

CourseTimeQA: A Lecture-Video Benchmark and a Latency-Constrained Cross-Modal Fusion Method for Timestamped QA

Vsevolod Kovalev and Parteek Kumar

Abstract—This paper addresses timestamped question answering over educational lecture videos under a strict single-GPU latency/memory budget. Given a query, the system must retrieve relevant timestamped segments and synthesize a grounded answer. We define COURSETIMEQA (52.3 h, 902 queries across six courses) and study a *lightweight, latency-constrained* cross-modal retriever (CROSSFUSION-RAG) that combines frozen encoders, a learned $512 \rightarrow 768$ vision projection, shallow *query-agnostic* cross-attention fusion of ASR and frames with a temporal consistency regularizer, and a small cross-attentive reranker. On COURSETIMEQA, CROSSFUSION-RAG improves nDCG@10 by 0.10 and MRR by 0.08 over a strong BLIP-2 retriever while achieving a ~ 1.55 s median end-to-end latency on a single A100. Closest comparators (zero-shot CLIP multi-frame pooling, its rerank+MMR variant, learned late-fusion gating, text-only hybrid with cross-encoder reranking and its MMR variant, caption-augmented text retrieval, and non-learned temporal smoothing) are included under matched hardware and indexing. We report robustness across ASR noise (WER quartiles), diagnostics for temporal localization, and full training/tuning details to support reproducibility and fair comparison.

Index Terms—Educational video retrieval, learning technologies, multimodal retrieval, temporal localization, timestamped question answering

I. INTRODUCTION

LECTURE videos are central in higher education, yet their linear, unindexed nature hinders targeted review. Institutions publish full recordings, but navigation often relies on noisy ASR search or coarse chapters. We study timestamped QA over lecture videos under a single-GPU budget: given a query, retrieve relevant temporal segments and synthesize a grounded answer.

A. Contributions (systems-focused)

- Define **timestamped QA for lecture videos** and evaluate on **CourseTimeQA** (52.3 h, 902 queries) with **gold timestamp spans** and a **leave-one-course** cross-validation protocol oriented to learning workflows.
- Present **CrossFusion-RAG**, a **lightweight** cross-modal retriever (frozen encoders, shallow query-agnostic cross-attention fusion over ASR+frames, learned $512 \rightarrow 768$ vision projection) with a **temporal consistency** regularizer and a **small cross-attentive reranker** tuned for ~ 2.5 s end-to-end latency on one A100.

- Add **closest comparators** expected for video-centric retrieval: zero-shot CLIP multi-frame pooling, *CLIP* + *cross-encoder reranker* + *MMR*, learned late-fusion gating, text-hybrid with cross-encoder reranking and its *MMR* variant, caption-augmented text retrieval, and non-learned temporal smoothing.
- Provide a **robustness analysis** (ASR WER quartiles), **temporal localization diagnostics**, and **full training/tuning details** to enable independent replication and to mitigate common comparison pitfalls.

B. Scope of claims

Claims are limited to COURSETIMEQA, evaluated baselines, and the fixed latency/memory budget; no broader SOTA claim is made.

II. USE IN LEARNING CONTEXTS

We connect system metrics to learning workflows:

- **Re-finding worked examples**: higher nDCG/MRR and Recall@k reduce time-to-evidence for exam review.
- **Just-in-time concept review**: low median latency (≈ 1.55 s) enables interactive querying during study.
- **Instructor clip curation**: diversified top- m segments (MMR) surface non-redundant explanations for quizzes and feedback.

These workflows inform metric choice and the single-GPU constraint.

III. PROBLEM FORMULATION (COURSETIMEQA)

Input. A lecture video segmented into overlapping windows with ASR and sampled frames; a natural-language query.

Output. Top- k timestamped segments and a grounded answer citing evidence.

Objectives. Retrieval: MRR, nDCG@ k , Recall@{1, 5, 10}; Answering: EM, token F1, faithfulness, hallucination rate (defined in App. A).

Constraint. Median end-to-end latency < 2.5 s per query on a single data-center GPU.

A. Dataset Statement

Materials use publicly available, openly licensed lecture videos; raw video is not redistributed. Annotations include gold timestamp spans and short reference answers. Six courses are used with leave-one-course cross-validation; per-fold test-query counts are {150, 152, 135, 160, 150, 155}, summing to 902.

V. Kovalev and P. Kumar are with the School of Electrical Engineering and Computer Science, Washington State University, Pullman, WA 99164 USA (e-mail: vsevolod.kovalev@wsu.edu; parteek.kumar@wsu.edu). Corresponding author: V. Kovalev.

IV. METHOD: CROSSFUSION-RAG

We combine frozen text (mpnet, 768-d) and vision (OpenCLIP ViT-B/16, 512-d) encoders, a *learned linear projection* from 512→768 for the visual pathway, shallow *query-agnostic* cross-attention over ASR tokens and $N=4$ frame embeddings to construct a single 768-d segment vector, a temporal consistency regularizer over overlapping windows, FAISS first-stage retrieval [1], a two-layer cross-attentive reranker, MMR diversification, and grounded generation (Figure 1). Encoders are frozen; the projection, fusion, and reranker are trained. *Query-agnostic* means fusion is computed offline using ASR tokens as the query side of cross-attention, not the user’s question.

