	A fully connected layer that takes the attended features and outputs the final classification results.
<b>Overall Model Architecture</b>	Input data is processed through layers that include preprocessing, main processing with residuals, attention application, and final classification.

Table A2: Baseline Multimodal Model Description

To establish a strong foundation for comparison, we implement a baseline multimodal model that leverages deep learning techniques to predict video popularity. The architecture of this baseline model is described in Table A2. The model consists of several key components, including a preprocessing layer to handle variable input sizes, main layers to learn complex representations, dimension matching layers to facilitate residual connections, an attention layer to weigh the importance of different parts of the input data, and a final classification layer to produce the output.

The baseline model utilizes various embeddings to represent the different features of the video data, as detailed in Table ???. For textual features, such as the title, description, and captions, the MPNet model is employed to generate embeddings of size 768. Visual features, including the thumbnail, are processed using the CLIP model, resulting in embeddings of size 512. Finally, the video content is represented using key frame aggregation and the VideoCLIP model, producing embeddings of size 512. The baseline multimodal model serves as a robust point of comparison for our proposed framework, allowing us to assess the performance improvements achieved through the integration of VLMs and LLMs in the video popularity prediction task.

A.8.1 Individual feature performance

We begin by examining the predictive power of each feature type in isolation. Table A3 presents the performance scores for individual features.

As shown in Table A3, video features and description features achieved the highest individual performance with a score of 0.79, followed closely by thumbnail features (0.77) and caption features (0.76). Title features showed the lowest individual performance at 0.75, suggesting that while titles contribute to prediction, they may not be as informative as other features when used alone.

Feature	Score
Thumbnail	0.77
Video	<b>0.79</b>
Title	0.75
Description	<b>0.79</b>
Caption	0.76

Table A3: Performance Scores for Individual Features

### A.8.2 Feature combination performance

Next, we explore the synergistic effects of combining different feature types. Table A4 illustrates the performance scores for various feature combinations.











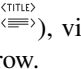

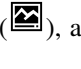
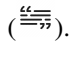
Combination				
-	<b>75</b>	<b>75</b>	73	74
	-	<b>76</b>	<b>76</b>	72
 & 	-	-	<b>77</b>	76
 &  & 	-	-	-	<b>78</b>

Table A4: Performance analysis of multimodal feature combinations for video popularity prediction. Icons represent title and description () , video content () , thumbnail () , and caption () . Bold numbers indicate the highest score in each row.

The combination of title features with other modalities consistently improved performance. Notably, the combination of thumbnail, description, and caption features achieved the highest score of 78%, demonstrating the complementary nature of these modalities in predicting video popularity.

### A.8.3 Complex feature combinations

Finally, we investigated the impact of combining four or five feature types. Table A5 shows the performance scores for these complex feature combinations.

Feature 1	Feature 2	Feature 3	Feature 4	Feature 5	Score
Videotext	Thumbnail	Video	Title	Description	0.81
Videotext	Thumbnail	Video	Title	Caption	0.81
<b>Videotext</b>	<b>Thumbnail</b>	<b>Video</b>	<b>Description</b>	<b>Caption</b>	<b>0.83</b>
Videotext	Thumbnail	Title	Description	Caption	0.82
Videotext	Video	Title	Description	Caption	0.80
Thumbnail	Video	Title	Description	Caption	0.82

Table A5: Performance Scores for Complex Feature Combinations

The combination of videotext, thumbnail, video, description, and caption features achieved the highest score of 0.83. However, it's important to note that this score is not significantly higher than some of the triple feature combinations, suggesting a point of diminishing returns in terms of prediction accuracy as we increase feature complexity.

These results highlight the importance of considering multiple modalities in video popularity prediction. While individual features provide valuable information, the combination of complementary features

leads to improved prediction accuracy. However, the experiments also reveal a trade-off between model complexity and performance gains, as the most complex feature combinations do not necessarily yield significantly better results than some simpler combinations. The analysis suggests that careful feature selection and combination can lead to efficient and effective video popularity prediction models. This combination likely captures a diverse range of information about the video content, including visual appeal, textual context, and spoken content.

These findings suggest that while incorporating multiple modalities can improve prediction accuracy, there is a point of diminishing returns. Future model development should focus on optimizing the balance between feature complexity and performance gains, potentially prioritizing the most informative feature combinations identified in this study.

### A.9 Comparison with Vision-Language Large Model

To further strengthen our contribution, we conducted additional experiments using an advanced Vision-Language Large Model (VLLM), Gemini. Specifically, we used the Gemini 1.5 Pro model, accessed as API in Vertex library. The input to the model included a combination of the summariser prompt and the final prediction prompt. The sequential frames were aligned with the video’s existing captions to ensure that visual and verbal elements were synchronized before being fed into the LLM to generate the video-to-frame summary for the entire video. This process closely follows the steps described in Section 3.2 where “Frame Extraction” step (extracting 5 frames per minute) and the “Caption Matching and Data Integration” step were employed to align captions and frames for LLM processing. Subsequently, "Frames to Text Conversion and Summarization" step generated the final video summary using an LLM call. For prediction, we used the same final prompt detailed in Figure A2, ensuring consistency with our primary method.

