



Dublin, Ireland 27-31.10.2025



Combating Online Misinformation Videos: Characterization, Detection, and Prevention

Qiang Sheng¹Peng Qi²Tianyun Yang³Yuyan Bu⁴Wynne Hsu² Mong Li Lee²Juan Cao¹¹Institute of Computing Technology, Chinese Academy of Sciences ²National University of Singapore³Shenzhen Research Institute of Big Data ⁴Beijing Academy of Artificial Intelligence

October 27, 2025 13:30-17:00 | Goldsmiths 1, Dublin Royal Convention Centre, Ireland

Agenda: 3-hour Talk plus 30-minute Break

Time	Section	Presenter
13:30-13:40	Introduction & Motivation	Qiang Sheng
13:40-13:55	Preliminaries: Video Editing & Generation	Qiang Sheng
13:55-14:10	Characterization of Misinformation Videos	Qiang Sheng
14:10-14:50	Detection Part I: Human-Edited Misinformation	Yuyan Bu
14:50-15:30	Detection Part II: AI-Generated Misinformation	Tianyun Yang
15:30-16:00	Coffee Break	/
16:00-16:40	Prevention Strategies	Peng Qi
16:40-17:00	Conclusion & Open Discussion / General QA	Qiang Sheng / All

Clarification questions are welcomed during the talk

Tutorial Outline

1. Introduction & Motivation

Background

Effects and Concerns

Goals of Our Tutorial

2. Preliminaries: Video Editing & Generation

Overview of Generative Models and Diffusion

Video Generation

Video Editing

Q+A/Discussion



Qiang Sheng

Tutorial Outline

3. Characterization of Misinformation Videos

Definition

Datasets Overview

Analysis on FakeSV

News Content

Social Context

Propagation

Q+A/Discussion



Qiang Sheng

Tutorial Outline

4. Detection Part I: Human-Edited Misinformation

Signal-based detection

Editing Traces

Generation Traces

Semantic-based detection

Seek clues within the sample's own content

Seek clues from external information

Intent-based detection

Social context

Clue integration for misinformation video detection

Parallel Integration

Sequential Integration

Q+A/Discussion



Yuyan Bu

Tutorial Outline

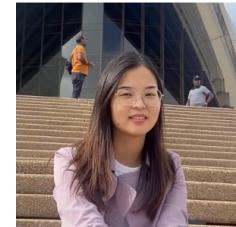
5. Detection Part II: AI-Generated Misinformation

Manipulated video detection

Generated video detection

Attributing AI-generated Content to the Source Model

Q+A/Discussion



Tianyun Yang

Tutorial Outline

6. Prevention Strategies

Creation Prevention

Embedding Tamper-proof Digital Identifier

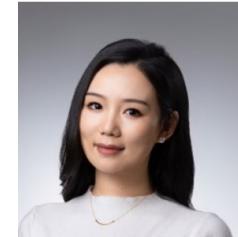
Mitigating Hallucination in Content Generation

Spread Prevention

Alert, Verification, and Resilience Building

Controlling the Spread of Misinformation

Promoting Truth and Debunking



Peng Qi

Tutorial Outline

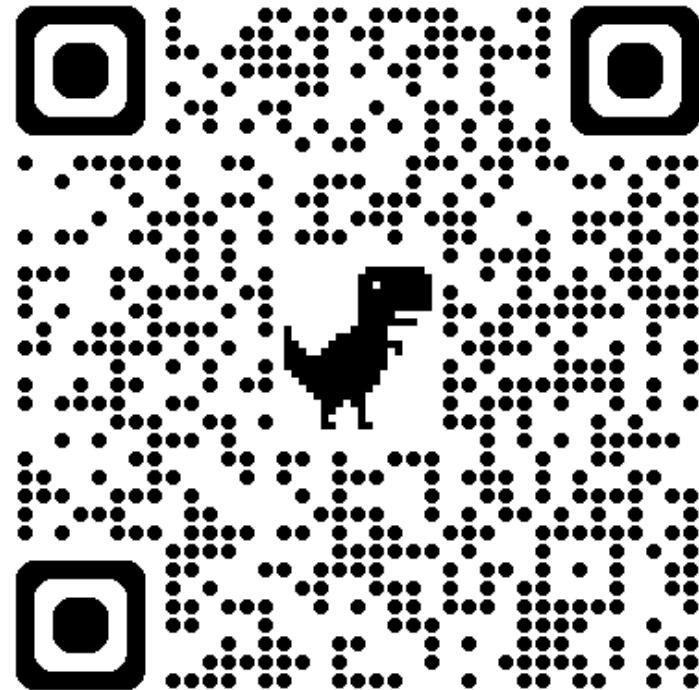
7. Conclusion & Open Discussion

plus General Q+A



Qiang Sheng

Slides & Reading List



<https://misinfo-video.github.io/>

Introduction & Motivation

Section 1

Tutorial Outline

1. Introduction & Motivation

Background

Effects and Concerns

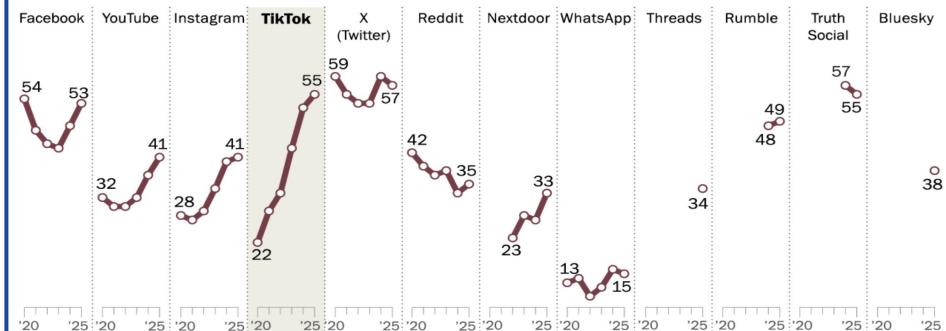
Goals of Our Tutorial

Video has been a popular form of news consumption

Short video-sharing platforms like TikTok and others has been important access to get daily news for users.

More than half of U.S. adult TikTok users get news there, up from 22% in 2020

% of each social media site's users who regularly get news there



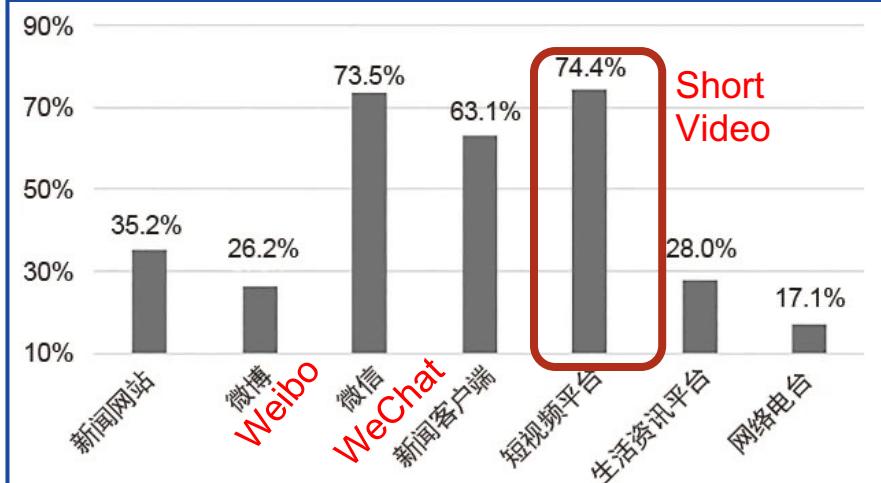
Note: The other response option was "No, don't regularly get news on this." Only respondents who indicated that they use each site were asked if they regularly get news on it. Social media sites are shown left to right in descending order by the share of U.S. adults who regularly get news there.

Source: Survey of U.S. adults conducted Aug. 18-24, 2025.

Longitudinal Comparison
 US adults that get news from
 TikTok 2020-2025: 22% -> 55%

Pew Research Center. 1 in 5 Americans now regularly get news on TikTok, up sharply from 2020.

CNNIC. Survey on the Current Status of News Access Channels for Chinese Netizens. (in Chinese)



Horizontal Comparison
 Short video platforms **surpass WeChat and Weibo**, becoming the main access for Chinese netizens.

Meanwhile, video generation techniques has rapid progress

Trend #1: More general. Video generation is not only face-swapping.

Before

Face swapping, editing, Reenactment



Recent

General video generation guided by text/image prompts



Meanwhile, video generation techniques has rapid progress

Trend #2: More vivid.

Video generation can be of high resolution with details.



~60s, multiple shots
physical details (though imperfect)



Large view
Dynamic shot

Tutorial Outline

1. Introduction & Motivation

Background

Effects and Concerns

Goals of Our Tutorial

Along with such a progress, video misinfo is easier to make

Misinformation was delivered as text-only, or image-included posts.
 Now, producing misinformation videos are easier than before.

Fact Check

Trump has been disqualified from receiving Nobel Prize? Here's the truth

An image supposedly showing a press release by The Associated Press spread the rumor.



BREAKING NEWS: Nobel Prize committee announces that @realDonaldTrump has been permanently disqualified from all future awards due to his renaming US Defense Dept "Department of War."
 #TrumpisaNationalDisgrace #NobelPeacePrize #gwb5

Fake Text

Trump is disqualified from Nobels?

<https://www.snopes.com/fact-check/trump-nobel-prize-disqualified/>

<https://apnews.com/article/fact-check-trump-NYPD-stormy-daniels-539393517762>

<https://www.youtube.com/watch?v=2wqWnoMg1dU>



Fake Image

Trump was arrested?
 (Using Generative AI)

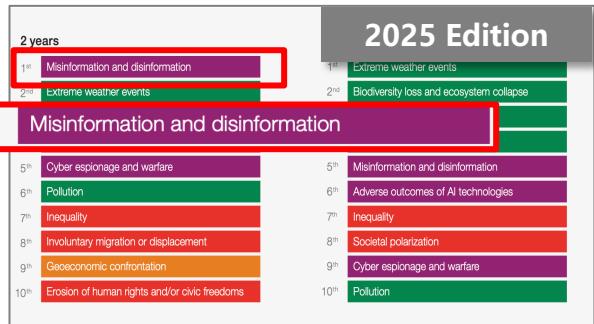


Fake Video

Trump fought with Zelensky?
 (Using Generative AI)

The Worrying Trend: AI Video Faking is Industrialized

Surveys and predictions show deep concerns regarding AI misinfo

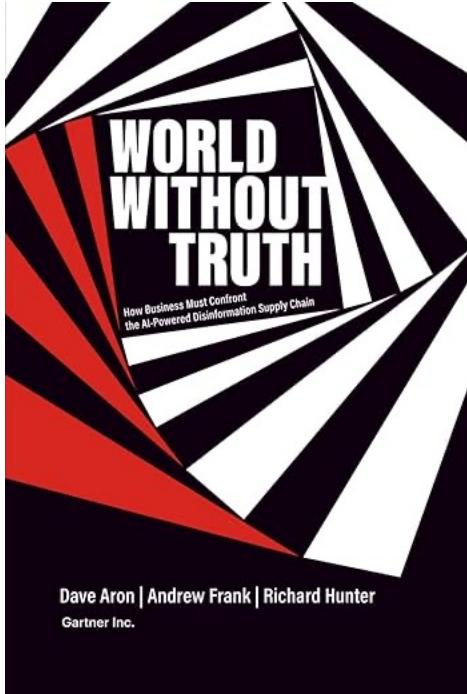


World Economic Forum, Global Risks Report:
Misinformation and disinformation is the TOP1 two-year risk and TOP1 ten-year technical risk



State of AI 2025 Report make a prediction:
A deepfake/agent-driven cyber attack triggers the first NATO/UN emergency debate on AI security.

Surveys and predictions show deep concerns regarding AI misinfo



- Disinformation is now **a sophisticated, organized business** with its own supply chain.
- **Generative AI accelerates the creation and spread** of convincing fake content, manipulating narratives and exploiting biases at scale.

With the support of recent AI techniques, misinformation video faking can be with a larger scale with lower cost



A Real-World Case Reported by China Central TV

A team at Zhejiang published **28.5k misinformation videos** by cutting & editing video clips, attracting **2.7 billion views**

Journal of Computer-Mediated Communication

Seeing Is Believing: Is Video Modality More Powerful in Spreading Fake News via Online Messaging Apps?

S. Shyam Sundar  ¹, Maria D. Molina  ², & Eugene Cho³

¹Media Effects Research Laboratory, Penn State University, University Park, PA 16802, USA

²Department of Advertising and Public Relations, Michigan State University, East Lansing, MI 48824, USA

³Department of Communication Studies, The College of New Jersey, Ewing, NJ 08628, USA

... It is clear from our findings that **video is causing individuals to perceive fake news as more credible than audio and text**, and increases the likelihood of them spreading it. ...

Fake news with videos
is more convincing for human
and easier to spread

As Good as a Coin Toss

Human Detection of AI-Generated Images, Video, Audio, and Audiovisual Stimuli

DI COOKE, Department of War Studies, King's College London, London, UK

ABIGAIL EDWARDS, Center for Strategic and International Studies, Washington D.C. USA

SOPHIA BARKOFF, Center for Strategic and International Studies, Washington D.C. USA

KATHYRN KELLY, Center for Strategic and International Studies, Washington D.C. USA

We find that **on average, people struggled to distinguish between synthetic and authentic media, with the mean detection performance close to a chance level performance of 50%**. We also find that accuracy rates worsen when the stimuli contain any degree of synthetic content, features foreign languages, and the media type is a single modality.

The synthetic media is hard to distinguish for ordinary people

For Detectors: New challenges posed by misinfo videos

High information heterogeneity brought by various modalities

1

- Requires a stronger and more comprehensive understandability
- Brings more uncertainty and even noise to the final prediction

2

Blurred distinction between misleading video manipulation and non-malicious artistic video editing

- Artistic editing like beautifying faces is general, making it hard to see the malicious manipulation traces.

3

New patterns of misinformation propagation due to the dominant role of recommendation systems on online video platforms

- Less social contexts than before on Twitter/Weibo. Brings new behaviors.
- Requires new detection and prevention methods.

Tutorial Outline

1. Introduction & Motivation

Background

Effects and Concerns

Goals of Our Tutorial

Goals of this Tutorial

Key Question

How to **Characterize, Detect, and Prevent Misinformation Videos?**

Survey Range

Covers typical or new works in this direction in recent five years

Best for

Audience that have knowledge about AI safety/multimedia content safety or is interested in combating misinformation issues

Preliminaries: Video Editing & Generation

Section 2

Tutorial Outline

2. Preliminaries: Video Editing & Generation

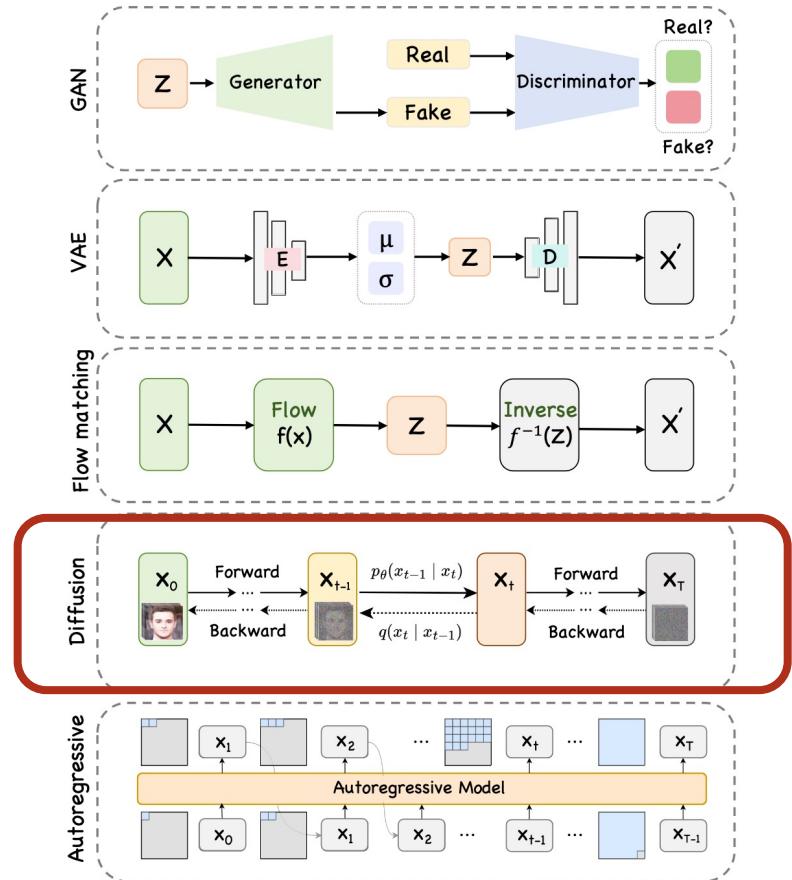
Overview of Generative Models and Diffusion

Video Generation

Video Editing

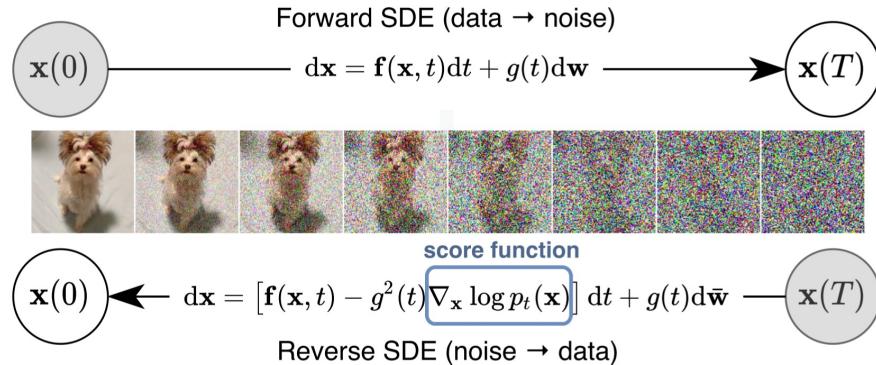
Q+A/Discussion

Generative Models



Our focus in this part

Image Diffusion Model



Basic Idea: Learning to add/remove noise

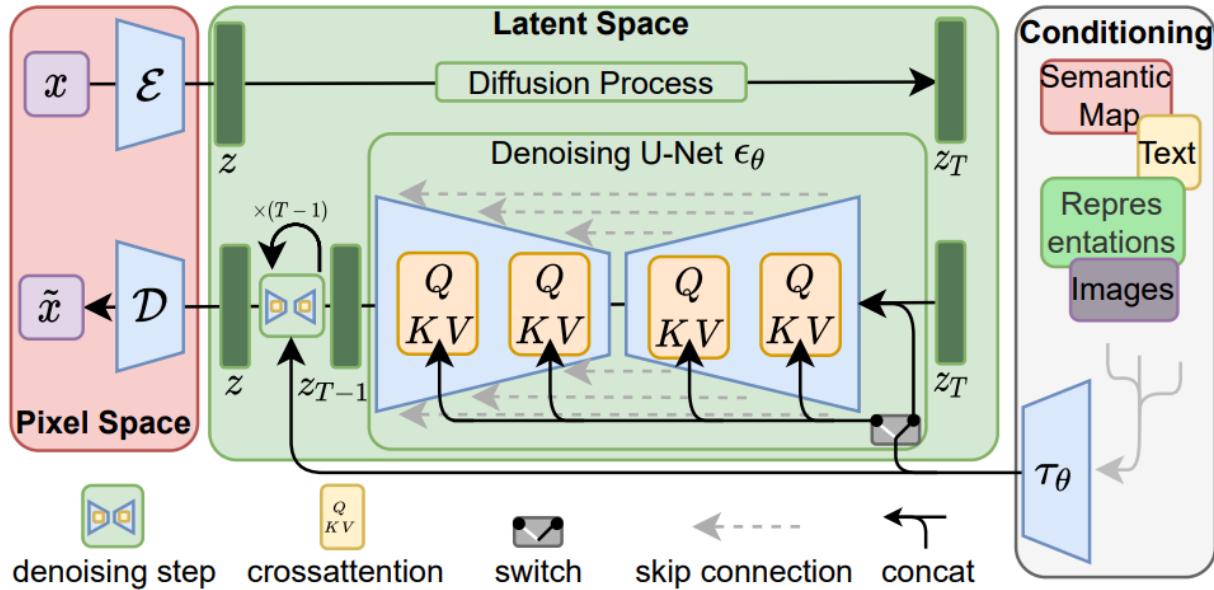
Generate an *image* as *building a house*.