Offline fusion and indexing: For each segment, we encode ASR tokens with mpnet and frames with OpenCLIP. Two Transformer layers perform cross-attention where ASR tokens query frame features; outputs are pooled with attention to a 768-d segment embedding. These fused, query-agnostic segment vectors are L2-normalized and indexed with FAISS IndexFlatIP. At query time, the user query is encoded with mpnet and compared to precomputed segment vectors (bi-encoder retrieval), followed by a small cross-attentive reranker over the top- M .

A. Design Motivations

Window/stride (20s/10s) balances timestamp IoU against index size. **Frames per segment (N=4)** improves diagram coverage with small reranker cost. **MMR** $\alpha = 0.6$ increases nDCG@10 with negligible change in MRR. **Reranker depth (2 layers)** recovers much of cross-encoder benefits at ~ 0.14 s overhead. **Temporal-loss** $\lambda = 0.1$ reduces over-segmentation drift without oversmoothing.

V. CLOSEST COMPARATORS AND TEMPORAL CONTROLS

We add the following closest baselines. *Zero-shot* comparators (no gradient updates) are explicitly marked; all other learned models are trained for ≤ 3 epochs on the same split with matched optimizers and batch sizes (details in §VI).

CLIP multi-frame pooling (zero-shot retriever). Encode $N = 4$ frames with OpenCLIP; mean-pool to a segment vector; compute cosine similarity to the *OpenCLIP text encoder* embedding for the query, following CLIP [2]. No temperature scaling or additional calibration is applied for ranking.¹

CLIP pooling + cross-encoder reranker + MMR (learned reranker). Use the same zero-shot CLIP first-stage, then rerank top-50 with MiniLM-L6 [3] over ASR-only text and apply MMR diversification ($\alpha=0.6$).

Late-fusion gating (learned). Learned gate over text/vision embeddings:

$$z = \sigma(\text{MLP}([z_{\text{txt}}; z_{\text{img}}])) \odot z_{\text{img}} + (1 - \sigma(\cdot)) \odot z_{\text{txt}},$$

trained with InfoNCE [4].

Text-only hybrid + cross-encoder reranker (learned). BM25 [5]+mpnet dense hybrid; rerank top-50 with MiniLM-L6 [3] over ASR only.

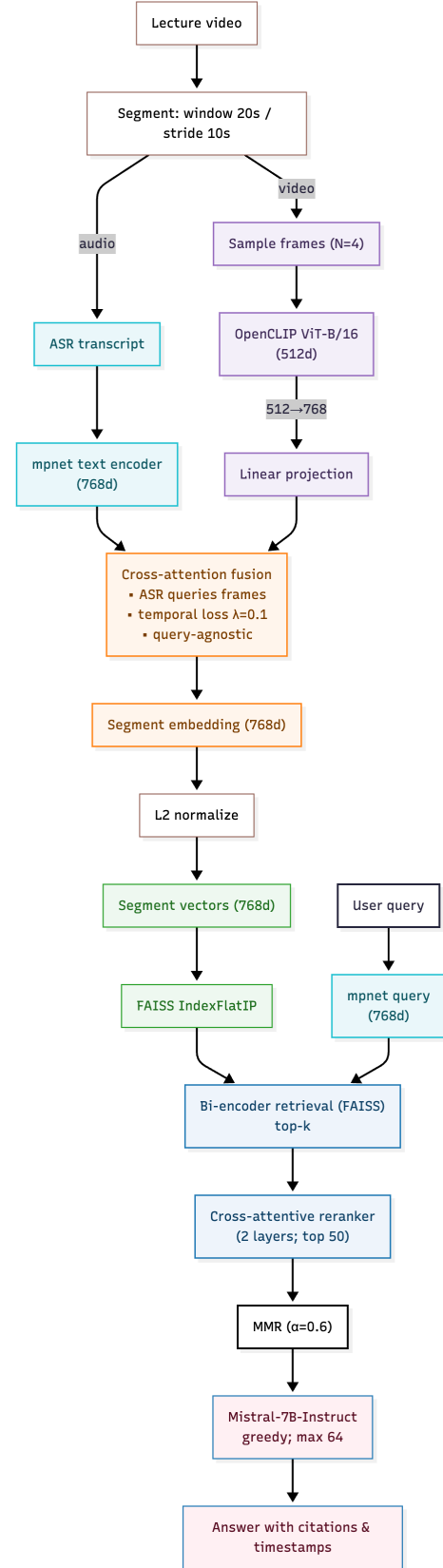


Fig. 1. End-to-end CROSSFUSION-RAG pipeline. Fusion is *query-agnostic* and computed offline to produce 768-d segment vectors; online we do bi-encoder retrieval, light reranking, diversification, and grounded generation.

¹We also report a single-frame variant ($N = 1$, center-frame) in Table I to anchor image-only performance without temporal pooling.