These experiments confirmed that our strategy—incorporating sequential prompting, hypothesis generation, and supervised signals—consistently improves prediction performance, even with this state-of-the-art model (see Figure A10. This underscores the generalizability and robustness of our approach across different model architectures.

Notably, this result demonstrates that more advanced models do not inherently outperform others across all aspects; rather, each model tends to excel in specific areas. This highlights the importance of strategically combining models based on their unique strengths to achieve optimal results. Additionally, the prompting strategies developed in this work offer practical guidance for designing effective multimodal solutions. These insights not only inform future research on multimodal learning but also emphasize the value of integrating pre-trained models tailored to specific task requirements.

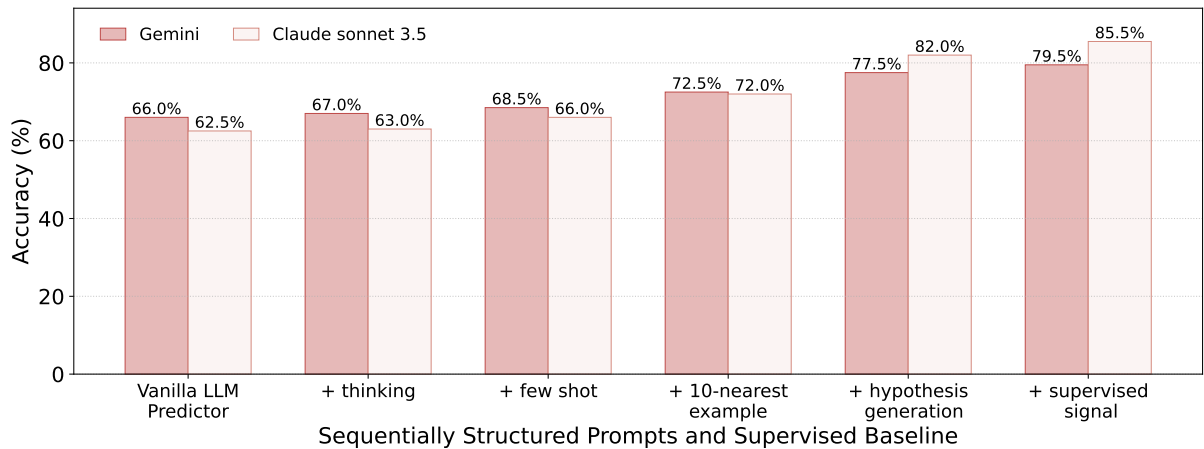


Figure A10: Performance evaluation of the Gemini 1.5 Pro model for video popularity prediction, demonstrating the impact of sequential prompting, hypothesis generation, and supervised signals.

## A.10 Addressing Language Imbalance in Target Classes

To assess potential language imbalance in classifying ‘local hit’ and ‘global big hit,’ we analyzed the dataset’s language distribution and its impact on model performance. Figure A11 shows the distribution of target classes across major languages in the entire dataset (red) and the subset used in our experiment (blue). English accounts for approximately 40% of the ‘global big hit’ category and 15% of the ‘local hit’ category, indicating an overrepresentation in the ‘global big hit’ class. However, non-English languages, including Spanish, Hindi, and French, contribute significantly to both categories, suggesting that the dataset maintains reasonable language diversity.

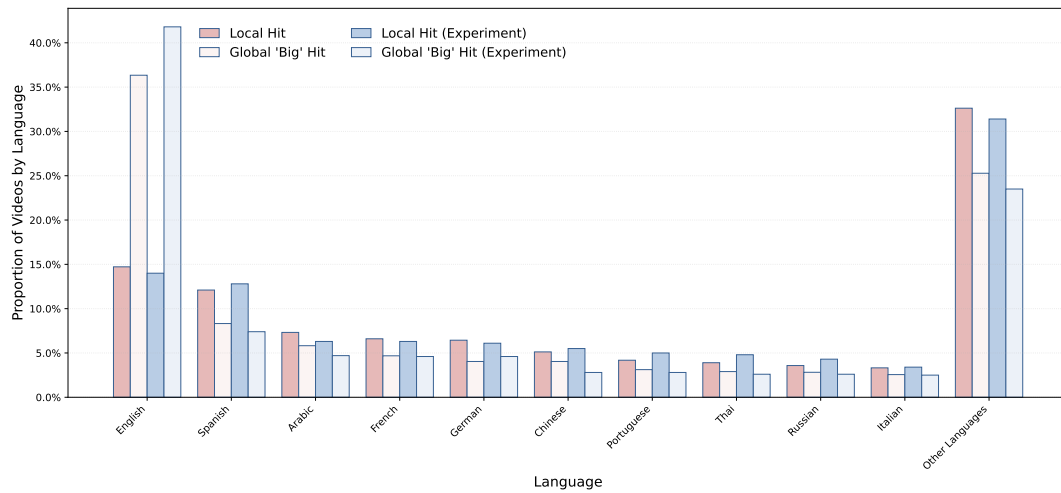


Figure A11: Distribution of ‘local hit’ and ‘global big hit’ categories across major languages. English accounts for 40% of the ‘global big hit’ category and 15% of the ‘local hit’ category, with other languages contributing substantially to both.