To learn it, separate the sample house, know each step, and then try to rebuild it

Too expensive to use in reality

If we want a 1024*1024 image, we need to input a 1024*1024 noise image, requiring high computation overhead.

Latent Diffusion Model



Common usage:

- Use VQVAE to transform into latent space
- A U-Net based diffusion model
- Controlled by cross-attention condition

Tutorial Outline

2. Preliminaries: Video Editing & Generation

Overview of Generative Models and Diffusion

Video Generation

Video Editing

Q+A/Discussion

From Image Generation to Video Generation



Image Generation (2022-)
Stable Diffusion, Dalle3...

Video Generation (2023-)
Stable Video Diffusion, Sora...

How to extend the usage of diffusion models to video?

Video Diffusion Models (VDM)

Build VDM Arch.

High-Impact VDM

Sora...



Make-A-video
(Meta, 2022)

Align-your-latent
(NVIDIA, 2023)

SVD
(2023)

Sora...

Focuses: High-quality, Controllable, and Arbitrary Length

Make-A-Video (Meta, 2022)

- Extend Text-to-Image Model to 3D Model
 - 2D Conv -> Pseudo 3D Conv
 - Spatial attention -> Pseudo 3D attention
- Frame Interpolation
 - Using mask prediction
- Spatial Super-Resolution
- Issue:
 - Low quality (a few seconds, 256*256, blurs)

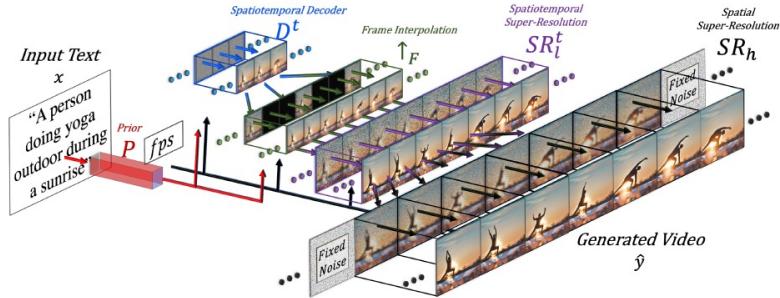


Figure 2: **Make-A-Video high-level architecture.** Given input text x translated by the prior P into an image embedding, and a desired frame rate fps , the decoder D^t generates 16 64×64 frames, which are then interpolated to a higher frame rate by \uparrow_F , and increased in resolution to 256×256 by SR_l^t and 768×768 by SR_h , resulting in a high-spatiotemporal-resolution generated video \hat{y} .

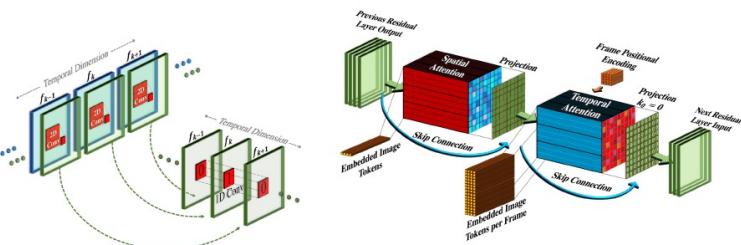
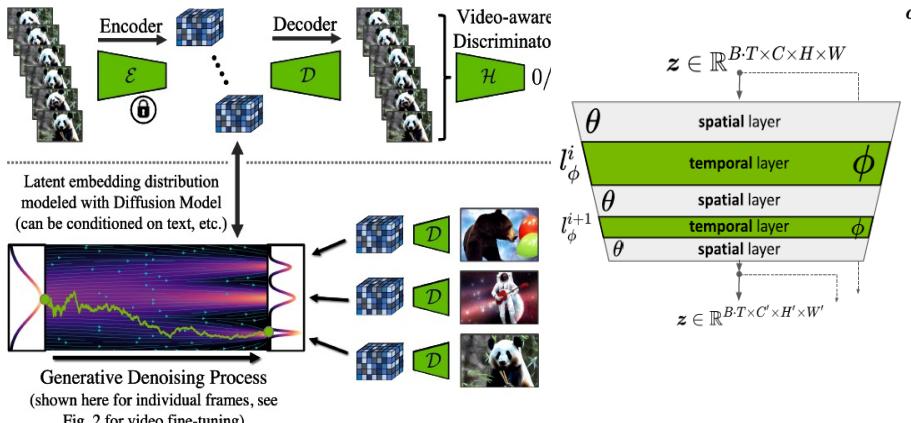


Figure 3: **The architecture and initialization scheme of the Pseudo-3D convolutional and attention layers, enabling the seamless transition of a pre-trained Text-to-Image model to the temporal dimension.** (left) Each spatial 2D conv layer is followed by a temporal 1D conv layer. The temporal conv layer is initialized with an identity function. (right) Temporal attention layers are applied following the spatial attention layers by initializing the temporal projection to zero, resulting in an identity function of the temporal attention blocks.

VideoLDM (NVIDIA, 2023)

- Practice **latent** VDM paradigm
- Can generate 2k resolution videos
- Long video generation: Keyframe + Mask Prediction



Video-level
VQ-VAE training

Text-to-Image ->
Text-to-Video

Align your Latents: High-Resolution Video Synthesis with Latent Diffusion Models

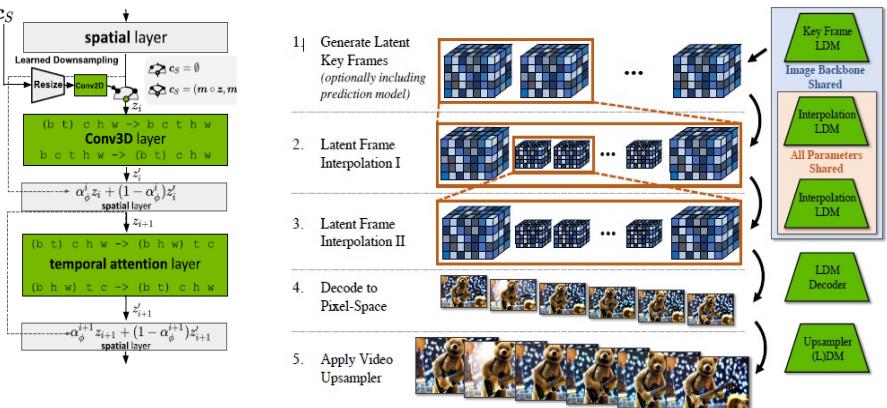
Align your Latents:

High-Resolution Video Synthesis with Latent Diffusion Models

Andreas Blattmann^{1 *†} Robin Rombach^{1 *†} Huan Ling^{2,3,4 *} Tim Dockhorn^{2,3,5 *†}
 Seung Wook Kim^{2,3,4} Sanja Fidler^{2,3,4} Karsten Kreis²

¹LMU Munich ²NVIDIA ³Vector Institute ⁴University of Toronto ⁵University of Waterloo

Project page: <https://research.nvidia.com/labs/torontolab/VideoLDM/>



Generate Keyframes ->
Latent Frame Interpolation->
Video Upsampling

Stable Video Diffusion (Stability AI, 2023)

- Based on VideoLDM architecture
- Better data cleaning for high-quality video data
 - With a useful video data processing pipeline:
 - Cut detection for video-cutting
 - Use CLIP/BLIP to obtain captions
 - Filter out static data based on optical flow
- High-quality video data -> High generation quality

Stable Video Diffusion: Scaling Latent Video Diffusion Models to Large Datasets

Andreas Blattmann* Tim Dockhorn* Sumith Kulal* Daniel Mendelevitch
 Maciej Kilian Dominik Lorenz Yam Levi Zion English Vikram Voleti
 Adam Letts Varun Jampani Robin Rombach
 Stability AI



"A robot dj is playing the turntables, in heavy raining futuristic tokyo, rooftop, sci-fi, fantasy"



"An exploding cheese house"



"A fat rabbit wearing a purple robe walking through a fantasy landscape."



Figure 1. Stable Video Diffusion samples. Top: Text-to-Video generation. Middle: (Text-to-)Image-to-Video generation. Bottom: Multi-view synthesis via Image-to-Video finetuning.

Sora (OpenAI, 2024)

- Breakthrough Achievements
 - 4K resolution
 - High spatial consistency
 - Diverse scenes with multiple shots
 - Longer video (60s)
- Quality issues ->
Commonsense/Physical violation
issues

Research

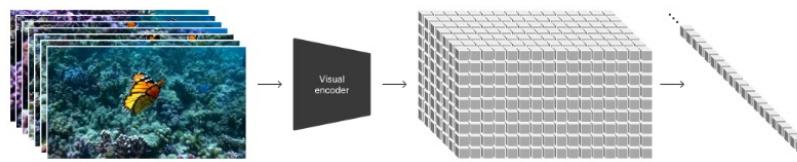
Video generation models as world simulators

We explore large-scale training of generative models on video data. Specifically, we train text-conditional diffusion models jointly on videos and images of variable durations, resolutions and aspect ratios. We leverage a transformer architecture that operates on spacetime patches of video and image latent codes. Our largest model, Sora, is capable of generating a minute of high fidelity video. Our results suggest that scaling video generation models is a promising path towards building general purpose simulators of the physical world.



Sora (OpenAI, 2024)

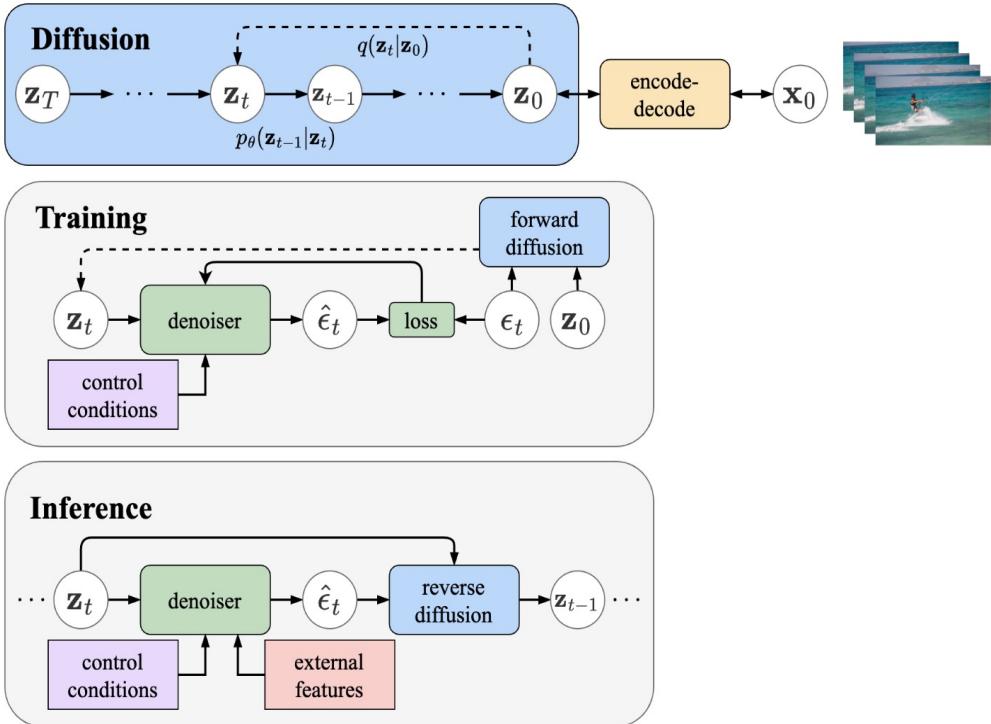
- **Cora Idea:** Turning visual data of all types into a unified representation that enables large-scale training of generative models
- **Turning visual data into patches**



- **Video compression network → Video VQ-VAE**
- **Spatiotemporal latent patches** → tokenization, arbitrary length and resolution
- **Scaling transformers for video generation**
 - Based diffusion transformer architecture (DiT)
 - Support

Diffusion-Based Video Editing

- Similar to Video Generation
 - Consider controllable signals more
- Text-base conditioning
- Point Conditioning
 - DragDiffusion, DragVideo...
- Pose Conditioning
 - Follow-Your-Pose...



Characterization of Misinformation Videos

Section 3

Tutorial Outline

3. Characterization of Misinformation Videos

Definition

Datasets Overview

Analysis on FakeSV

Signal Level

Semantics Level

Intent Level

Q+A/Discussion

Definition & Taxonomy

- **Misinformation Video**

- A video post that conveys false, inaccurate, or misleading information.

- **Misinformation Video Detection**

- Predict whether the video post v contains misinformation given all the accessible features

$$\mathcal{F} : \mathcal{E} \mapsto \{0, 1\}$$

- For some recent works, the output also contains a natural language text beyond the binary classification labels to provide human-understandable explanations.



Title: Dems voting for Nikki
 #nikkihaley #democrats

Explanation: The title suggests that Democrats are voting for Nikki Haley, which is misleading. The audio and video summaries indicate that the speaker, while acknowledging Haley's capabilities, still prefers Biden and expresses uncertainty about supporting Haley. This suggests that the title exaggerates the level of Democratic support for Haley.

Definition

Text

Microchipped my the government from the Covid vaccine!
#covid #vaccine

Username
Additional Name
114 Following 644 Followers 17.5K Likes
304 ♦
Taken 🔒

Social Context (Publisher Profile)

Engaged User
not true!!! its been proven false!!!
2021-5-22 6

Social Context (User Response)

Video

Covid vaccine = microchipped my government 🤡

Audio

Category	Fields
Content	video, cover image, title, published time
Response	# of likes/stars/comments, top 100 comments (with reviewed time, # of likes and # of sub-comments)
Publisher	info_verified, info_introduction, current IP location, # of fans/subscribes/likes/videos and top 100 published videos' covers

Tutorial Outline

3. Characterization of Misinformation Videos

Definition

Datasets Overview

Analysis on FakeSV

News Content

Social Context

Propagation

Q+A/Discussion

Datasets Overview

Real-world data:

Table 1: Summary of datasets of fake news video detection. Metadata refers to basic statistics such as # of likes/stars/comments. “-” represents open-domain. Names of sources are abbreviated for simplicity (YT: YouTube, TW: Twitter, FB: Facebook, TT: TikTok, BB: Bilibili, DY: Douyin, KS: Kuaishou).

Dataset	Features					Instances (fake/real)	Domain	Language	Released	Source
	Video	Title	Metadata	Comment	User					
(Papadopoulou et al. 2018)	✓	✓	✓	✓		2,916/2,090	-	En,Fr,Ru,Ge,Ar	Y	YT,TW,FB
(Palod et al. 2019)	✓	✓	✓	✓		123/423	-	En	Y	YT
(Hou et al. 2019)	✓		✓			118/132	prostate cancer	En	N	YT
(Medina et al. 2020)	✓	✓		✓		113/67	COVID-19	En	N	YT
(Choi and Ko 2021)	✓	✓		✓		902/903	-	En	N	YT
(Shang et al. 2021)	✓	✓				226/665	COVID-19	En	N	TT
(Li et al. 2022)	✓		✓			210/490	health	Ch	N	BB
FakeSV	✓	✓	✓	✓	✓	1,827/1,827	-	Ch	Y	DY, KS
FakeTT	✓	✓			✓	1,172/819	-	En	Y	TT

- **FakeSV and FakeTT follows highly-similar data curation and annotation pipelines**
- **Widely-used in existing works**

Dataset	Time Range	Avg Duration (s)	#Fake	#Real	#All
FakeSV	2017/10-2022/02	39.88	1,810	1,814	3,624
FakeTT	2019/05-2024/03	47.69	1,172	819	1,991

FakeSV: A Multimodal Benchmark with Rich Social Context for Fake News Detection on Short Video Platforms
 FakingRecipe: Detecting Fake News on Short Video Platforms from the Perspective of Creative Process

Datasets Overview

Examples from FakeSV



Title: How did this 3.8 kg gold nuggets get taken out by this Chinese lad? Let's take a look!