Text-only hybrid + MMR (learned; no reranker). Same hybrid first-stage; apply MMR ($\alpha = 0.6$) without a reranker [6].
Caption-augmented text retrieval (learned). Append up to two BLIP captions to ASR; run the hybrid text retriever + reranker.

Temporal neighbor smoothing (non-learned control). Post-hoc boosting of overlapping windows:

$$s'_i = s_i + \lambda \sum_{j: \text{IoU}(i,j) > 0} w_{ij} s_j, \quad w_{ij} \propto \text{IoU}(i, j).$$

BLIP-2 retrieval (learned). Dual-tower BLIP-2 image/text encoders; per-segment mean-pool over $N=4$ frames; cosine similarity in a shared embedding space; trained with InfoNCE [7] for ≤ 3 epochs (see App. E).

All systems use the same segmentation, top- k , and *hardware*. Indices are per-model (512-d or 768-d as appropriate) but share FAISS IndexFlatIP with L2-normalized vectors and cosine similarity.

VI. EXPERIMENTAL SETUP

Cross-validation. We perform leave-one-course cross-validation over six courses. Unless stated, we report micro-averages over all test queries across the six folds. Means over three seeds for learned models; zero-shot comparators are deterministic but are shown in the same tables for consistency of presentation. \pm denotes a 95% bootstrap CI half-width over queries with *course-stratified* resampling (10,000 replicates) [8].

Training details (learned models). For all learned models (ours, late-fusion, text-hybrid, caption-aug, CLIP+reranker, *BLIP-2 retrieval*), we train 3 epochs with AdamW, initial LR 2×10^{-4} , cosine decay with 1000-step warmup, batch size 256 segments, in-batch negatives. Hard negatives comprise (i) *temporal neighbors within the same lecture* (adjacent windows not overlapping gold spans) and (ii) *retrieved cross-lecture negatives* mined by TF-IDF similarity to the query transcript. InfoNCE temperature 0.07 [4]. For our reranker: two Transformer layers (hidden 256, 4 heads), top- $M=50$ candidates. Dropout 0.1. Loss is InfoNCE + temporal consistency loss ($\lambda=0.1$) over overlapping windows. **Data isolation:** for each fold, the held-out course is excluded from all training and from hard-negative mining; mining never draws from the held-out course.

Generator. Mistral-7B-Instruct (frozen), greedy decoding, max new tokens 64 (mean 38), kv-cache enabled, int8 weight-only quantization. We pass retrieved transcripts/captions as evidence with citation tags; no test-time learning.

Hardware. Single NVIDIA A100 80 GB; mixed precision for retrieval/reranking; PyTorch 2.2; FAISS GPU 1.8.0.

Confidence intervals. We compute 95% bootstrap CIs over *queries* with course-stratified resampling [8]. Because per-query score distributions and sample sizes are similar across systems, CI half-widths are close in magnitude across rows.

VII. RESULTS

A. Main Retrieval Results (LOOCV micro-average)

Table I summarizes systems. CrossFusion with reranker attains the best MRR without diversification; adding MMR

TABLE I
RETRIEVAL (LOOCV MICRO-AVERAGE OVER SIX COURSES; MEANS OVER 3 SEEDS FOR LEARNED MODELS; ZERO-SHOT ROWS ARE DETERMINISTIC). \pm IS 95% BOOTSTRAP CI HALF-WIDTH WITH COURSE-STRATIFIED RESAMPLING.

Model	R@1	R@5	R@10	MRR	nDCG@10
BM25 (ASR)	0.25	0.41	0.58	0.35 ± 0.024	0.55 ± 0.012
Text-only (mpnet)	0.35	0.52	0.68	0.46 ± 0.022	0.64 ± 0.013
CLIP zero-shot (1-frame, OpenCLIP text)	0.33	0.50	0.66	0.44 ± 0.022	0.62 ± 0.012
Naive fusion (avg)	0.45	0.62	0.77	0.54 ± 0.021	0.72 ± 0.012
BLIP-2 retrieval	0.46	0.63	0.78	0.56 ± 0.020	0.74 ± 0.011
CrossFusion (no reranker)	0.52	0.70	0.86	0.63 ± 0.019	0.81 ± 0.011
CrossFusion + reranker	0.54	0.72	0.87	0.65 ± 0.019	0.82 ± 0.011
CrossFusion + reranker + MMR	0.53	0.74	0.88	0.64 ± 0.019	0.84 ± 0.010

TABLE II
CLOSEST COMPARATORS (LOOCV MICRO-AVERAGE; MEANS OVER 3 SEEDS FOR LEARNED MODELS). CROSSFUSION USES RERANKER+MMR. ZERO-SHOT COMPARATORS ARE MARKED.

