
Algorithm 1: Confidence-based Event-centric Online Video Question Answering.

Inputs : Online video stream \mathcal{O} , question q , answer space \mathcal{A}

Output: Answer α to the question

1 Initialize confidence thresholds c_{min}, c_{max} ;

2 **for** Every frame $f \in \mathcal{O}$ **do**

3 $c \leftarrow \text{VideoTextEncoder}(q, f)$;

4 **if** $c \geq c_{max}$ **then**

5 $\mathcal{T}_B := \text{Backward-traverse frames till } c < c_{min}$;

6 $\mathcal{T}_F := \text{Forward-traverse frames till } c < c_{min}$;

7 $\mathcal{T} := \mathcal{T}_B \cup \mathcal{T}_F$;

8 Break;

9 **end**

10 **end**

11 Video features $V := \text{VideoEncoder}(\mathcal{T})$;

12 Question features $Q := \text{QuestionEncoder}(q)$;

13 $M := \text{MultiModalEncoder}(Q, V)$;

14 $\alpha \leftarrow \text{AnswerDecoder}(M, \mathcal{A})$;