Title: Emergency by West Lake in Hangzhou this morning, waiting for reinforcements! #special police



Description: In Qinzhou, Guangxi, the wife was brutally beaten in the street for refusing to pay off her husband's gambling debts... What a heinous act!

On-screen Text: A man demands money from his wife to pay off gambling debts. Upon refusal, he assaults her, dragging and then body-slamming her on the road. Family violence again! Urgent need for intervention!



Description: #Shaanxi Shangluo residents transport anti-epidemic supplies with mules: A publicity stunt or full of sincerity?

#TopVQuickComment

On-screen Text: Shaanxi Shangluo locals use mules to deliver pandemic supplies: seeking attention or full of compassion? Someone question the efficiency. That's a mutual aid from all sides in times of trouble. Truly Touching.

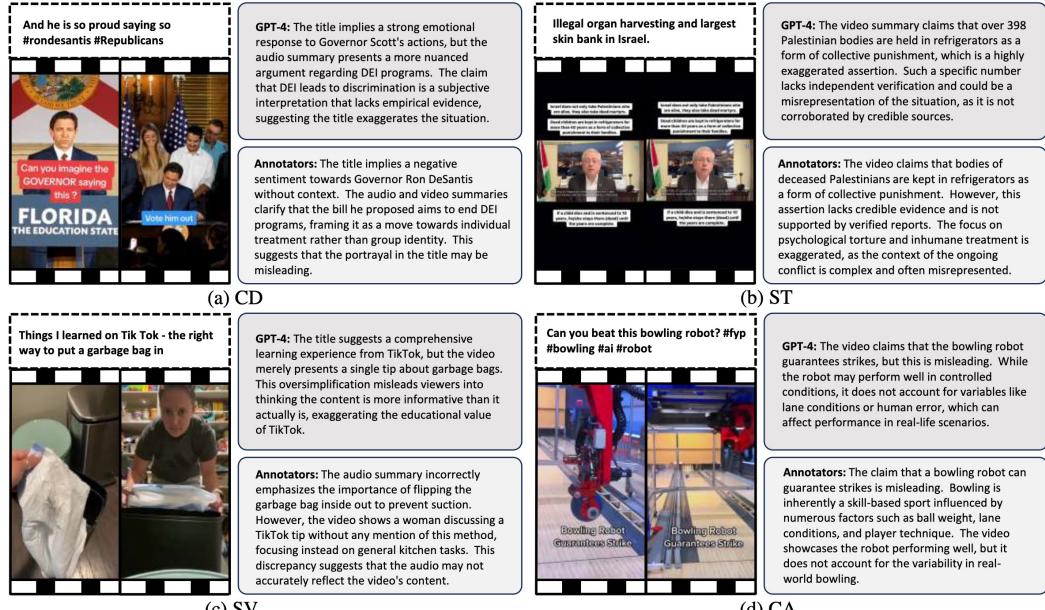
Datasets Overview

Real-world data (cont'd):

- FakeSV and FakeTT only provide binary classification labels
- The recent dataset considers natural language explanations

①FakeVE

- 2,672 samples based on FakeTT and FMNV
- Four types of explanations
 - *Contextual Dishonesty*
 - *Splice Tampering*
 - *Synthetic Voiceover*
 - *Contrived Absurdity*



Datasets Overview

Real-world data (cont'd):

- FakeSV and FakeTT only provide binary classification labels
- The recent dataset considers natural language explanations

② TRUE

- 1,097+1,828 samples from Snopes.com
- **Two types of rationales:**

- Original rationale by human
- Summary rationales by LLMs

Pioneering Explainable Video Fact-Checking with a New Dataset and Multi-role Multimodal Model Approach

 Claim	 Label						
A video shows an F-18 Super Hornet breaking the sound barrier and creating a sonic boom.							
 Video Content		 Video Information					
	Video Headline	F18 Super Hornet - Jones Beach AirShow ...					
	Video Date	25 May. 2009	Platform				
	Video Transcript	... not to exceed the speed of sound ...					
 Original Rationale							
main rationale	... the air show pilots didn't break the sound barrier.						
additional rationale 1	... evidence of the plane going supersonic, it's not.						
 Summary Rationale							
synthesized rationale	... it did not break the sound barrier as confirmed ...						
detailed reason 1	... F-18 did not exceed the speed of sound ...						
detailed reason 2	... cone is explained to be a natural phenomenon ...						
detailed reason 3	... regulations banning supersonic flight over land ...						
detailed reason 4	... sonic booms causing widespread damage ...						
 Evidences							
evidence1	The vapour cones are created by a shockwave that is						

Figure 1: A sample in the proposed TRUE Dataset. It includes the claim, video, and video background information. Besides, three types of annotations are provided: 1) label, 2) evidences, and 3) original and summary rationales.

Datasets Overview

Real-world data (cont'd):

- FakeSV and FakeTT only provide binary classification labels
- The recent dataset considers natural language explanations

③ GroundLie360

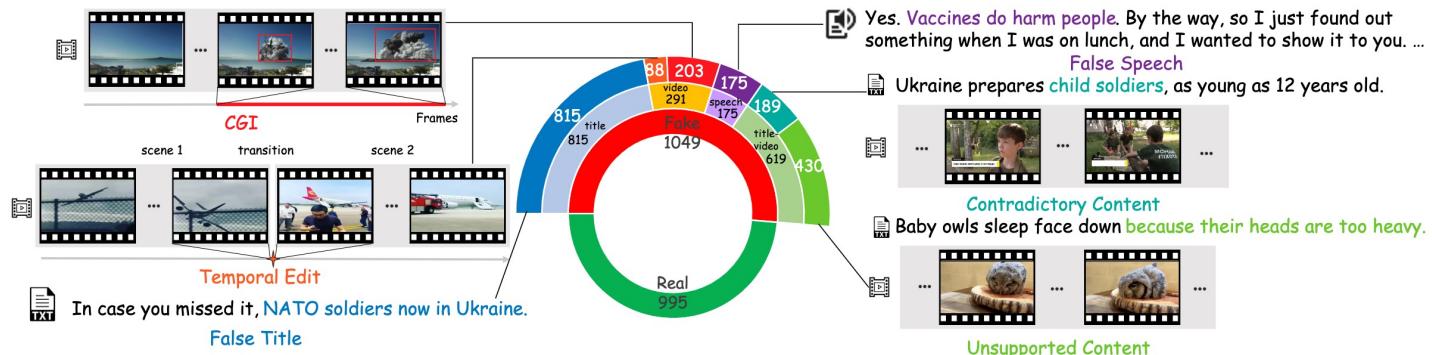


Figure 1: Overview of the GROUNDLIE360 Dataset. Our multi-modal benchmark contains 2,000+ fact-checked videos with fake type and grounding annotations. Fake types include: (1) **False Title/False Speech** - video title or spoken content containing demonstrably false claims; (2) **Temporal Edit** - videos altered to distort event chronologies or fabricate deceptive narratives; (3) **CGI** - digitally manipulated or generated synthetic media; (4) **Contradictory Content** - text-video semantic mismatches; and (5) **Unsupported Content** - headlines lacking evidentiary support in video content. The dataset offers a unified benchmark for fake content classification and localization.

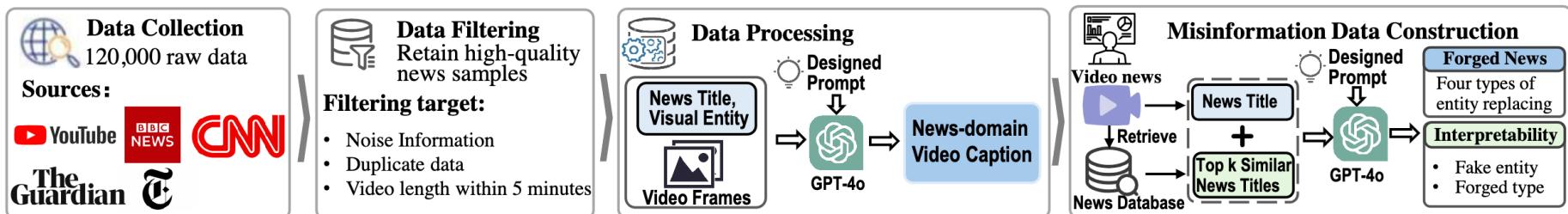
Datasets Overview

Synthesized data

- Starting from the real news samples, these datasets modify the semantic descriptions to construct fake news samples
- Of larger scale than real-world ones due to the convenience of (M)LLM-assisted data generation

① FakeVV

- 100k samples from 2006 to 2025
- Substitute with representation-similar entities based on the retrieved results
- Input Top3 samples to GPT-4o for generation



Datasets Overview

Synthesized data

- Starting from the real news samples, these datasets modify the semantic descriptions to construct fake news samples
- Of larger scale than real-world ones due to the continence of (M)LLM-assisted data generation

② Official-NV

- Use 5000 news samples from authoritative sources
- Modify the text description by four strategies

Original Text	Modified Text
China's booming tea industry imbued with new momentum	China's tea industry surges forward with rejuvenated vitality (TT)
The stunning many-coloured landscapes of Xinjiang	The stunning many-coloured landscapes of Anhui (FT in position)
Palestinian death toll from Israeli attacks in Gaza, West Bank nears 20,000	Palestinian death toll from Israeli attacks in Gaza, West Bank more than 30,000 (FT in quantity)
China seeks to build world's largest national park system	China aims to dismantle extensive national park network (FT in action)
Ready... set... GO! This cat sure knows how to win a sprint race	Ready... set... GO! This dog sure knows how to win a sprint race (FT in object)

Tutorial Outline

3. Characterization of Misinformation Videos

Definition

Datasets Overview

Analysis on FakeSV

News Content

Social Context

Propagation

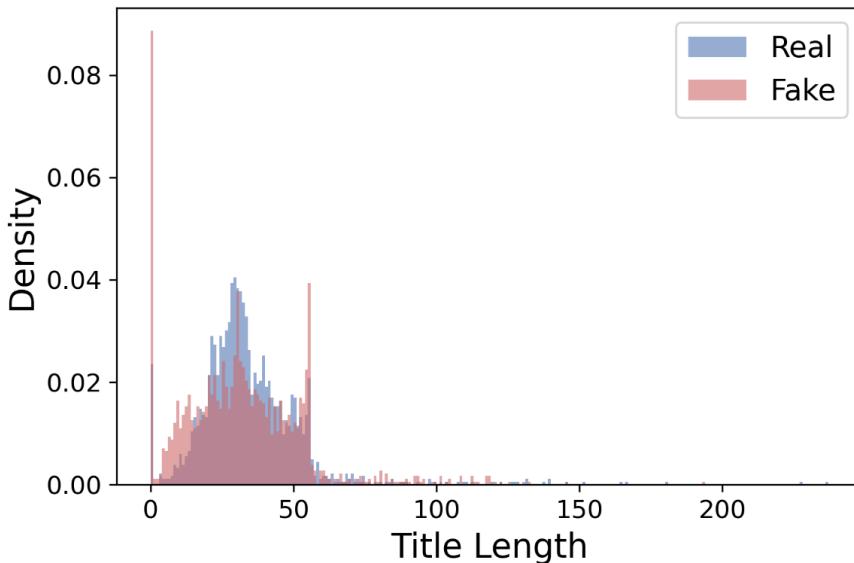
Q+A/Discussion

Analysis on FakeSV

- Take FakeSV as an example. This part shows the statistical difference between real and fake news videos.

1. News Content -> Text Length

Fake news videos have shorter and more empty titles, providing less information compared with real news.

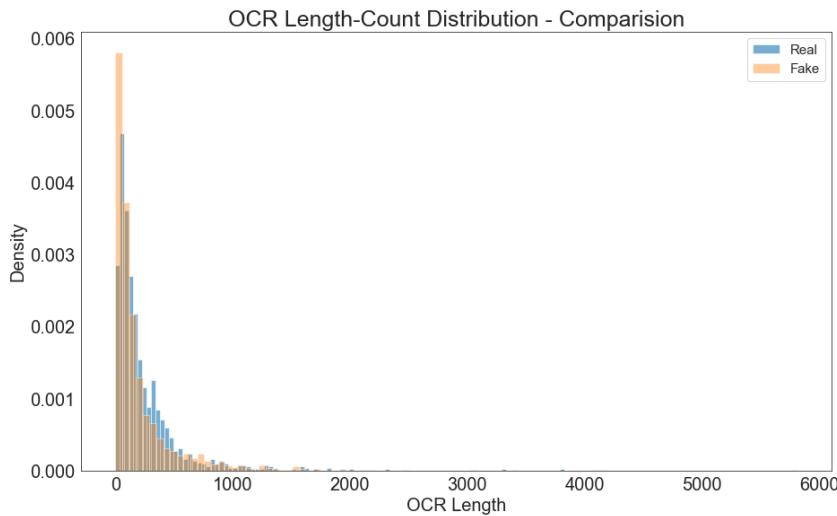


Analysis on FakeSV

- Take FakeSV as an example. This part shows the statistical difference between real and fake news videos.

1. News Content -> Text Length

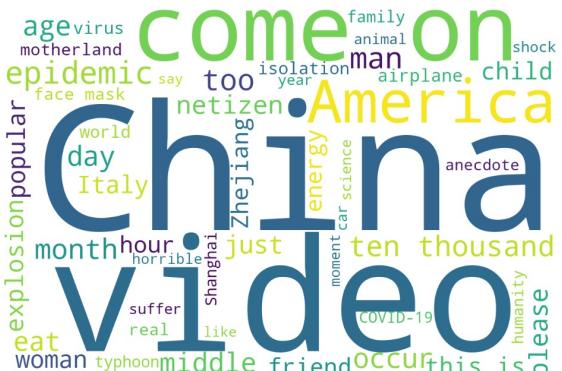
For on-screen texts extracted by OCR tools, fake and real news videos show similar distribution, with real news videos more likely to have longer on-screen texts, showing its informativeness.



Analysis on FakeSV

1. News Content → Emotion and Words

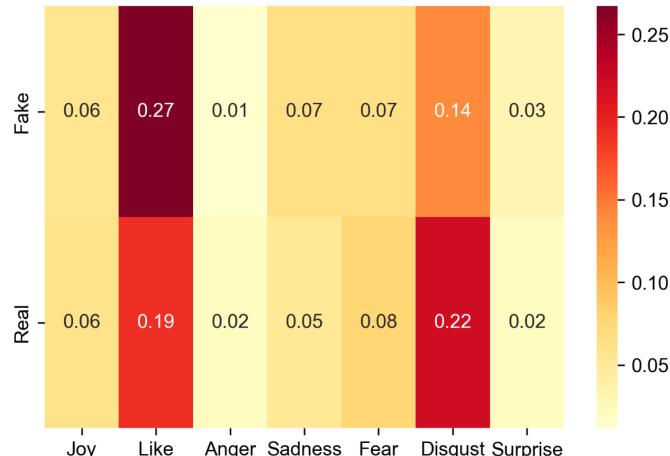
- Fake news titles emphasize the word “video” much, **prefer emotional and spoken words**, and cover diverse topics. Real news videos use **more journalese** and focus more on accidents and disasters.
 - Based on Affective Lexicon Ontology, we find that fake news titles show **more like** while real news titles show **more disgust**.



(a) Fake



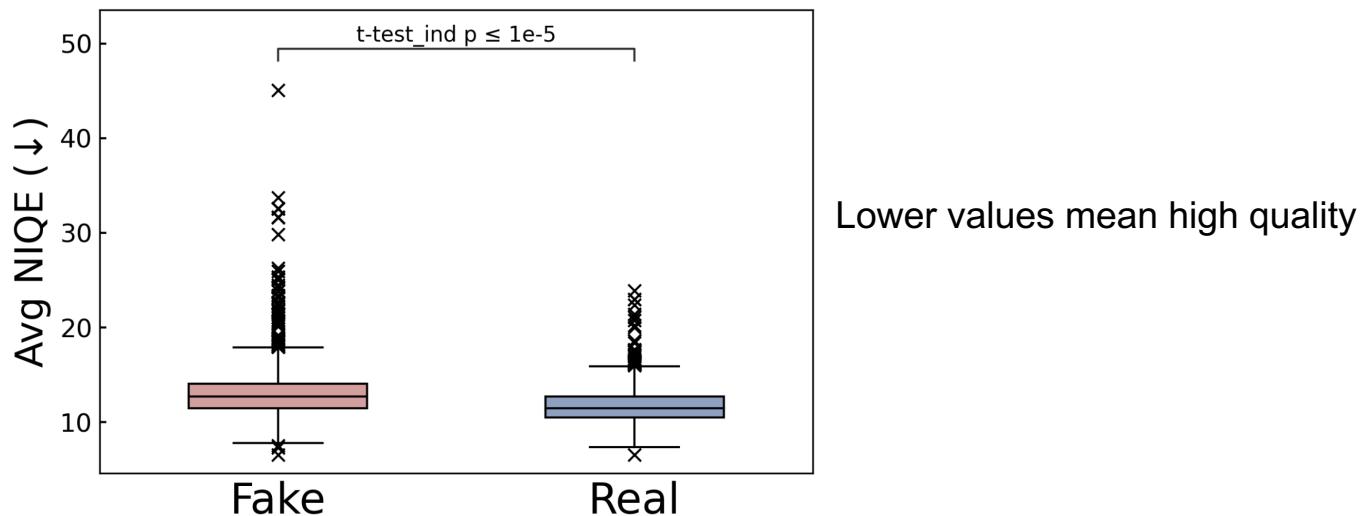
(b) Real



Analysis on FakeSV

1. News Content -> Video

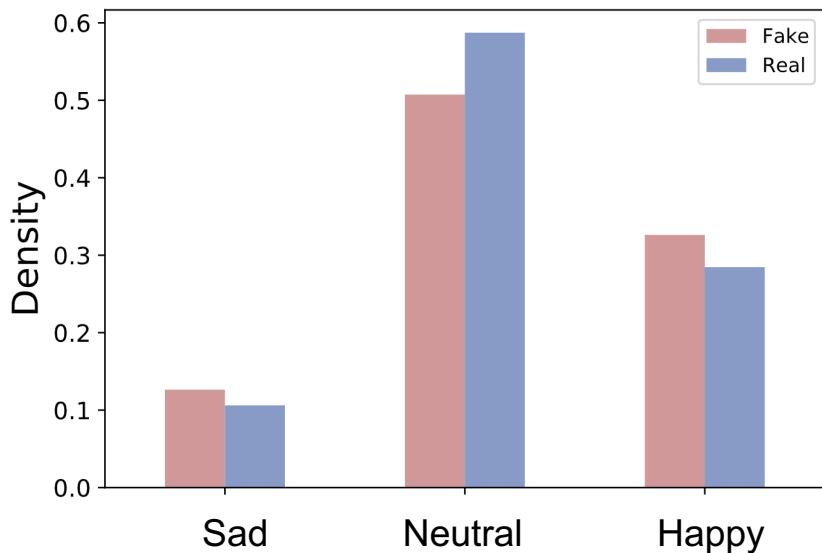
- By employing NIQE on video frames to indirectly measure video quality, we see that **fake news videos have lower quality than real news and contain videos with particularly poor quality**
- This is because the materials are often from unprofessional devices or simply old. AI-generated ones may go to another end: **It might be too clear to be unrealistic**