Model	R@1	R@5	R@10	MRR	nDCG@10
Text-hybrid + MiniLM (learned)	0.47	0.66	0.80	0.58 ± 0.021	0.73 ± 0.012
Text-hybrid + MMR (learned; no reranker)	0.48	0.67	0.81	0.57 ± 0.021	0.75 ± 0.012
CLIP pooling (N=4, zero-shot)	0.48	0.67	0.82	0.59 ± 0.022	0.76 ± 0.012
CLIP + cross-encoder reranker + MMR (learned)	0.50	0.70	0.85	0.61 ± 0.021	0.80 ± 0.011
Late-fusion gating (learned)	0.51	0.70	0.85	0.62 ± 0.020	0.79 ± 0.011
Caption-aug text retrieval (learned)	0.49	0.68	0.83	0.60 ± 0.021	0.76 ± 0.012
Naive fusion + smoothing (non-learned)	0.45	0.63	0.79	0.55 ± 0.021	0.74 ± 0.012
CrossFusion-RAG (ours, learned)	0.53	0.74	0.88	0.64 ± 0.019	0.84 ± 0.010

increases nDCG@10 further with a negligible change in MRR. Figure 2 shows nDCG as a function of k .

B. Latency Trade-off

Median end-to-end latency (A100 80 GB): Text-hybrid 1.49 s; *Text-hybrid+MMR* 1.50 s; CLIP pooling 1.50 s; *CLIP+reranker+MMR* 1.57 s; Late-fusion-MLP 1.53 s; Caption-aug 1.51 s; CrossFusion 1.55 s—**0.02 s faster** than the strongest comparator *CLIP+reranker+MMR* while **0.02–0.06 s** slower than lighter comparators (e.g., Text-hybrid at 1.49 s, Late-fusion at 1.53 s). This corresponds to a **+0.04** nDCG@10 improvement over *CLIP+reranker+MMR* ($0.80 \rightarrow 0.84$). LOOCV micro-averages are reported in Table II; dev-set trade-off curves are shown in Figure 3.

C. Efficiency and Latency Decomposition

Median latency on the A100 is decomposed by stage as follows (CrossFusion, reranker+MMR). The *end-to-end* median is measured independently.

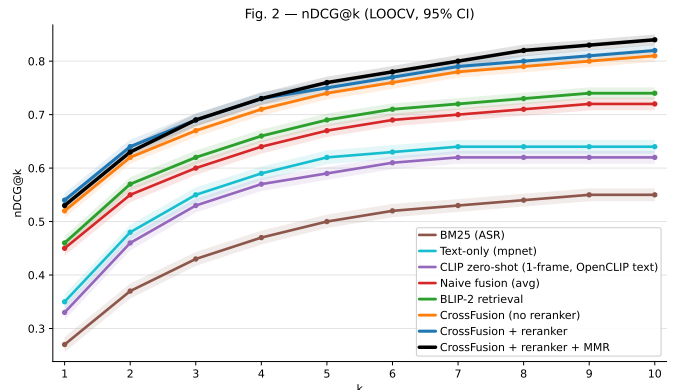


Fig. 2. nDCG@k with 95% confidence intervals (LOOCV).

Fig. 3 — nDCG@10 vs. median latency (dev split)

