

Motivation

- Creators lack tools to predict short video engagement before publication
- Limited editing experience makes social media video optimization challenging
- Need for automated video refinement to boost engagement potential

Challenges

Data & Annotation: Lack of fine-grained labels linking editing actions (e.g., story flow, music, text, transitions) to engagement

Evaluation & Generalizability: Existing engagement prediction methods depend on platform-specific, post-publication metrics whose insights are not generalizable

Feedback & Interpretability: Current SOTA can understand actions and objects (pixel-level) information in a short-video but cannot ground editing operations to user behavior (engagement) for efficient editing-driven interpretable feedback

SmartEdit

We propose SmartEdit, which analyses short videos on the basis of video editing operations.

It is able to

- Predict an **engagement score**
- Provide a **ranking of editing operations (Edit Signals)** that positively or negatively influence engagement
- Suggest **interpretable feedback** for creators to enhance content to improve engagement

References

Dasong Li et al., "Delving deep into engagement prediction of short videos," in ECCV 2025.

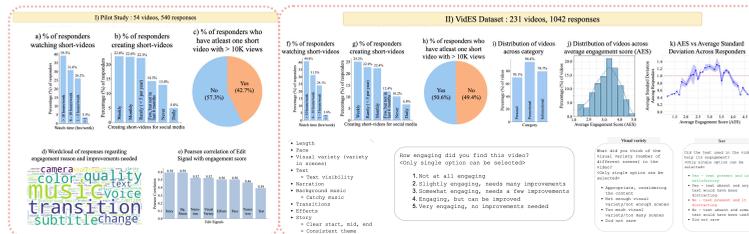
Haooning Wu et al., "Q-align: Teaching llms for visual scoring via discrete text-defined levels," in ICML 2024

Zhanyu Wang et al., "Gpt4video: A unified multimodal large language model for instruction-followed understanding and safety-aware generation," in ACM ICM 2024.



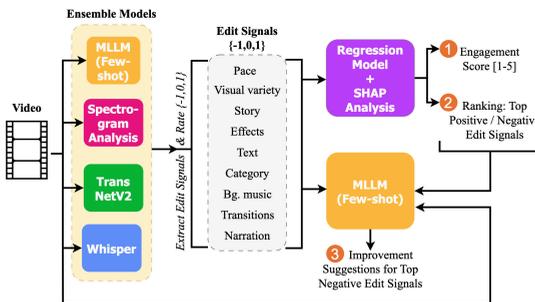
Proposed Approach

- Curate a **novel dataset, VidES**, containing short-videos, engagement scores, and detailed human evaluations of Edit Signals, establishing a foundation for data-driven approaches to short-video refinement.

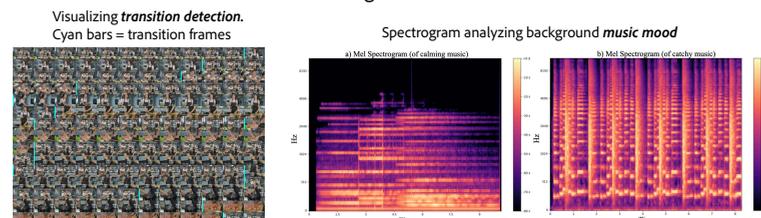


Validation study & full-scale dataset collection along with snippets of survey questions asked to the users

- Propose **SmartEdit**, a **two-stage framework** that decomposes short-video engagement prediction into interpretable video editing operations (Edit Signals):



- Short-video Engagement Prediction:** Random forest regressor trained on VidES data to predict an engagement score based on extracted edit signals
- Interpretable feedback for improvement:** a) SHAP analysis to rank top engagement driving and harming Edit Signals, and b) few-shot MLLM to generate actionable improvement suggestions grounded in structured user feedback.



Results

- SmartEdit **outperforms baselines** (random, Q-Align, GPT-4o zero/few-shot) in **engagement prediction and feedback accuracy**
- Ablation studies show the importance of both good and bad examples, and reasoning in prompting
- Qualitative examples demonstrate iterative improvement of videos using SmartEdit's feedback

a) Video link: <https://rb.gy/d7but0>

1 Engagement Score: 2.11
Ground Truth: 1.40

2 Top Edit Signals: Story, Effects, Visual variety, Pace, Bg. music, Narration

3 Improvement Feedback: What needs to improve?
a) Clear start, middle, end
b) Face
c) Effects missing and some effects would have been helpful
d) Zoom (e.g. close-ups, pan)

b) Video link: <https://rb.gy/um8ora>

1 Engagement Score: 4.18
Ground Truth: 4.40

2 Top Edit Signals: Text

3 Improvement Feedback: What needs to improve?
a) Less interference of the text with the visuals
g) Enhance visibility of the text (e.g., contrast with background)

Performance on Engagement Prediction & Improvement Feedback

Method	Engagement Prediction			Improvement Feedback	
	ES Decomp.	MAE ↓	RMSE ↓	clfAcc ↑	F1 ↑
Random model	✗	1.45	1.77	0.21	0.27
Q-Align [20]	✗	0.79	0.97	0.39	N/A
GPT-4o (Zero-shot)	✗	1.19	1.39	0.13	0.36
GPT-4o (Few-shot)	✗	1.02	1.23	0.28	0.41
GPT-4o (Few-shot)	✓	0.75	0.88	0.38	0.45
SmartEdit	✓	0.64	0.84	0.49	0.51

Ablation study testing different backbone for engagement prediction model

Regression Model	MAE ↓	RMSE ↓	clfAcc ↑
XGBoost	1.15	1.33	0.28
Multi-Layer Perceptron (MLP)	0.91	1.18	0.37
Random Forest Regressor	0.64	0.84	0.49

Ablation study on feedback improvement

#Good	#Bad	Reasoning?	Accuracy ↑	F1 ↑
0	4	✗	0.44	0.59
0	4	✓	0.47	0.63
2	2	✗	0.45	0.62
2	2	✓	0.51	0.67