Analysis on FakeSV

1. News Content -> Audio

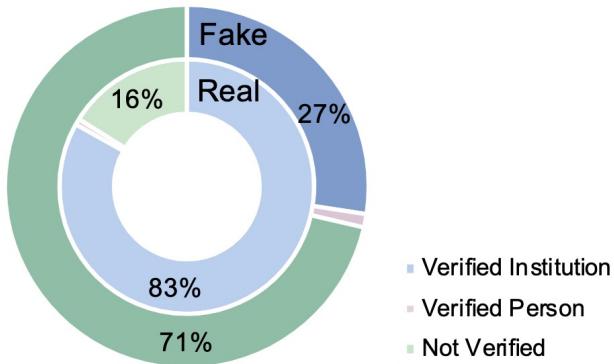
- We analyze the speech emotion by the pre-trained wav2vec model
- **Speech in fake news videos shows more obvious emotional preferences than real news**



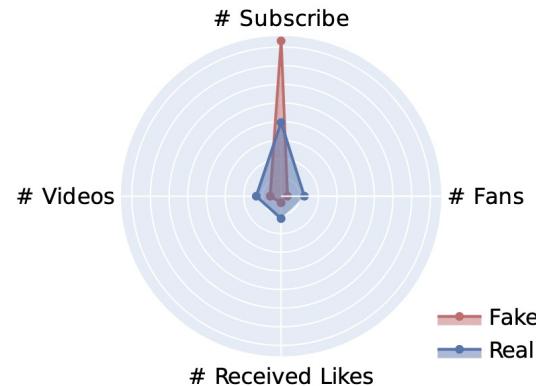
Analysis on FakeSV

2. Social Context-> Publisher Profiles

- The distribution is similar to what have been found on conventional social media like Weibo.
- **Most publishers of real news are verified accounts while most fake news publishers are not**
- **Fake news publishers have more “consuming” behaviors (subscribes) and less “creating” behaviors (published videos, received likes, and fans) than real news publishers**



(a) Authority

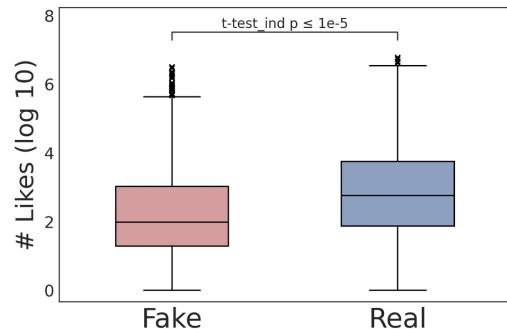


(b) Statistics

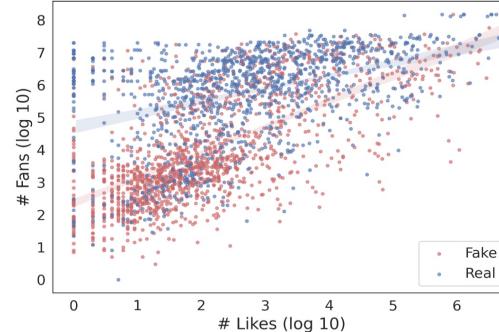
Analysis on FakeSV

2. Social Context-> User Responses

- Real news videos **receive more likes** than fake news, which is intuitive considering that real news publishers have more followers.
- **However, fake news videos receive more likes than real news when their publishers have a similar number of followers, which illustrates that fake news videos are more attractive than real news**



(a) Number of likes

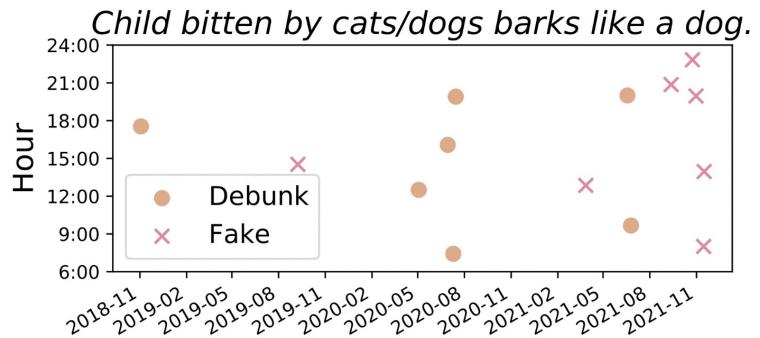
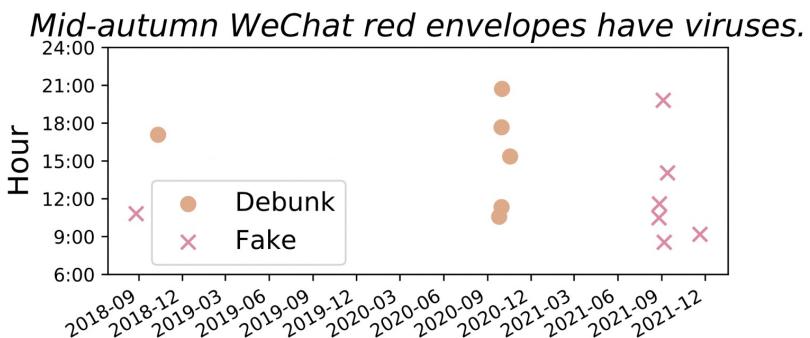


(b) Relationship between the number of publisher fans and likes.

Analysis on FakeSV

3. Propagation-> Temporal Distribution

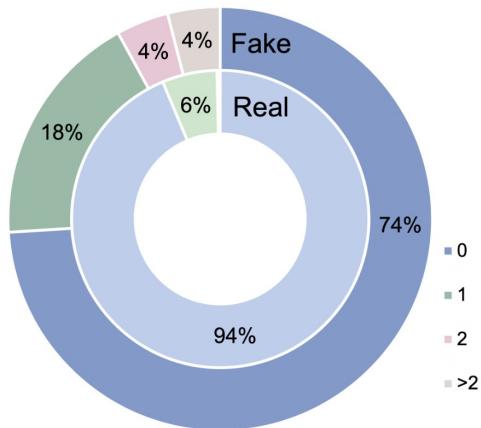
- Fake news that has been previously debunked can still spread
- **For 434 events with debunking videos, 39% of them have fake news videos emerged after the debunking videos were posted, especially the current or long-standing hot event**



Analysis on FakeSV

3. Propagation-> Video duplication

- With the convenience of video editing functions provided by these platforms, people tend to edit and re-upload the videos, usually with no mention of the source.
- **Using pHash on the video covers, we find that fake news videos have higher repetition while real news videos are more diverse.** This is due to the fact that real events will receive various images/videos from different witnesses and sources.



Detection Part I: Human-Edited Misinformation

Section 4

Tutorial Outline

Detection Part I: Human-Edited Misinformation

Signal-based detection

Editing Traces

Generation Traces

Semantic-based detection

Seek clues within the sample's own content

Seek clues from external information

Intent-based detection

Social context

Clue integration for misinformation video detection

Parallel Integration

Sequential Integration

Q+A/Discussion

Signal-based detection

Misinformation videos often contain manipulated or generated video and audio content in which the forgery procedure often leads to traces in underlying digital signals.

Editing

Alterations on existing data of video and audio modality



Generation

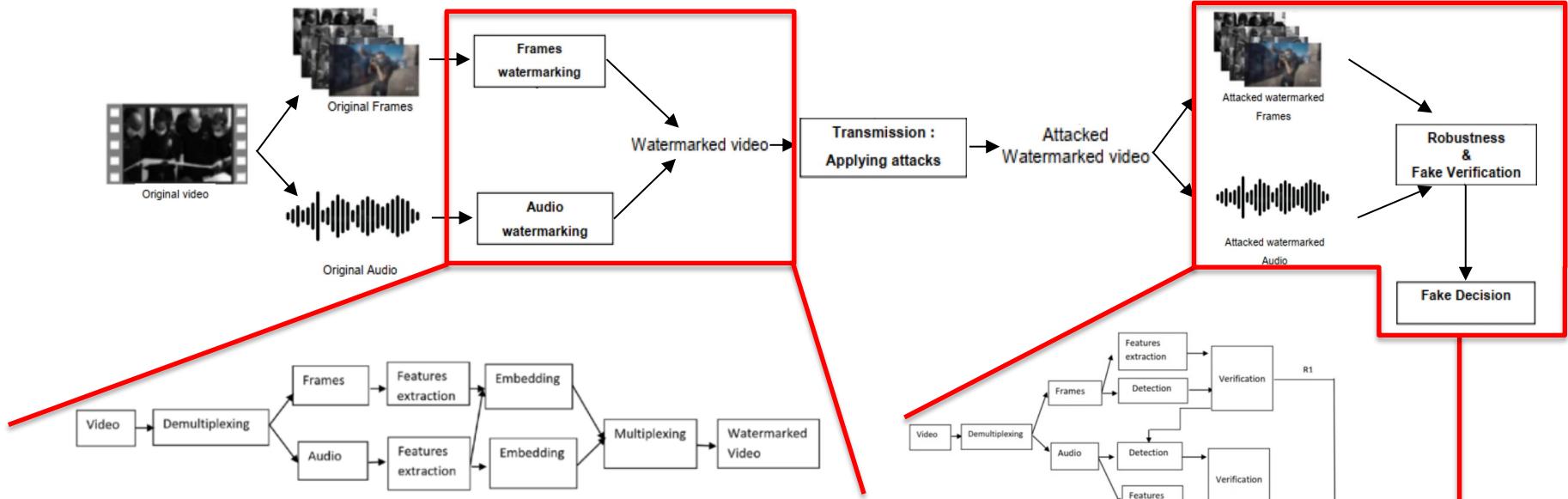
Directly generate complete vivid videos with forged human faces or voices



Editing Traces

Active detection

Pre-embed and extract digital watermarks for detection

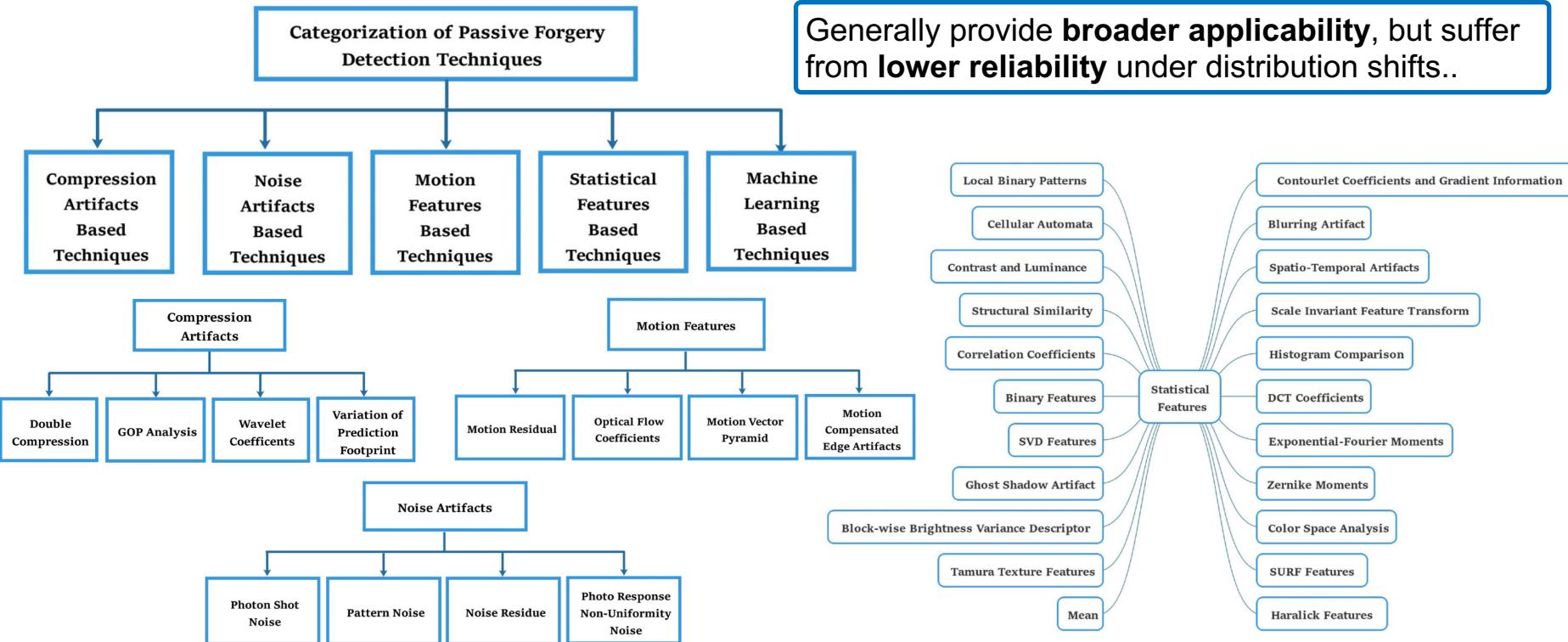


Generally provide **quicker responses** and **more accurate judgments**, but suffers from coverage gaps.

Editing Traces

Passive detection

Use the characteristics of the digital video itself for detection



Generation traces

Mining spatial-temporal and spectral-prosodic traces from video and audio for detection

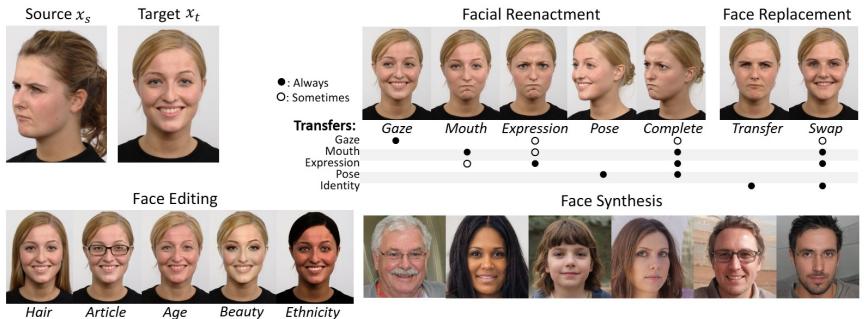


Fig 1. Visual deepfakes examples.

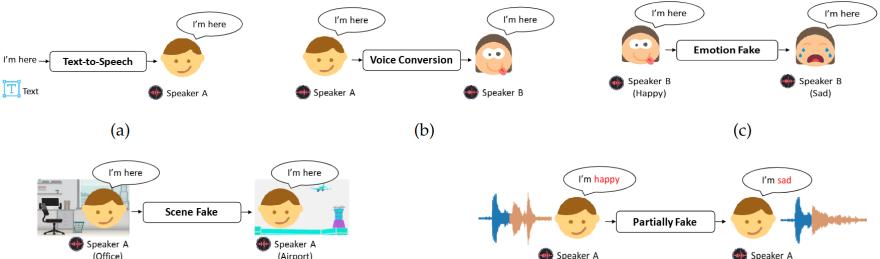


Fig 2. Deepfake audio types.

- Deepfake Video Detection Clues:
 - Boundary artifacts between fake face & background
 - Inconsistent lighting, warping, background
 - Generator or sensor “fingerprints”
 - Unrealistic motion or emotion patterns
 - Missing biological signals
 - Temporal inconsistency
 - Deep features learned from end-to-end models
 - Short-term spectral features
 - Phase / group-delay features
 - Long-term spectral features
 - Prosodic features
 - intrinsic (a)synchronization between video and audio frames
 -

Tutorial Outline

Detection Part I: Human-Edited Misinformation

Signal-based detection

Editing Traces

Generation Traces

Semantic-based detection

Seek clues within the sample's own content

Seek clues from external information

Intent-based detection

Social context

Clue integration for misinformation video detection

Parallel Integration

Sequential Integration

Q+A/Discussion

Semantic-based detection

The falsehood is conveyed through **incorrect semantic changes against the truth**.

 Even if **editing or generation traces** are detected, it does **not necessarily** mean that the video **conveys misinformation**.

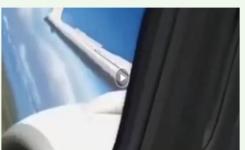
A video that is technically untampered can be employed in a deceptive manner:

- Fact Distortion
- Misleading Substitution
- Groundless Fabrication
-

 In the UGC era, most misinformation arises from such semantic manipulations

Original Real News

> Simulation Video of the Ethiopian Airlines Flight ET302 crash incident

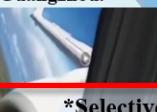


Fact Distortion > The black box from the crashed Ethiopian Airlines plane has been found. This is the last 10 seconds recorded by the automatic recorder... the screams of people on the brink of death. Being alive is the greatest luck.

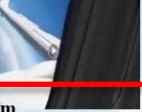


Misleading Substitution

> A China Eastern Boeing 737 aircraft lost contact and crashed over Wuzhou, Guangxi, while operating a flight from Kunming to Guangzhou.



Groundless Fabrication: Lion Air pilot, facing financial problems, sent away the co-pilot and crashed the plane. Footage before the crash revealed!



Selective Editing > Footage from the Ethiopian Airlines Crash Site



* The Ethiopian Airlines Flight ET302 crash occurred on March 11, 2019, local time.