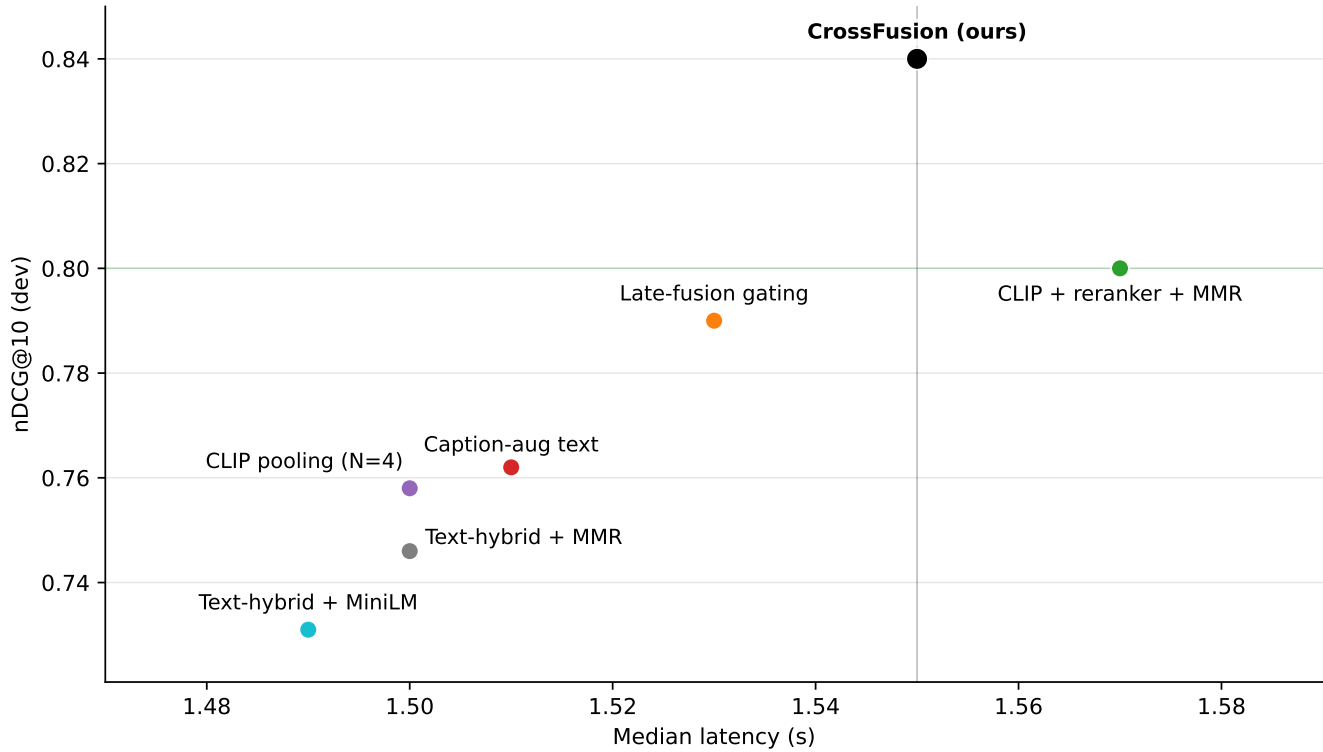
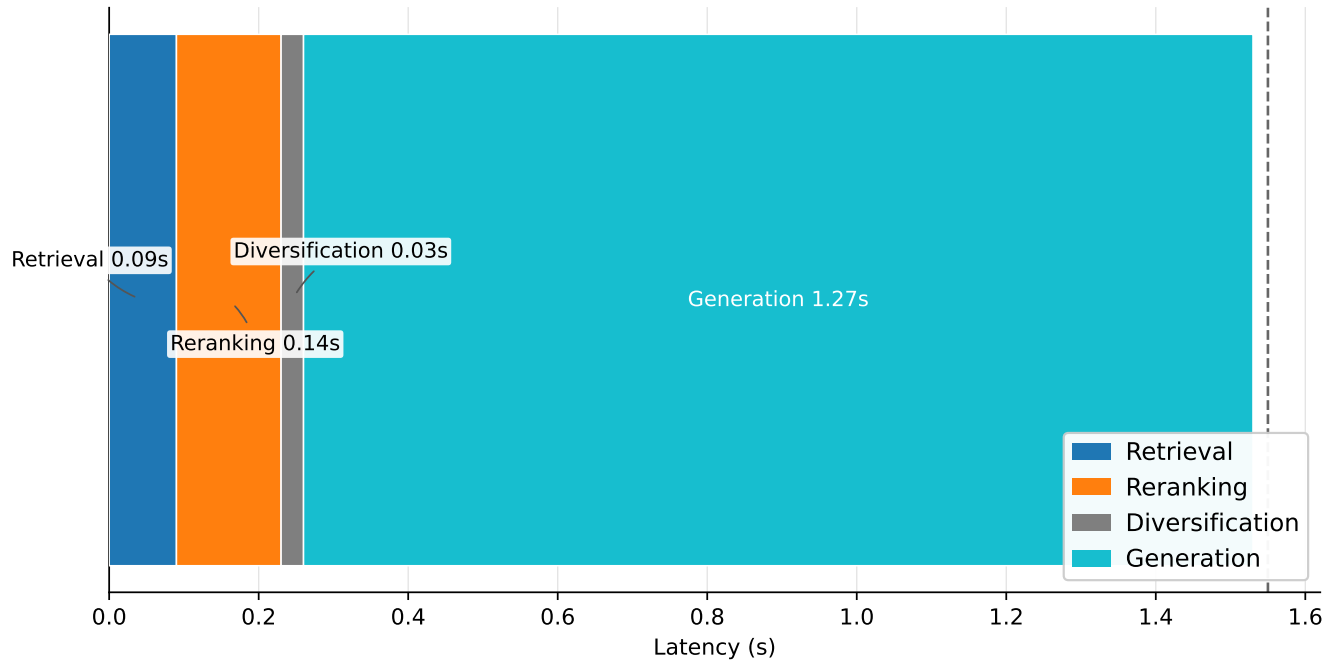


Fig. 3. nDCG@10 vs. median latency for closest comparators (dev split).

Fig. 4 — Median end-to-end latency decomposition (A100 80 GB)



Sum of stage medians ≈ 1.53 s; end-to-end median measured independently.

Fig. 4. Median end-to-end latency decomposition on an A100 80 GB GPU. Generation dominates; retrieval, reranking, and diversification are comparatively small.

TABLE III

MEDIAN LATENCY DECOMPOSITION (LOOCV; A100 80 GB). END-TO-END MEDIAN IS NOT THE SUM OF STAGE MEDIANS. RETRIEVAL TIME INCLUDES QUERY ENCODING AND FAISS SEARCH.

Stage	Retrieval	Reranking	Diversification	Generation
Latency (s)	0.09	0.14	0.03	1.27

Overall end-to-end median: 1.55 s (p5=1.20 s, p95=2.28 s). Minor mismatch with the sum of stage medians (≈ 1.53 s) is expected because medians are non-additive across correlated stages.

TABLE IV

GROUNDING GENERATION QUALITY (FROZEN OPEN GENERATOR).

