① Video to Text Retrieval

Video

Query

Question: When was the design of the plane in the video proposed?

Extract key frames from the video.

CLIP

Generate embeddings for each key frame using CLIP.

Select documents with high embedding similarity.

$cos(.,.)$ 0.32 ✓
$cos(.,.)$ 0.28 ✗
$cos(.,.)$ 0.22 ✗

Title: ...
Content: ...

Title: ...
Content: ...

Title: ...
Content: ...

Title: Airbus A350
Title: Airbus A380
Title: Shuttle Carrier Aircraft
Contnet: The Shuttle Carrier Aircraft (SCA) are two extensively...

Generate embeddings of document titles using CLIP.

Retrieve top documents for each key frame, aggregate as augmented knowledge.

1
2
3

② Speculate Generate

VLM Drafter

V Q + 1 → $e_1$ $r_1$ $a_1$
V Q + 2 → $e_2$ $r_2$ $a_2$
V Q + 3 → $e_3$ $r_3$ $a_3$

Each document is independently processed by a small VLM in parallel to generate an answer $a_i$, along with the identified video entities $e_i$ and reasoning statements $r_i$.

VLM Verifier

$Score_1^{reliabilityility}$
$Score_2^{reliabilityility}$
$Score_3^{reliabilityility}$

Compute the reliability score of $a_i$ using a larger VLM.

Select high-reliability answers $a_i$ with corresponding $e_i$

$a_1$ $e_1$
$a_3$ $e_3$

Compute the similarity between $e_i$ and video $V$ as the alignment score for $a_i$.

V

CLIP

$Score_1^{alignment}$
$Score_3^{alignment}$

argmax

$a_1$

Select the answer $a$ with the corresponding entity $e$ most similar to video $V$ as the final answer.