Seek clues within the sample's own content

Early works most leverage textual information for misinformation detection

Textual information:

- Video description
- Title
- Subtitles
- Transcriptions
-

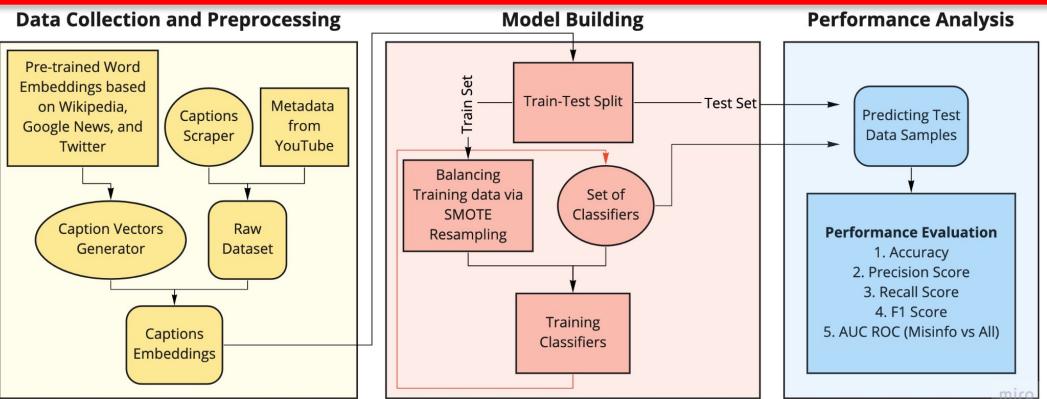
Hand-crafted features:

- Basic statistical attributes
- Specific expressions
- Corpus-aware features(N-grams, TF-IDF, LIWC)

Continuous representation features:

- Pre-trained static embeddings
- Task-specific encoders

From video description 05: text length 06: number of words 07–08: contains question/exclamation mark (Boolean) 09–10: contains 1st/3rd person pronoun (Boolean) 11: number of uppercase characters 12–13: number of positive/negative sentiment words 14: number of slang words 15: has ":" symbol (Boolean) 16–17: number of question/exclamation marks	Title&Transcript	text length the number of words contains question/exclamation mark (Boolean) contains 1st/3rd person pronoun (Boolean) the number of positive/negative sentiment words has ":" symbol (Boolean) the number of question/exclamation marks has clickbait phrase (Boolean) sentiment polarity the number of modal particles the number of personal pronouns tf-idf Ngrams LIWC
---	------------------	--



Seek clues within the sample's own content

Aggregate heterogeneous information across different modalities (beyond textual content)

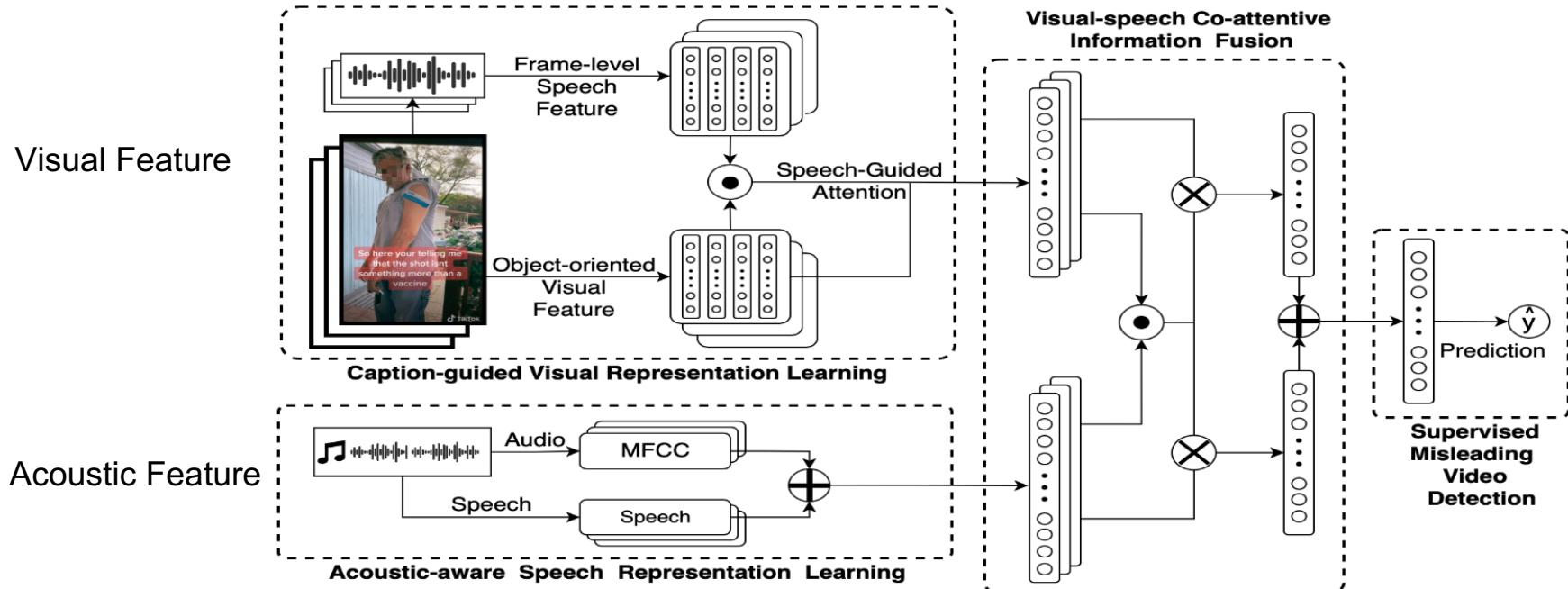
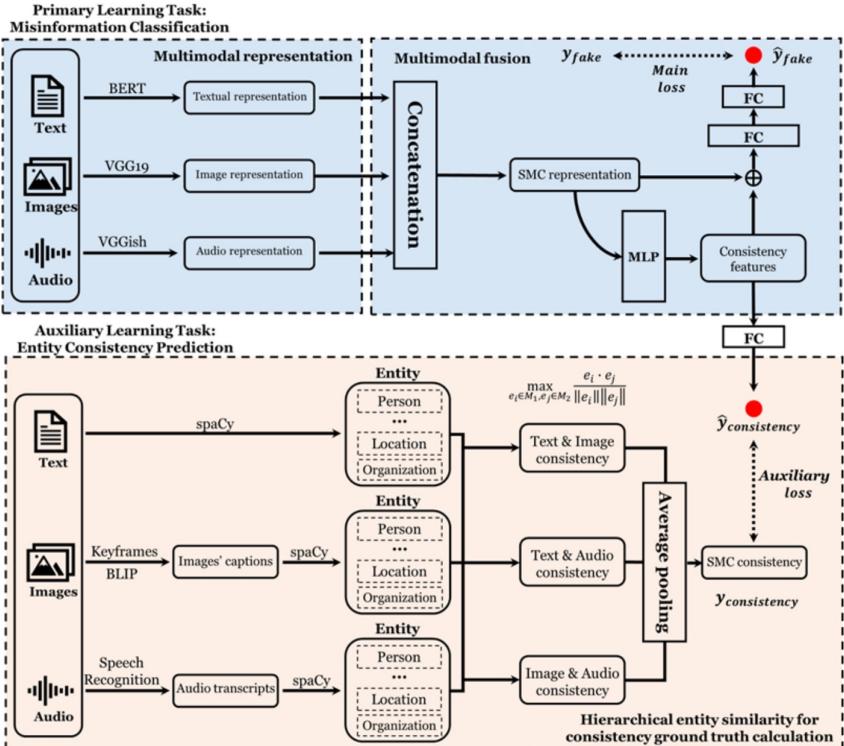
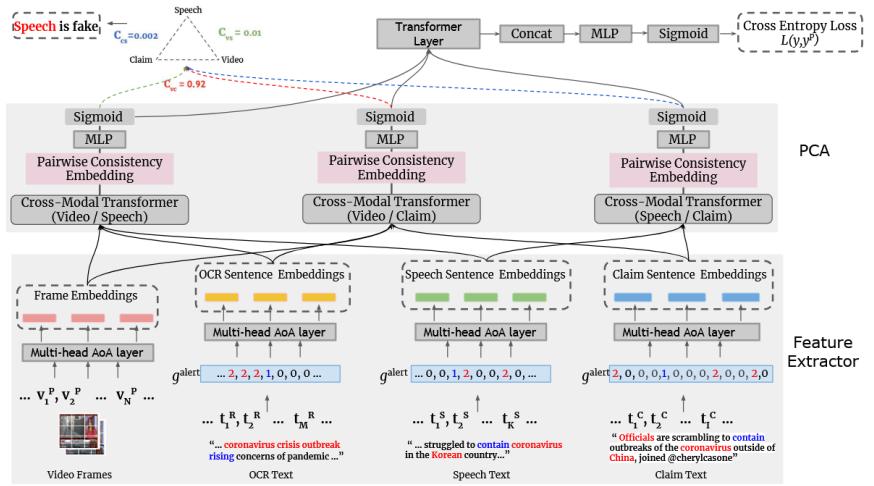


Figure 3: Overview of the TikTec Framework

Seek clues within the sample's own content

Leverage cross-modal correlation: Find mismatches between modalities (video-text-audio)



👉 Utilize **embedding consistency** across three different modalities for detection

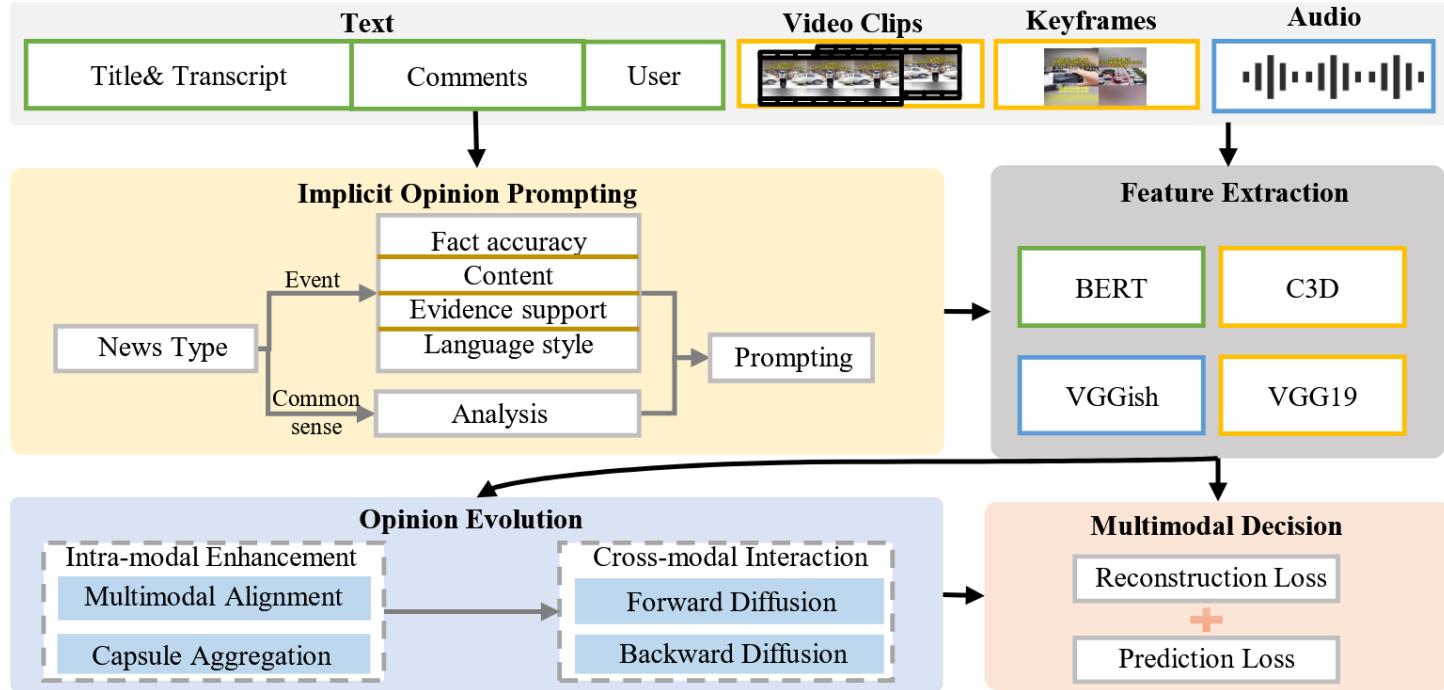
👉 Utilize **entity consistency** across three different modalities for detection

Seek clues within the sample's own content

Deepen cross-modal correlation clues mining: Uncover implicit opinions across modalities

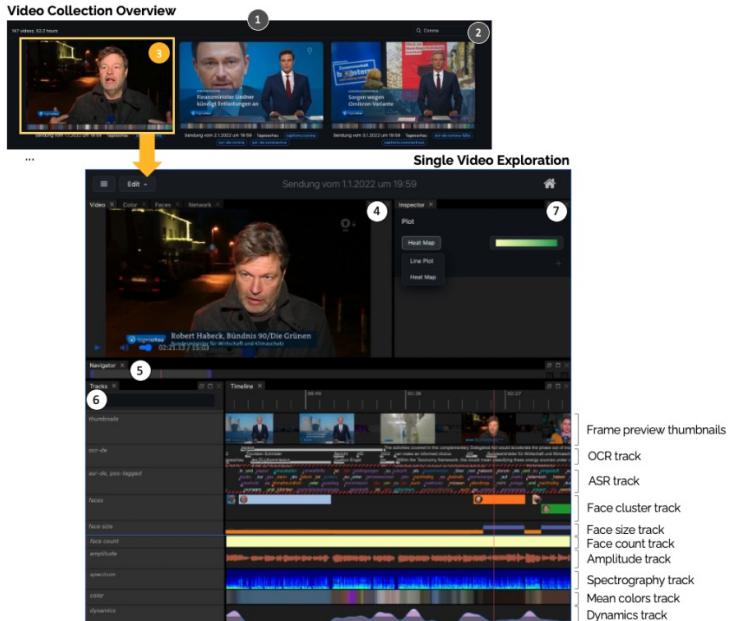
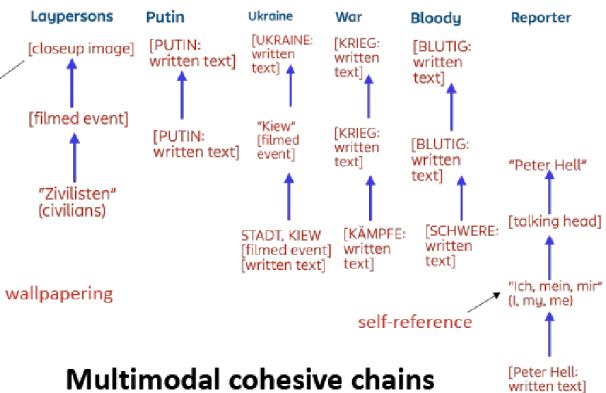
💡 Each modality carries an **implicit credibility attitude**; without cross-modal propagation, **local** deceptive cues fail to influence the **global** decision.

👉 Simulate mutual reinforcement among modal opinions through a **diffusion process**



Seek clues within the sample's own content

Another analysis perspective: News as Narrative



💡 Two pivotal aspects of Narrative Theory: analyzing the “what” (the content of the story) and the “how” (the strategy of storytelling)

👉 Disinformation TV news videos has distinguish narratives.

Seek clues within the sample's own content

Enhance misinformation video detection by analyzing the narrative creation process.

👉 Analyze the creative process behind misinformation videos:

Phase I – Material Selection:
 Fake news exhibits **emotional bias** and **semantic selectivity** when choosing video materials.

Phase II – Material Editing:
 When **spatially** imposing text and **temporally** splicing materials, fake news tend to adopt a relatively **simple arrangement**.

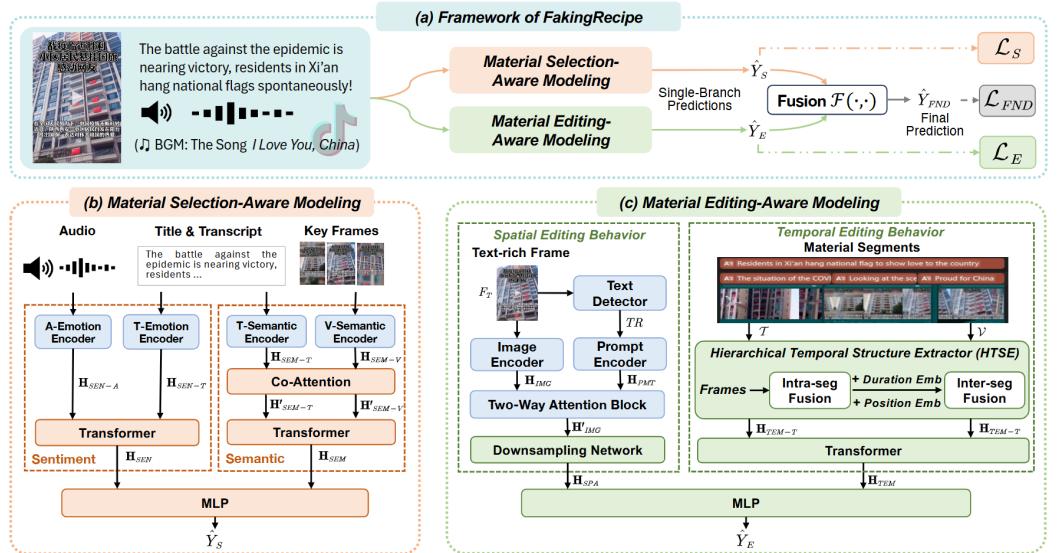
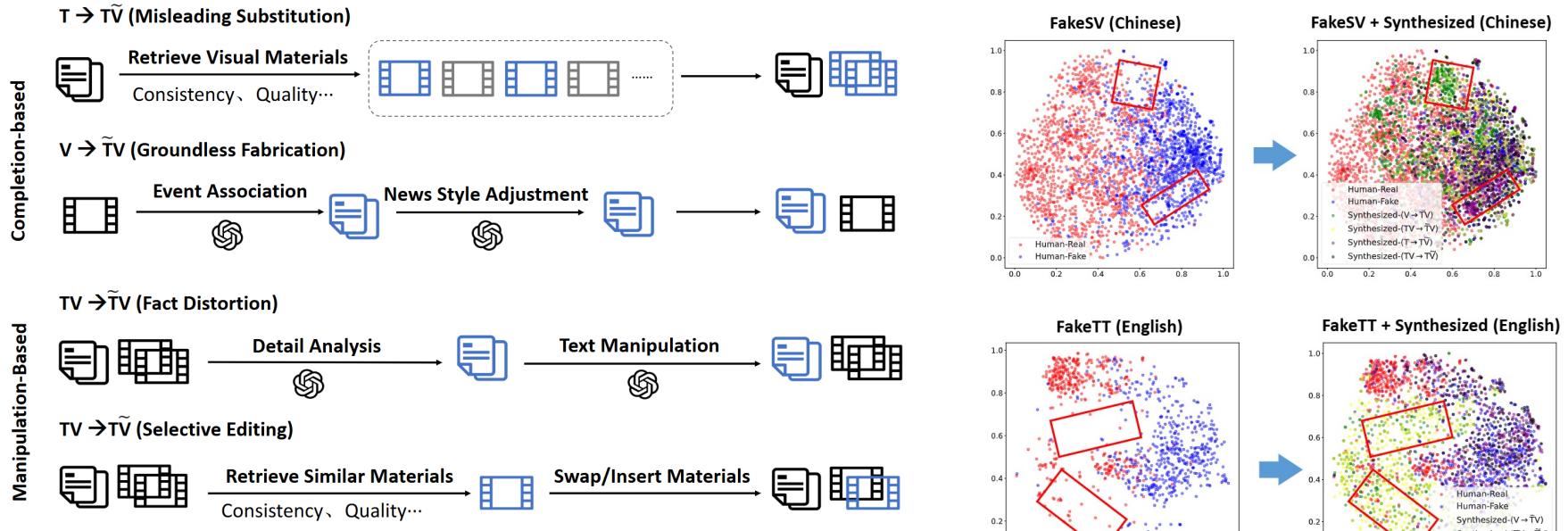


Figure 6: Overview of the proposed FakingRecipe model. (a) Overall framework: The news video is processed through dual perspectives, with a late fusion strategy employed to integrate clues for final prediction. (b) Material Selection-Aware Modeling (MSAM) module: Extracts clues from both sentimental and semantic aspects. (c) Material Editing-Aware Modeling (MEAM) module: Extracts clues based on spatial and temporal aspects. $\mathcal{F}(\cdot, \cdot)$ denotes the fusion function. The parameters in the modules in blue are frozen and others are trainable. The overall model is trained under the supervision of the loss functions \mathcal{L}_S , \mathcal{L}_{FND} , and \mathcal{L}_E .