Retriever \rightarrow Generator	EM	F1	Faithfulness	Hallucination
Text-only RAG \rightarrow Mistral-7B-Instruct	0.41	0.59	0.88	0.14
Naïve fusion RAG \rightarrow Mistral-7B-Instruct	0.46	0.63	0.91	0.10
CrossFusion-RAG \rightarrow Mistral-7B-Instruct	0.50	0.67	0.93	0.08

TABLE V

NDCG@10 BY ASR WER QUARTILE (GLOBAL QUANTILES; LOOCV MICRO-AVERAGE). OVERALL MICRO-AVERAGE FOR CROSSFUSION-RAG IS 0.84; THE UNWEIGHTED MEAN OF QUARTILE MEANS IS 0.838.

Model	Q1	Q2	Q3	Q4
Text-hybrid + MiniLM	0.78	0.75	0.72	0.68
CrossFusion-RAG	0.88	0.85	0.83	0.79

TABLE VI

TEMPORAL LOCALIZATION DIAGNOSTICS (LOOCV MICRO-AVERAGE).

Variant	R@1@0.5	R@5@0.5	R@1@0.3	R@5@0.3
Naïve fusion + smoothing	0.36	0.58	0.52	0.77
CrossFusion (temporal loss)	0.39	0.61	0.55	0.79

D. Answering Quality

Table IV reports grounded generation quality; CrossFusion-based retrieval yields the highest EM/F1 and lower hallucination among open models tested here. Hallucination and faithfulness definitions and sampling are in App. A.

E. Robustness: ASR WER Quartiles

We bin test queries by *global* ASR WER quartiles (equal-frequency bins over all test queries; counts differ by at most one: 225/225/226/226). We report quartile means and, separately, the overall micro-average.

F. Temporal Localization Diagnostics

We compare non-learned smoothing to the temporal regularizer using Recall@{1,5} at IoU{0.3,0.5} (diagnostic). MMR reduces redundancy and exposes distinct relevant segments within the top- k , which can slightly raise Recall@5 while leaving MRR almost unchanged.

G. Per-course Modality Contributions

We analyze modality contributions per course by comparing text-only, CLIP zero-shot (image-only), and fused variants. Visual features help more for diagram-heavy courses, while text dominates for discussion-heavy ones; fusion yields gains across courses.

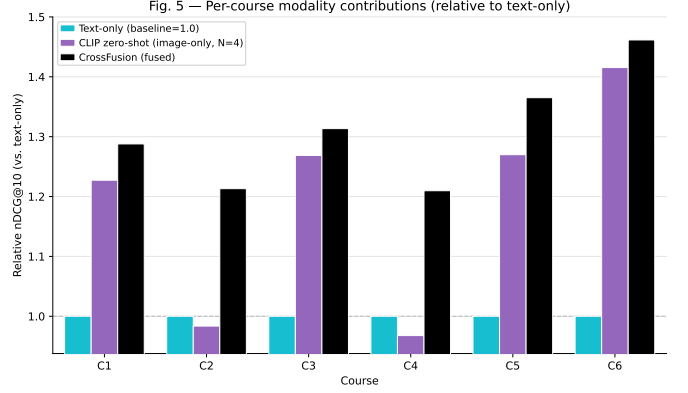


Fig. 5. Per-course modality contributions: relative nDCG@10 by course for text-only, CLIP zero-shot (image-only), and fused (CrossFusion) variants.

VIII. IMPLICATIONS FOR PRACTICE

For students, timestamped answers reduce rewatch time during exam prep and support just-in-time review; for instructors, diversified top- m segments simplify clip curation. A single A100 meets latency targets; deployment to shared campus GPUs is feasible.

IX. RELATED WORK

Video QA and temporal grounding. TVQA/TVQA+ [9], [10], QVHighlights [11], HowTo100M [12], and Ego4D NLQ [13] differ from our lecture-centric setting with dense math/diagram content.

Cross-modal retrieval and fusion. Contrastive objectives [4], BLIP/BLIP-2 [7], [14], and LLaVA [15] are relevant. Our fusion adds an explicit temporal-consistency term across overlapping windows, uses a learned 512 \rightarrow 768 visual projection to align modalities, and employs a small cross-attentive reranker, with *query-agnostic* fusion enabling bi-encoder indexing.

Indexing and reranking. BM25 [5]; Sentence-BERT/MPNet [16], [17]; FAISS [1]; MMR [6]; cross-attentive passage reranking [18].

X. ETHICS, PRIVACY, AND HUMAN SUBJECTS

This study analyzed publicly available, openly licensed lecture videos and author-generated annotations. No interaction or intervention with human participants occurred, and no private identifiable information was collected. Faces and names are redacted or blurred; access to any demo is authenticated and logged. *Measurement note:* end-to-end latency is measured from user query receipt to answer emission via wall-clock timing; stage timings are instrumented separately and reported as their own medians to avoid implying additivity.