To evaluate potential language bias, we analyzed prediction accuracy across languages. As shown in Figure A12, the model performs comparably for both major (e.g., English, Spanish) and minor (e.g., Swedish, Indonesian) languages, with no significant bias favoring English or other dominant languages. These results suggest that language representation does not disproportionately affect the model’s predictions.

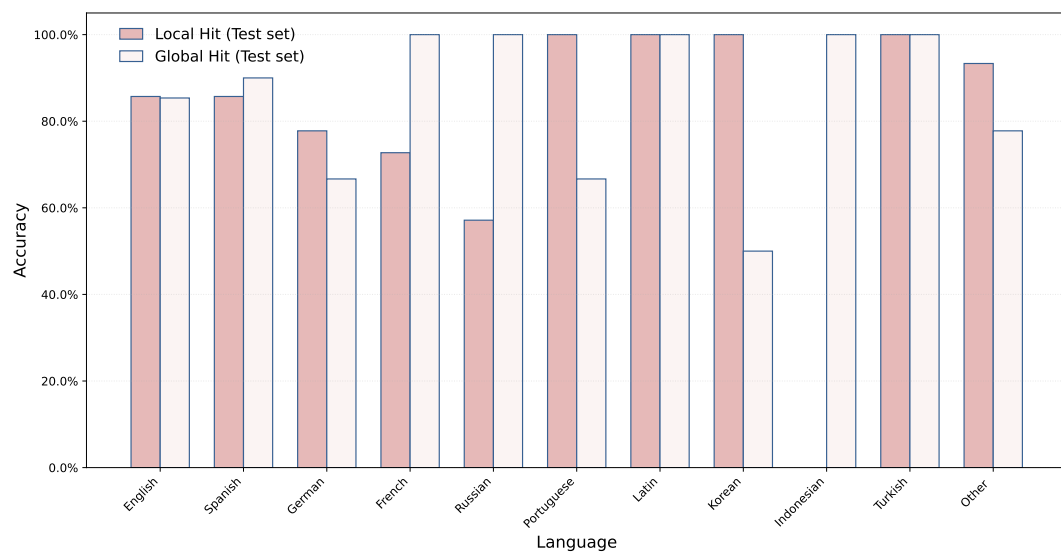


Figure A12: Prediction accuracy for ‘local hit’ and ‘global big hit’ categories across languages. The results show comparable model performance across major (e.g., English, Spanish) and minor (e.g., Swedish, Indonesian) languages, indicating that prediction accuracy is not significantly influenced by language representation.



Figure A13: Illustration of our LLM-based pipeline for explainable video popularity prediction. Shown are two successful predictions (outlined in black) and one misclassification (outlined in red), highlighting both the strengths and challenges of the approach.

### A.11 Additional Qualitative Analysis

In this section, we present a deeper qualitative analysis of our LLM-based framework by highlighting two successful predictions and one misclassification (Figure A13). These examples offer additional insights into how our method processes visual and textual information to generate final predictions and accompanying explanations.

**Positive Examples** Two videos, “Mexico vs. Brazil Highlights” and “Minecraft Survivor VS 3 Hitmen,” illustrate how our pipeline effectively identifies factors driving popularity. In the football highlights, the model attributes the video’s success to the involvement of the internationally recognized Brazilian national team and star players such as Vinicius and Richarlison. This demonstrates the framework’s capacity to capture both *engagement intensity* (high view counts driven by star power) and *geographic spread* (global appeal of Brazil’s national team and the event, FIFA World Cup). Similarly, for the Minecraft video, the model recognizes a unique speedrun challenge and engaging personalities, reflecting its ability to detect cross-border relevance within gaming communities.



**Negative Example** The misclassification of the “Dad Jokes” video reveals potential challenges in capturing cultural or emotional nuances. Although the pipeline identifies the universal appeal of humor and the involvement of popular creators, it overestimates the video’s worldwide reach. This suggests that subtle context factors (e.g., culturally specific humor) can be difficult for the model to fully grasp, pointing to areas where additional sentiment analysis or more nuanced cultural embeddings might prove valuable.

1281  
1282  
1283  
1284  
1285

**Discussion** Overall, these qualitative examples demonstrate that our pipeline can effectively integrate textual summaries, captions, and high-level contextual signals to predict video popularity. While the misclassification underscores the difficulty of handling culturally specific or emotionally resonant content, the two correct predictions illustrate the framework’s ability to capture global appeal and engagement drivers across diverse video genres. Future work could delve more deeply into cultural features, sentiment modeling, or real-time data streams to further refine the approach.

1286  
1287  
1288  
1289  
1290  
1291