Seek clues within the sample's own content

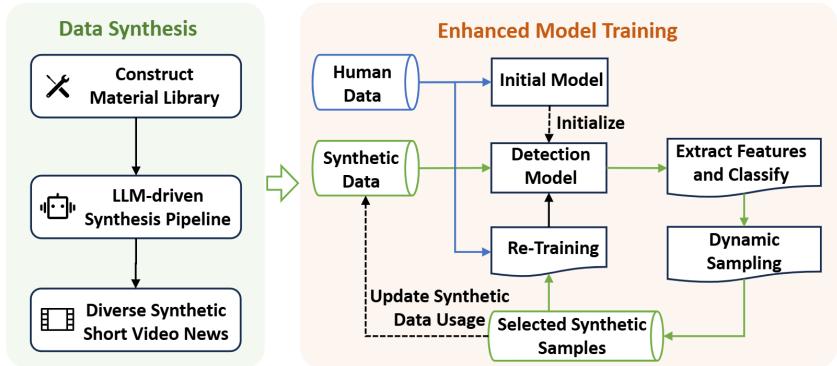
Enhance misinfo video detection via data augmentation: Simulating typical creative processes



👉 The synthesized samples not only align with human-crafted fakes but also enrich underrepresented regions of the feature space

Seek clues within the sample's own content

Enhance misinfo video detection via data augmentation: Simulating typical creative processes



Model	FakeSV				FakeTT			
	Acc	F1	Prec.	Rec.	Acc	F1	Prec.	Rec.
MMVD	75.83	75.33	75.51	75.21	67.50	66.20	66.51	68.43
w/AgentAugRAN	73.80	73.22	73.33	73.13	63.57	61.69	61.84	63.11
w/AgentAugBAL	76.01	75.58	75.86	75.45	66.43	65.32	65.55	66.94
w/AgentAugAL	77.12	76.69	76.93	76.55	68.57	67.26	67.56	69.71
FANVM	78.41	77.89	78.25	77.70	71.57	70.21	70.21	72.63
w/AgentAugRAN	78.78	77.63	80.08	77.16	67.22	66.77	69.18	71.42
w/AgentAugBAL	79.34	78.31	80.42	77.84	68.23	67.99	71.57	73.70
w/AgentAugAL	81.37	80.42	82.73	79.88	75.25	74.02	73.84	76.65
SVFEND	80.88	80.54	80.83	80.51	77.14	75.63	75.12	77.56
w/AgentAugRAN	82.66	81.94	83.54	81.44	79.72	77.72	77.33	78.24
w/AgentAugBAL	81.92	80.85	83.96	80.23	77.58	76.00	75.43	77.40
w/AgentAugAL	83.76	82.98	85.20	82.38	80.43	78.61	78.12	79.29
SVRPM	81.34	81.11	81.38	80.97	81.79	79.42	79.67	79.19
w/AgentAugRAN	79.57	79.33	79.61	79.20	80.00	77.34	77.49	77.19
w/AgentAugBAL	81.73	81.58	82.06	81.49	82.14	80.28	81.36	79.58
w/AgentAugAL	83.10	82.89	83.16	82.70	82.86	80.57	80.75	80.41

👉 By integrating **active learning** to select potentially useful augmented samples,
 👉 the framework **consistently boosts** short-video fake news detection performance.

Seek clues within the sample's own content

More robust detection: Consider incomplete modality conditions

Text: Heavy rainfall in Zhengzhou caused water to flood into subway stations. The water level inside subway cars on Line 5 reached over one meter. The subway station was temporarily closed, and operations were suspended.

Image/Audio/Video: Missing (Resource Expiration)

Pattern 1: (T, MI, MA, MV)

(a) Resource Expiration



Text/Image: Incomplete(Blank Contents)

Audio: 

Pattern 2: (MT, MI, A, V)

(b) Blank Contents



Text: When a surge occurred on the Ganjiang River embankment in Nanchang, a commando team of 15 party members was the first to jump into block the surge.
Audio/Video: Missing(Audio Corruption)

Pattern 3: (T, MA, I, V)

(c) Audio Corruption



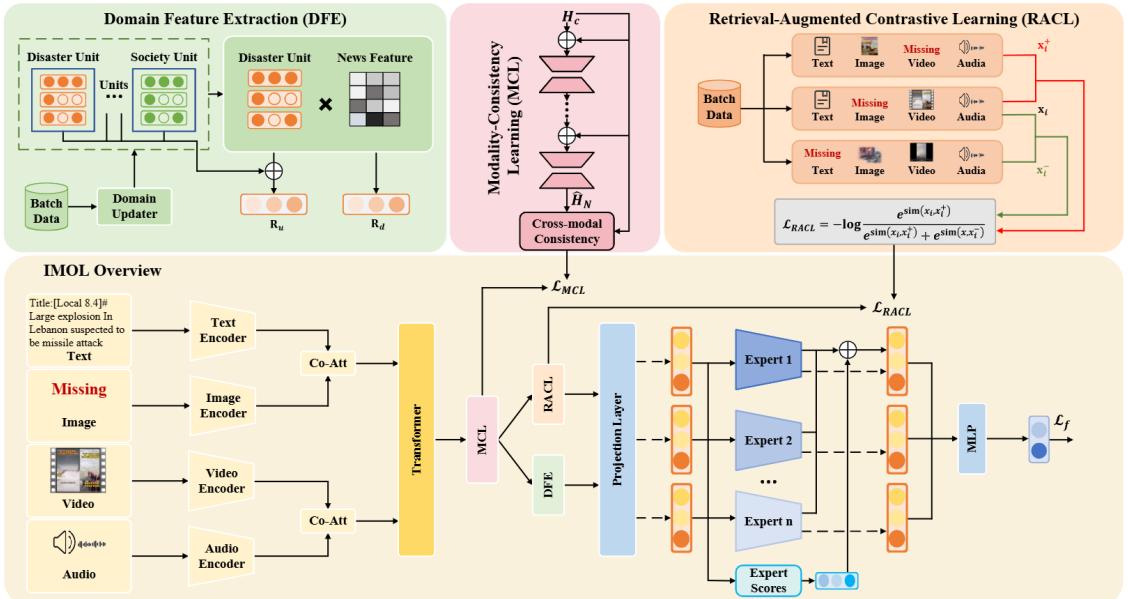
Text: Yiwu faces strong winds and heavy rain, leaving San Ting Road night market in a devastating state.
Audio/Video: Missing(Contents Damage)

Pattern 4: (T, I, MA, MV)

(d) Contents Damage

Figure 1: Examples of modality patterns in news videos. Each news piece consists of four modalities: Video (V), Audio (A), Text (T), and Image (I). However, real-world news is often *modality-incomplete*, where one or more modalities may be missing (MV: Missing Video, MA: Missing Audio, MT: Missing Text, MI: Missing Image).

Jointly modeling cross-modal reconstruction and cross-sample reasoning → robust and generalizable fake news video detection.



Seek clues within the sample's own content

More robust detection: Mitigate modality bias & Adaptively select reliable feature

👉 Debiases static, dynamic, and social views through causal and counterfactual reasoning

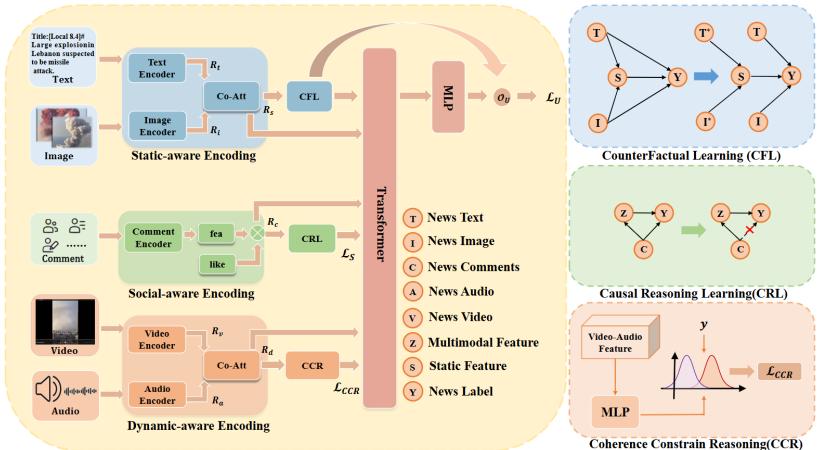


Figure 4: Overview of proposed Multimodal Multi-View Debiasing framework. The CFL, CCR and CRL mitigate the static, dynamic and social biases during multimodal fusion, respectively. Then the MMVD is learned to determine whether the news video is fake or not.

👉 Adaptively trusts the most reliable modalities and selecting the according modality experts

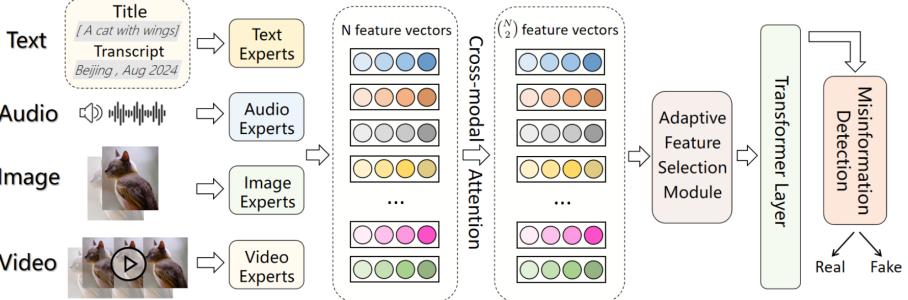


Figure 1: Architecture of the multimodal misinformation detection framework MisD-MoE

Seek clues from external information

Integrate the neighborhood relationship of new videos belonging to the same event

 News videos from different perspectives regarding the same event contain complementary or contradictory information

Utilize debunking videos to rectify false negative predictions

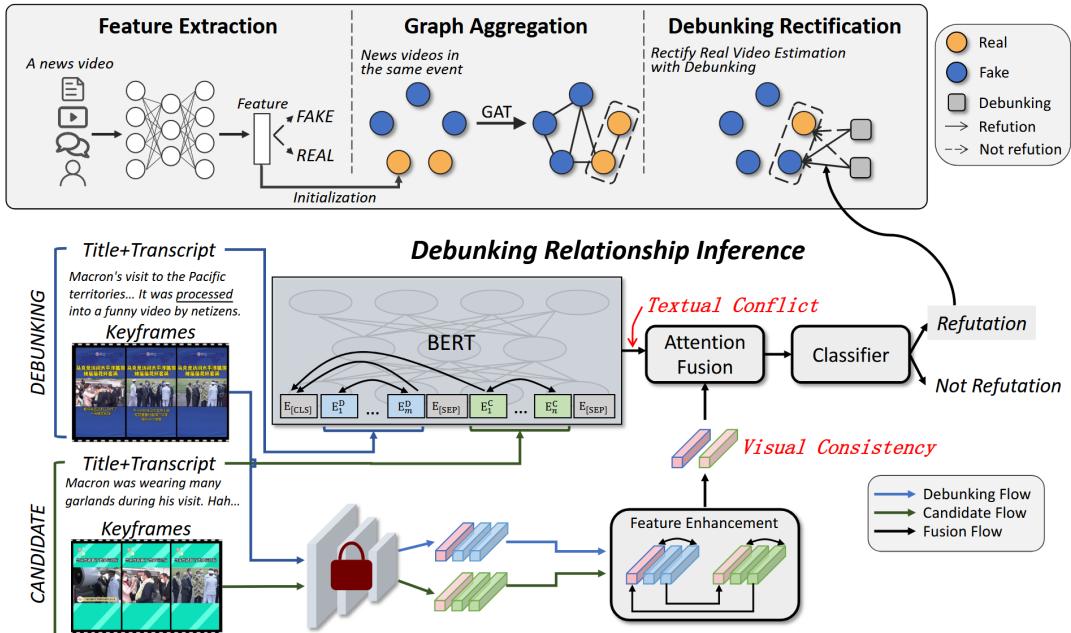


Figure 2: Architecture of the proposed framework NEED. The first row indicates the three stages in NEED, including feature extraction, graph aggregation, and debunking rectification. To realize the debunking rectification, debunking relation inference (the second row) is introduced to determine the refutation relationship.

Seek clues from external information

Model latent social and cascade relationships for fake news detection

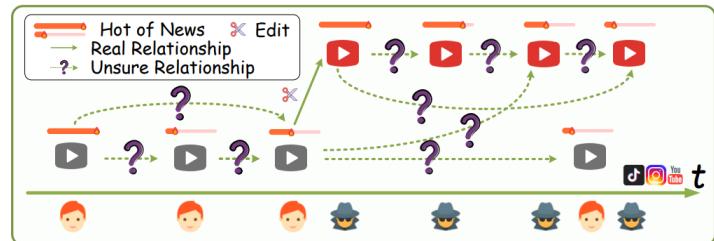


Fig 1. News dissemination relationship is unsure in TikTok, Instagram, and YouTube

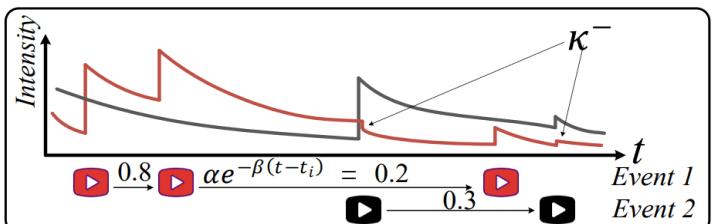
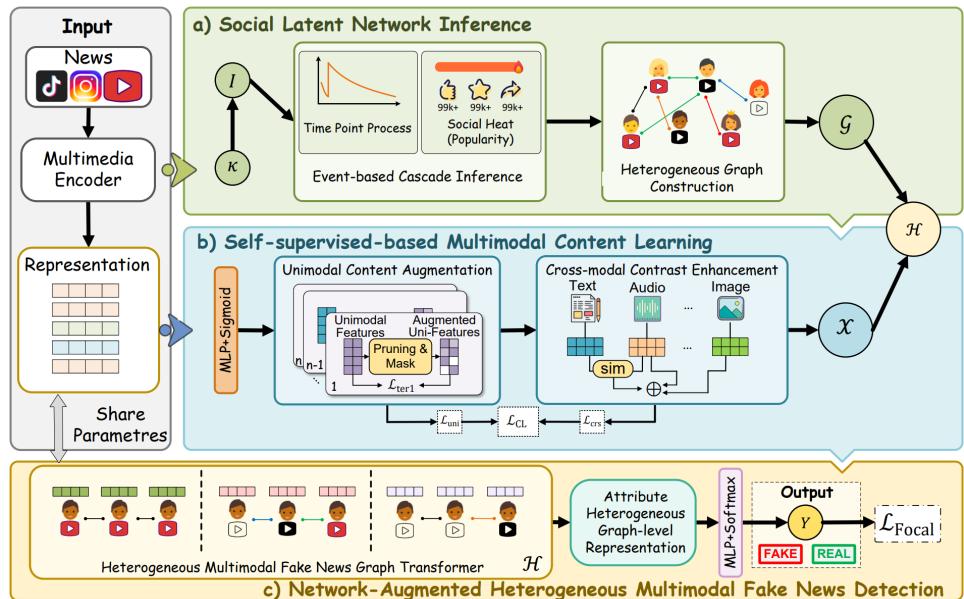