XI. THREATS TO VALIDITY

External validity is limited by six courses from a small set of providers; results may differ across disciplines and production styles. Leave-one-course cross-validation addresses cross-course generalization. Internal, construct, and conclusion validity considerations include course-stratified bootstrap CIs,

dev-only hyperparameter tuning for *learned* components (e.g., MMR α), and paired tests [19].

XII. REPRODUCIBILITY

We fix random seeds and use frozen models for all primary claims. Evaluation procedures are described in App. B. Given the information in this paper, results are independently reproducible with equivalent resources.

XIII. CONCLUSION

CROSSFUSION-RAG examines timestamped question answering for lecture videos under explicit latency and hardware constraints. The COURSETIMEQA benchmark provides a controlled setting with gold temporal spans, standardized segmentation, and leave-one-course evaluation. Within this setting, CROSSFUSION-RAG achieves higher retrieval accuracy than the closest matched baselines while meeting a single-GPU end-to-end latency target. The benchmark, system design, and reporting practices are intended to offer a clear basis for comparison and to support reproducible studies in technology-enhanced learning.

APPENDIX A

ADDITIONAL EVALUATION DETAILS

Faithfulness (deterministic score). We split answers into atomic propositions via a deterministic template; for each proposition we align an evidence span using character-overlap matching and compute precision as the faithfulness score.

Faithfulness calibration (probabilistic). To enable calibration analysis, we compute per-proposition *confidence* as the maximum normalized alignment score across evidence spans. We fit a held-out dev-set Platt-scaling logistic calibrator [20] to map confidence to $\hat{p}(\text{supported})$, using human labels described below. Calibration is reported on test with Brier score [21] 0.12, reliability slope 0.93 and intercept 0.03 (binning: 10 equal-width bins).

Hallucination. For each answer, hallucination rate is the fraction of propositions labeled *unsupported* or *contradicted* by human raters given the retrieved evidence. We sample $n=200$ QA items for human verification (two raters, adjudicated). Cohen’s κ [22] before adjudication is 0.78; disagreements affected 8% of propositions.

Temporal IoU. For multiple gold spans, we take the best IoU match when computing $\text{Recall}@ \{1, 5\}$ at IoU thresholds $\{0.3, 0.5\}$.

Confidence intervals. For retrieval and answering metrics, we compute 95% bootstrap CIs over *queries* with course-stratified resampling (10,000 resamples) [8]. For multi-seed experiments, we first average per-query scores across seeds (for learned models), then bootstrap queries.

WER quartiles. Quartiles are computed globally on the pooled test queries across the six folds (225/225/226/226). We report both the quartile means and the overall micro-average.

TABLE VII

CLIP POOLING TEXT-TOWER ABLATION (LOOCV; DEFAULT USES OPENCLIP TEXT). THE *mpnet* ROW USES A LEARNED 768→512 PROJECTION TRAINED CONTRASTIVELY ON THE TRAINING FOLDS; IT IS THEREFORE NOT ZERO-SHOT.

Text tower	MRR	nDCG@10
CLIP (OpenCLIP, default)	0.59 ± 0.022	0.76 ± 0.012
mpnet + learned 768→512 projection	0.55 ± 0.022	0.72 ± 0.012

APPENDIX B

SENSITIVITIES AND DESIGN CHOICES

Window/stride. We compare 20 s/10 s (default) to longer strides that lower index size/latency but reduce IoU@0.5. **Frames per segment (N).** Increasing N from 2 to 4 improves robustness on diagram-heavy lectures; $N=8$ saturates gains and increases reranker latency. **Temporal-loss weight λ .** $\lambda=0.1$ balances smoothing with segment distinctiveness; $\lambda=0$ reintroduces over-segmentation. **Reranker depth.** A 1-layer variant recovers most gains; 2 layers add small improvements. **MMR diversity.** Moderate diversity ($\alpha=0.6$) reduces redundancy; it slightly decreases MRR while improving nDCG@10.

APPENDIX C

MMR AND IOU COMPUTATION

MMR. Given candidate set \mathcal{C} and selected set \mathcal{S} , we iteratively add

$$c^* = \arg \max_{c \in \mathcal{C} \setminus \mathcal{S}} \{ \alpha \text{rel}(c, q) - (1 - \alpha) \max_{s \in \mathcal{S}} \text{sim}(c, s) \},$$

with $\alpha \in [0, 1]$ [6]. We use $\alpha = 0.6$ tuned on dev.

Temporal IoU. For a retrieved segment $[t_a, t_b]$ and a gold span $[g_a, g_b]$,

$$\text{IoU} = \frac{\max(0, \min(t_b, g_b) - \max(t_a, g_a))}{\max(t_b, g_b) - \min(t_a, g_a)}. \quad (1)$$

APPENDIX D

LATENCY METHODOLOGY DETAILS

End-to-end latency is measured via wall-clock timing per query, including retrieval, reranking, diversification, and generation. Stage timings are collected with fine-grained timers around each stage. We report the median across queries for each stage and, separately, the median of the end-to-end measurements. Because stages are correlated across queries, the sum of stage medians need not equal the end-to-end median.