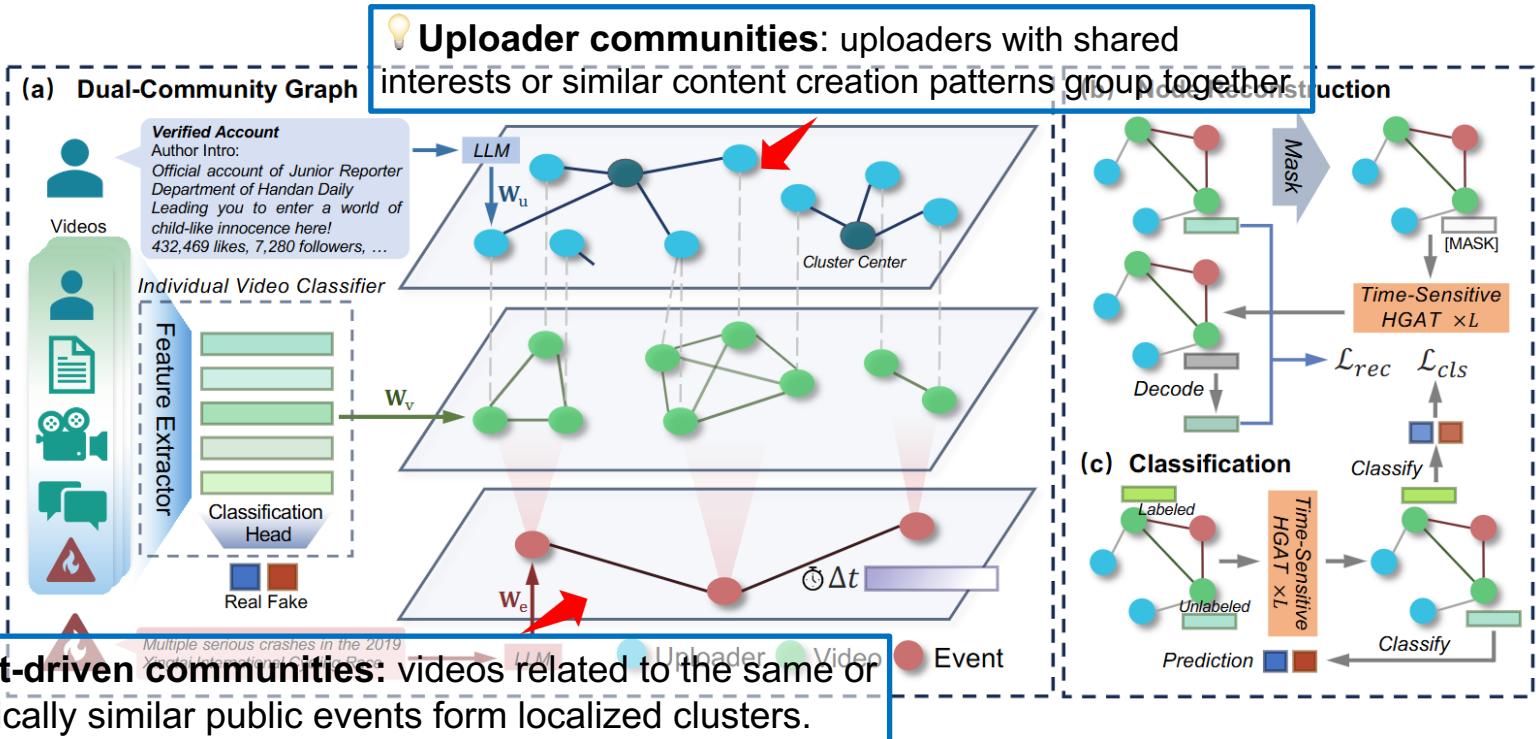
Fig 2. Illustrate of event-based cascade influence.



👉 Infer latent social cascades via event-based temporal modeling and heterogeneous graph construction

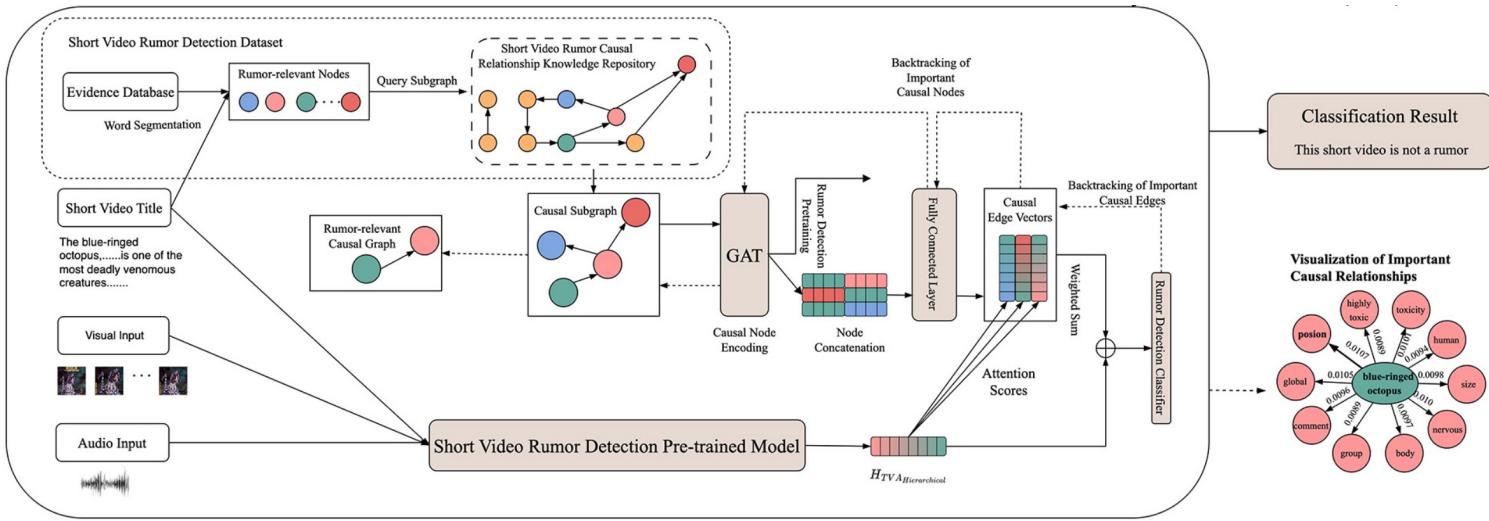
Seek clues from external information

Integrates dual-community patterns: **Uploaders'** and **Event-driven** communities



Seek clues from external information

Introduce external knowledge through **causal relationship graphs**



 Construct causal relationships between entities and integrates the causal subgraphs for the interpretation of knowledge distortion

Seek clues from external information

Introduce external knowledge through Large Language Model

 LLMs serve as flexible knowledge bridges, transforming external information into structured reasoning evidence for multimodal verification.

 Design multi-stage pipelines where LLMs refine content, retrieve domain knowledge, and reason to generate verifiable explanations.

 Build multi-role LLM frameworks that decompose fact-checking into retrieval, verification, and reasoning stages.

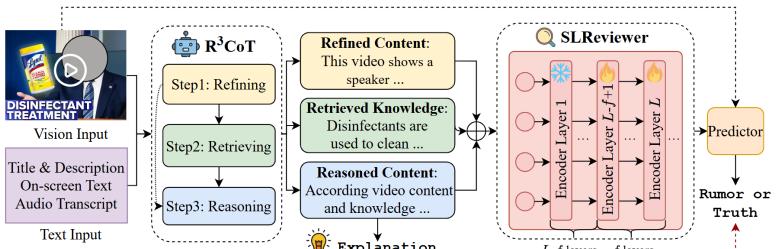


Figure 2: The structure of our proposed ExMRD framework. (1) The R³CoT process prompts MLLMs to refine the video content, retrieve domain knowledge, and reason to provide explanations. (2) The SLReviewer is to distill the explainable evidence from R³CoT to facilitate reliable rumor detection.

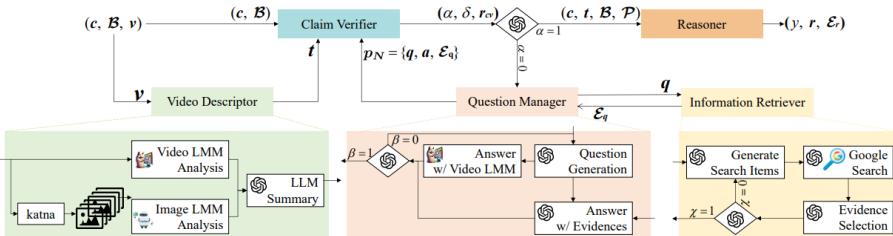
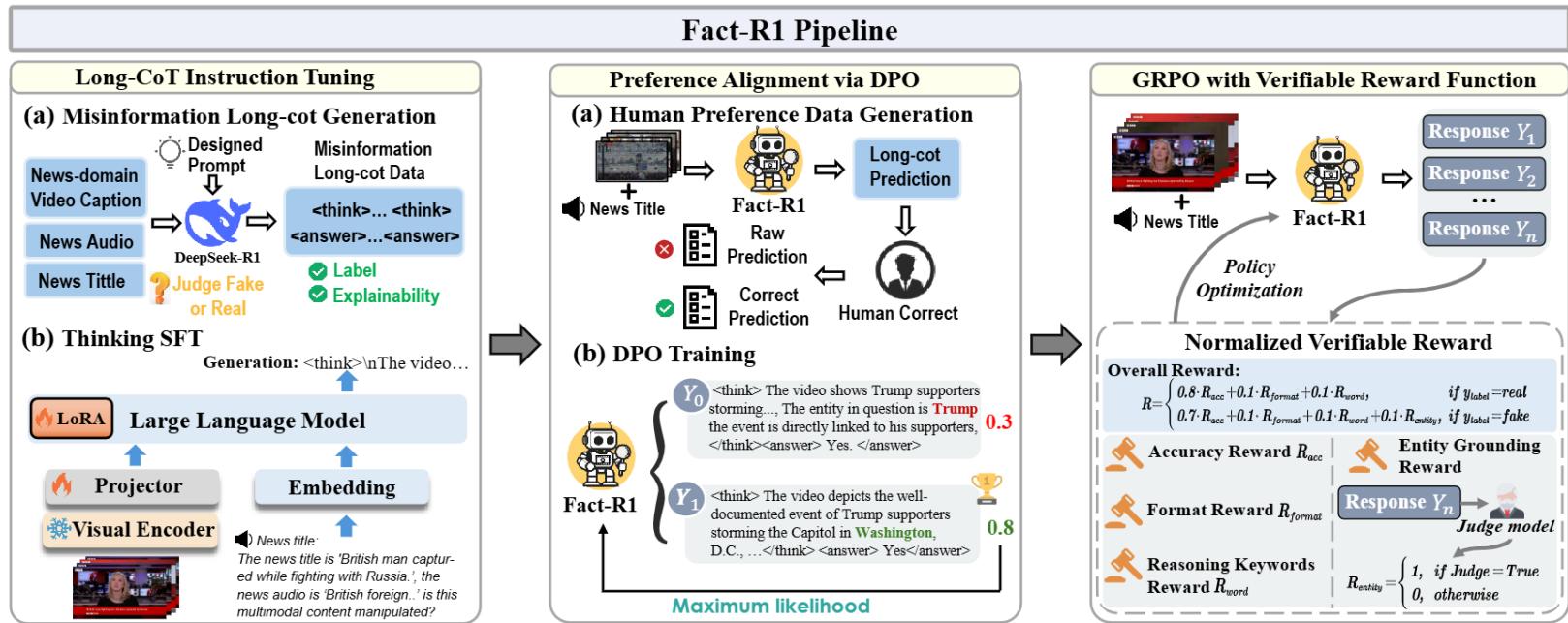


Figure 4: Overview of the proposed 3MFact framework, comprising five components: Video Descriptor (video-to-text conversion), Claim Verifier (assesses evidence sufficiency), Question Manager (generates questions and retrieves answers), Information Retriever (searches for evidence), and Reasoner (synthesizes judgment with rationale and evidence).

Seek clues from external information

Train domain-aligned LLMs for factual reasoning and verifiable judgment



From pipeline usage to intrinsic training, LLMs evolve from external assistants to internalized reasoners for video misinformation detection.

Seek clues from external information

Generate and reason with debunking evidence via multi-agent LLM and diffusion collaboration

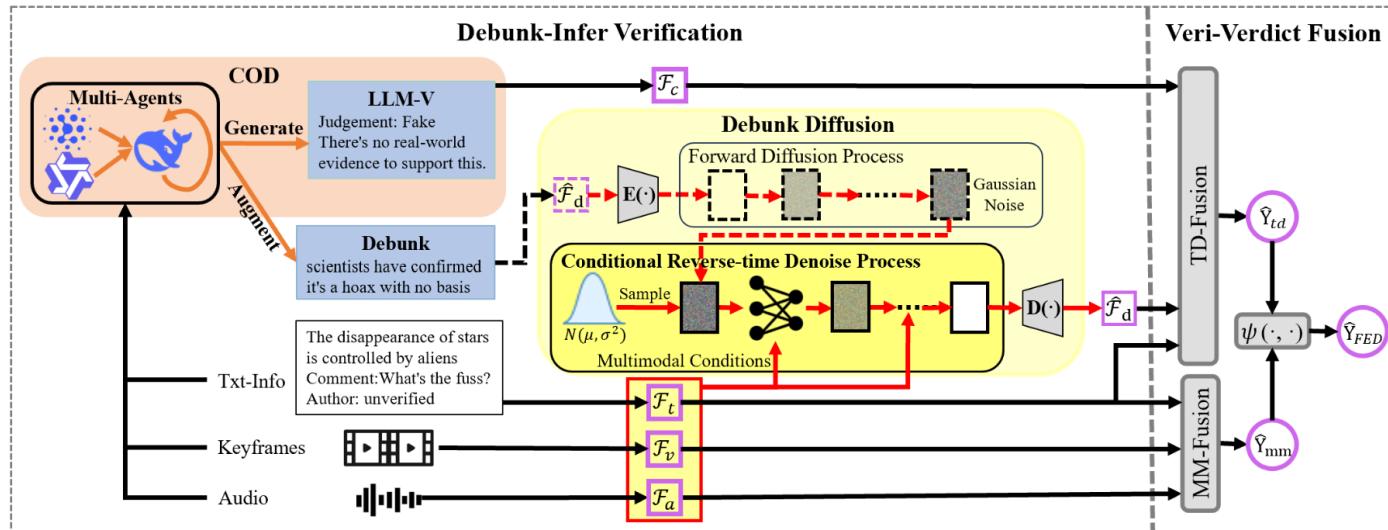


Fig. 1. Overview of DIFND Framework. The DIFND framework consists of Debunk-Infer Verification and Veri-Verdict Fusion. Debunk Diffusion generates debunking cues conditioned on multimodal inputs and is trained on the LLM-augmented dataset, while multi-agent LLMs perform chain-of-debunk for reasoning-rich verification named LLM-V. Final decisions are made via attentive fusion of features from all modules. The dashed arrows indicate paths that are used only during training and are not involved during inference. The textual information with blue background is generated or enhanced by LLMs.

 LLMs act as reasoning agents to generate debunking evidence and verify multimodal claims through diffusion-based inference.

Tutorial Outline

Detection Part I: Human-Edited Misinformation

Signal-based detection

 Editing Traces

 Generation Traces

Semantic-based detection

 Seek clues within the sample's own content

 Seek clues from external information

Intent-based detection

 Social context

Clue integration for misinformation video detection

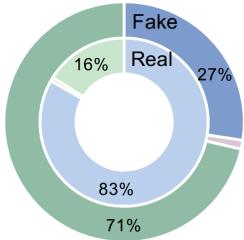
 Parallel Integration

 Sequential Integration

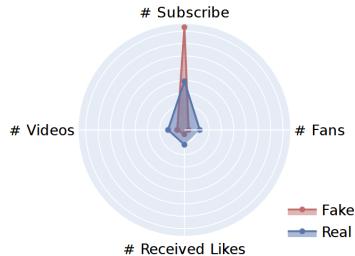
Q+A/Discussion

Intent-based detection

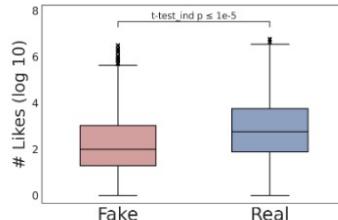
Misinformation reflects underlying user intents, which can be distinguished through social behaviors and engagement patterns.



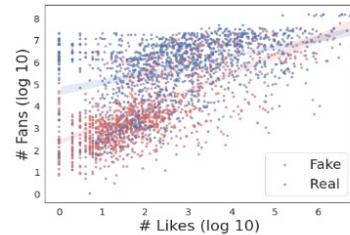
(a) Authority



(b) Statistics



(a) Number of likes



(b) Relationship between the number of publisher fans and likes.

Metadata

- the number of comments
- the number of likes
- the video duration in seconds
- the number of videos that the publisher uploaded
- the follower-following ratio

Comments

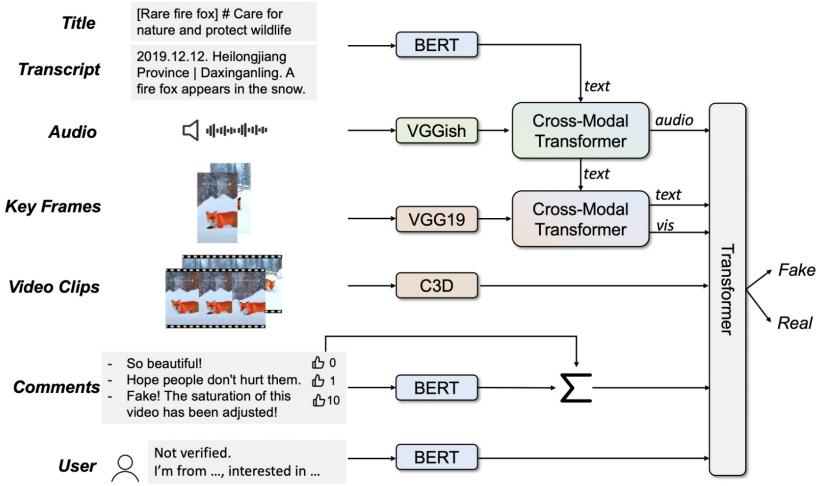
- comments fakeness ratio
- comments inappropriateness ratio (swear words)
- comments conversation ratio (at least one reply)
- top 100 comments tf-idf
- top three popular comments: sentiment polarity, the number of modal particles, the number of personal pronouns and text length



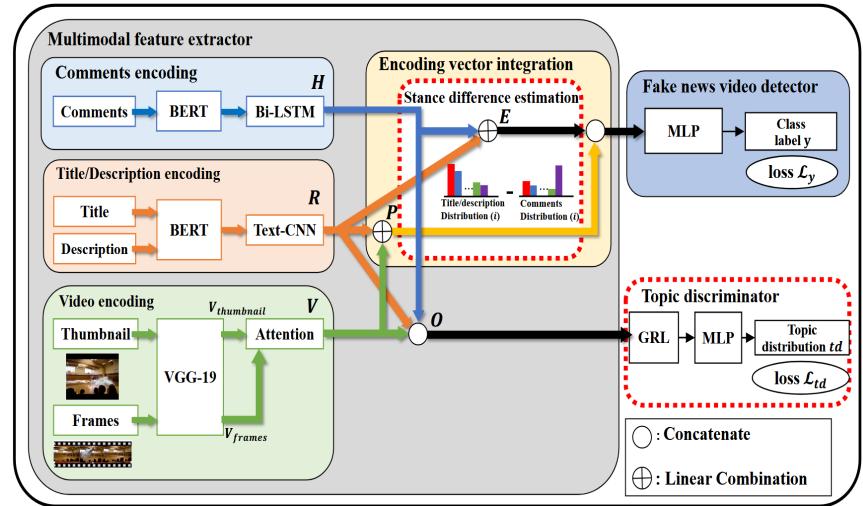
Statistical and social engagement features reveal discriminative behavioral patterns.

Intent-based detection

Model social context features to capture intent-related cues from users and comments.



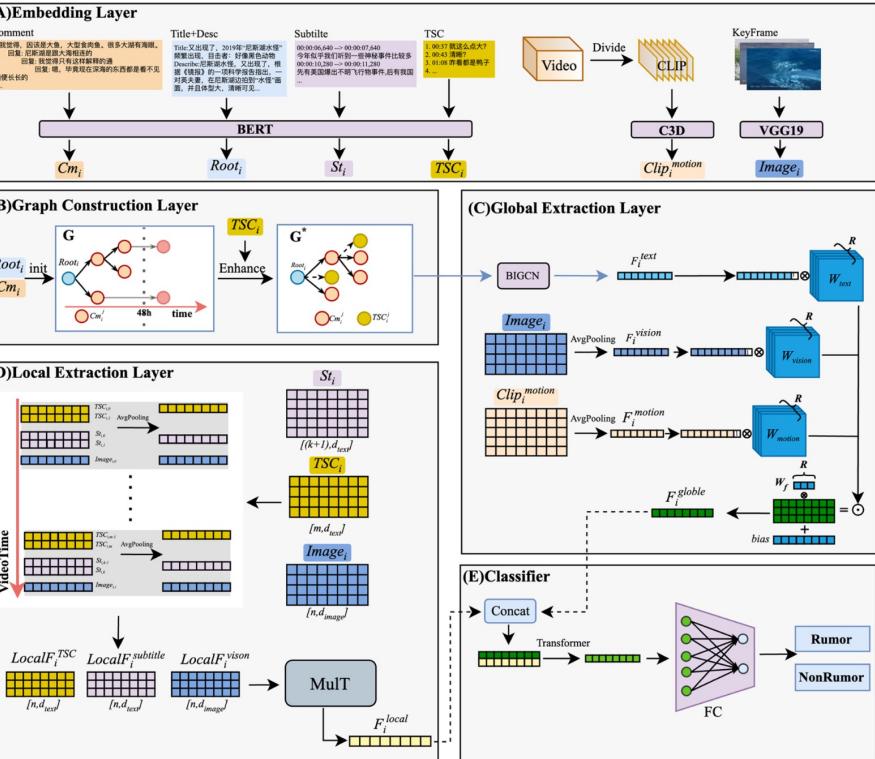
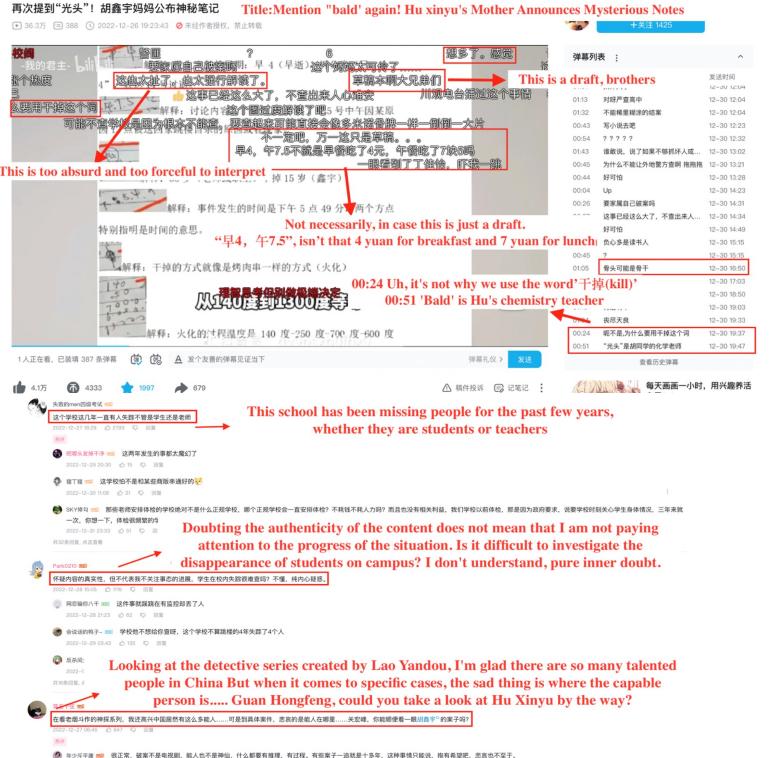
👉 Incorporates **user attributes** and multimodal **comment features** via cross-modal transformers.



👉 Models **stance** and **topic** context between **comments** and video text.

Intent-based detection

Incorporate **time-synchronized comments** to capture dynamic social context



Tutorial Outline

Detection Part I: Human-Edited Misinformation

Signal-based detection

 Editing Traces

 Generation Traces

Semantic-based detection

 Seek clues within the sample's own content

 Seek clues from external information

Intent-based detection

 Social context

 Clue integration for misinformation video detection

 Parallel Integration

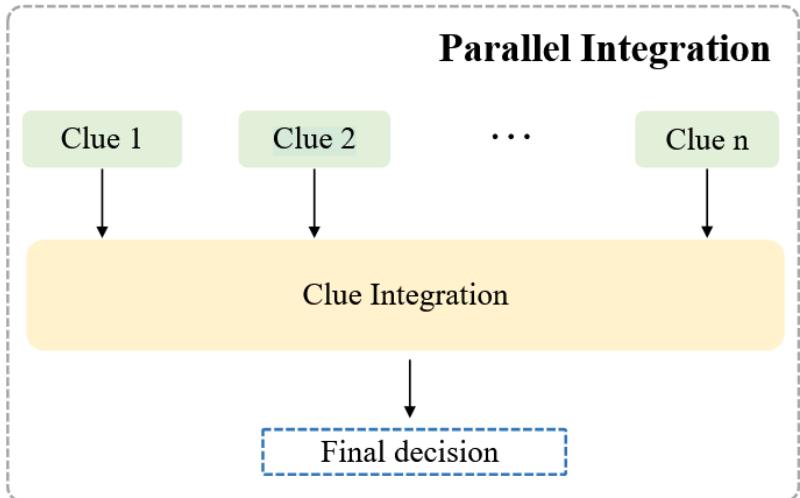
 Sequential Integration

Q+A/Discussion

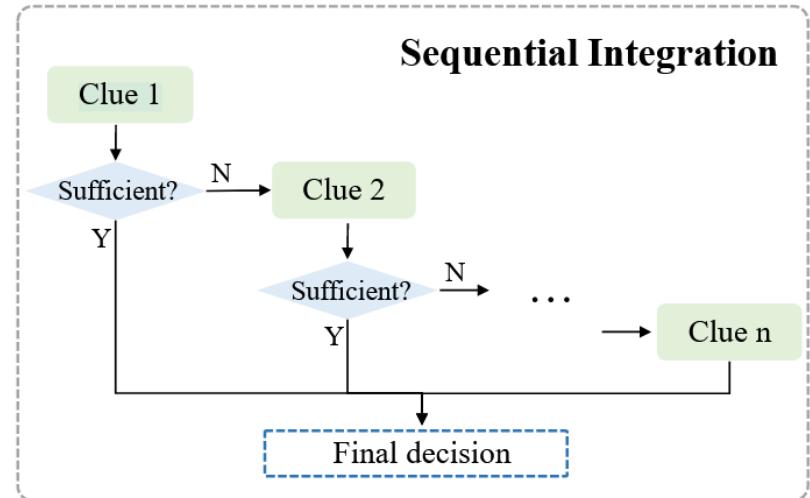
Clue integration for misinformation video detection

Two major paradigms for combining multiple features from different modalities

- **Parallel integration:** all clues from different modalities contribute to the final decision-making process

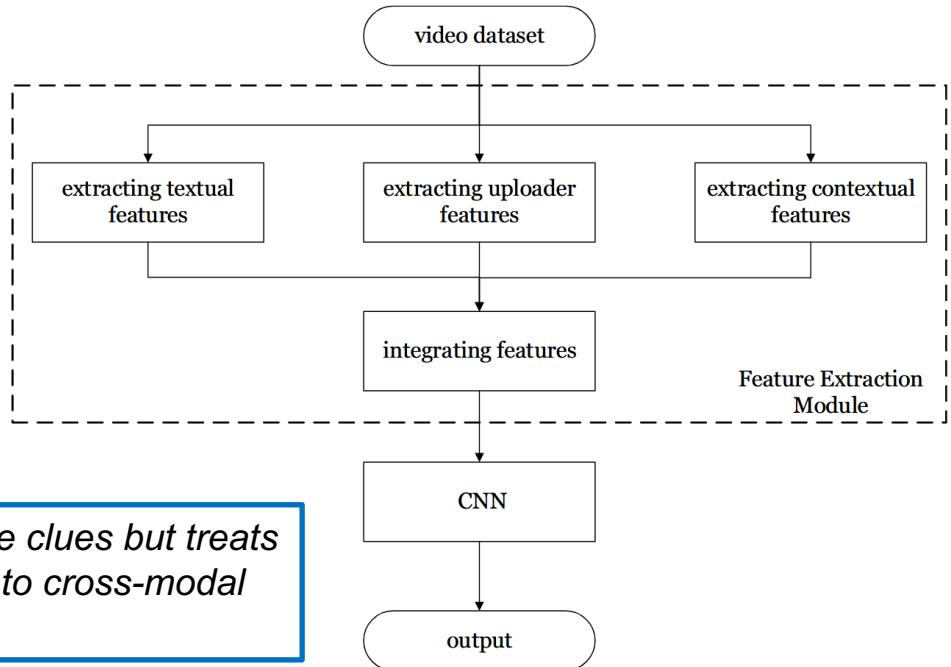
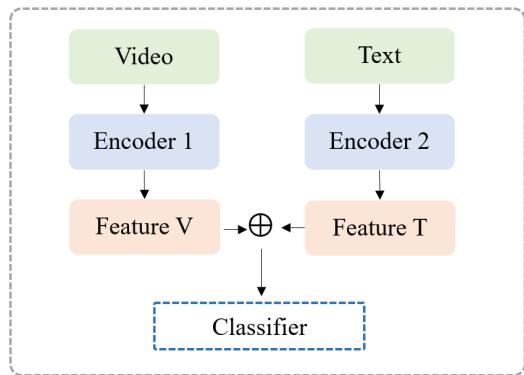


- **Sequential integration:** clues from different modalities are combined in a step-wise manner with each modality contributing incrementally to the final decision.



Concatenation-Based

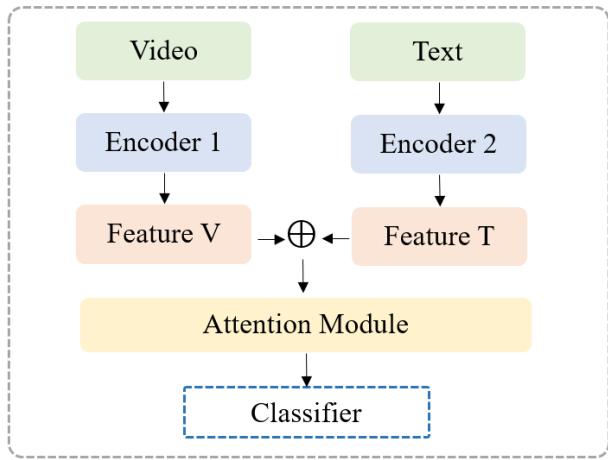
Feature fusion technique: Direct concatenation of multi-modal representations.



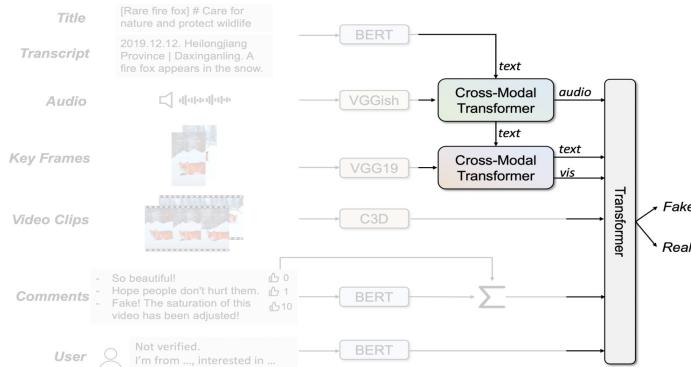
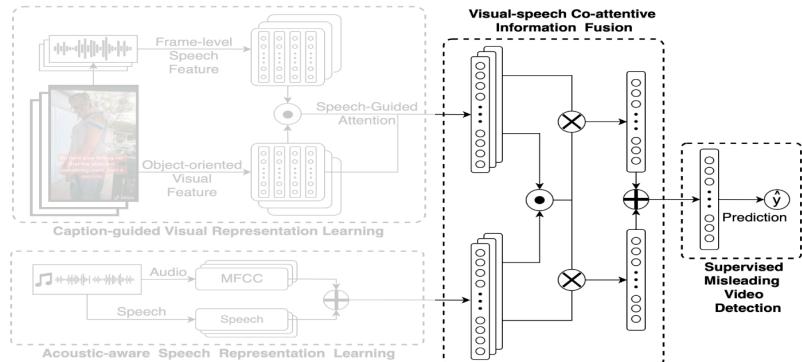
 Concatenation preserves all available clues but treats them equally — efficient yet insensitive to cross-modal dependency.

Attention-Based

Feature fusion technique: Focus on informative clues via attention.



 *Attention-based fusion highlights informative modalities and captures inter-modal interactions dynamically.*

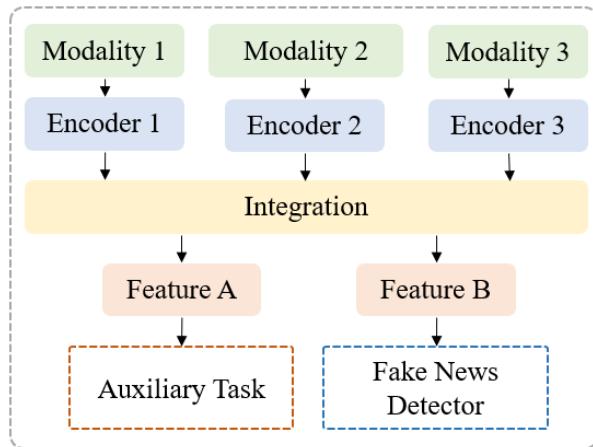


Shang, Lanyu, et al. "A multimodal misinformation detector for covid-19 short videos on tiktok." IEEE Big Data, 2021.

Qi, Peng, et al. "Fakesv: A multimodal benchmark with rich social context for fake news detection on short video platforms." AAAI 2023.

Multitask-Based

Fusion through auxiliary tasks to enhance generalization and consistency.



 **Multitask fusion aligns representations across modalities by leveraging auxiliary supervision — improving robustness under distribution shifts.**

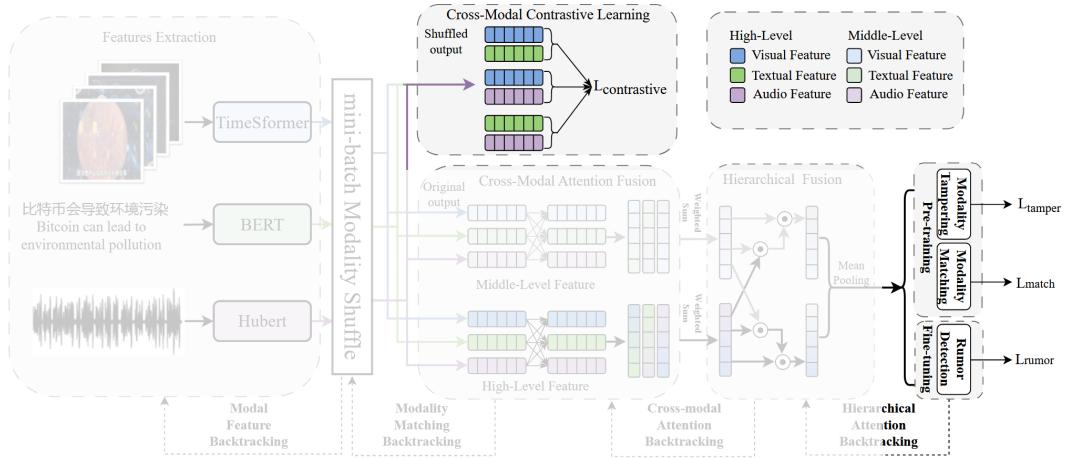


Figure 1: Model Architecture Overview of SVRPM. The model consists of five main modules: (a) Feature extraction: Extract visual, textual, and audio modal features using different encoders, respectively. (b) Mini-batch Modality Shuffle: Randomly shuffle their corresponding modal features for a minibatch. (c) Cross-modal Contrastive Learning: Constructing Positive and Negative Samples using Modality Shuffle Module. (d) Cross-modal Fusion and Hierarchical Fusion. (e) Modality Tampering Backtracking: Using attention backtracking operation to obtain the local features which may be tampered.

Pipeline-Based

Sequential integration: Step-wise reasoning over heterogeneous clues.

 Pipeline-style sequential integration mimics human reasoning — verifying clues in stages.

(Enhance interpretability and efficiency under missing or redundant modalities, but may accumulate early-stage errors).