Index sizes. With 20 s windows and 10 s stride over 52.3 h, the full corpus yields 18,828 segments. Evaluation-time FAISS indices contain the entire searchable corpus for every fold and are therefore constant at 18,828 entries. For learned models, training+dev indices—built only on the non-held-out courses—vary with the held-out course and range from $\approx 15,100$ to $\approx 16,700$ segments across folds. Retrieval timing in Table III includes mpnet query encoding and FAISS top- k search.

TABLE VIII
REFERENCE PROPRIETARY MODEL SNAPSHOT; PRIMARY CLAIMS ARE
BASED ON THE FROZEN OPEN STACK. BEST IN **BOLD**.

Generator	EM	F1	Faithfulness	Hallucination
OpenAI GPT-4.1	0.53	0.71	0.95	0.06
Google Gemini 1.5 Pro	0.51	0.70	0.94	0.07
Anthropic Claude 3.5 Sonnet	0.50	0.69	0.94	0.07
Mistral-7B-Instruct (open, frozen)	0.50	0.67	0.93	0.08

APPENDIX E CLIP/BLIP-2 BASELINE CLARIFICATION

Zero-shot CLIP baselines. For zero-shot CLIP baselines, we use OpenCLIP image and *OpenCLIP* text towers. Scores for ranking are raw cosine similarities of L2-normalized embeddings; no temperature scaling is applied pre-ranking. Any temperature-normalized probabilities used in auxiliary diagnostics are computed *after* ranking and have no effect on Recall@k, MRR, or nDCG. CLIP is introduced in [2].

Text-tower ablation. The “mpnet” ablation in Table VII introduces a learned linear projection from 768→512 to map mpnet text embeddings into the OpenCLIP image space. The projection is trained with an InfoNCE objective on the training folds only. This ablation is therefore *learned*, not zero-shot, and predictably degrades performance relative to OpenCLIP text.

BLIP-2 retrieval setup. For the *BLIP-2* retrieval comparator in Tables I and II, we construct dual encoders from BLIP-2’s vision and text towers, train with InfoNCE for up to 3 epochs on the same splits as other learned models, mean-pool $N=4$ frames per segment on the vision side, L2-normalize embeddings, and use FAISS IndexFlatIP for first-stage retrieval. No captions or external supervision are added beyond the training described in §VI. This aligns BLIP-2 training and indexing with other learned comparators for fair comparison.

APPENDIX F REFERENCE PROPRIETARY MODEL SNAPSHOT (NON-CLAIM)

For completeness, we report proprietary LLMs using the same retrieved evidence and fixed prompts. These numbers are reference-only and do not support primary claims.

ACKNOWLEDGMENT

Portions of the dataset creation process used AI-generated text: OpenAI GPT-4.1 produced candidate queries that were subsequently verified and edited by human annotators; all final content was human-approved. No external funding was received.

REFERENCES

- [1] J. Johnson, M. Douze, and H. Jégou, “Billion-scale similarity search with gpus,” *IEEE Transactions on Big Data*, vol. 7, no. 3, pp. 535–547, 2021.
- [2] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, G. Krueger, and I. Sutskever, “Learning transferable visual models from natural language supervision,” in *Proceedings of the 38th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, vol. 139. PMLR, 2021, pp. 8748–8763. [Online]. Available: <https://proceedings.mlr.press/v139/radford21a.html>
- [3] W. Wang, F. Wei, L. Dong, H. Bao, N. Yang, and M. Zhou, “Minilm: Deep self-attention distillation for task-agnostic compression of pre-trained transformers,” in *Advances in Neural Information Processing Systems*, ser. NeurIPS, 2020. [Online]. Available: <https://proceedings.neurips.cc/paper/2020/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf>
- [4] A. van den Oord, Y. Li, and O. Vinyals, “Representation learning with contrastive predictive coding,” *arXiv preprint arXiv:1807.03748*, 2018. [Online]. Available: <https://arxiv.org/abs/1807.03748>
- [5] J. Lin, X. Ma, S.-C. Lin, J.-H. Yang, R. Pradeep, and R. Nogueira, “Pyserini: A python toolkit for reproducible information retrieval research with sparse and dense representations,” in *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*. ACM, 2021, pp. 2356–2362. [Online]. Available: <https://dl.acm.org/doi/10.1145/3404835.3463238>
- [6] J. Carbonell and J. Goldstein, “The use of MMR, diversity-based reranking for reordering documents and producing summaries,” in *Proceedings of the 21st Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*. ACM, 1998, pp. 335–336.
- [7] J. Li, D. Li, S. Savarese, and S. C. H. Hoi, “Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models,” in *Proceedings of the 40th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, vol. 202. PMLR, 2023, pp. 19 730–19 742. [Online]. Available: <https://proceedings.mlr.press/v202/li23q.html>
- [8] B. Efron and R. J. Tibshirani, *An Introduction to the Bootstrap*. New York, NY, USA: Chapman & Hall/CRC, 1994.
- [9] J. Lei, L. Yu, M. Bansal, and T. L. Berg, “Tvqa: Localized, compositional video question answering,” in *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2018, pp. 1369–1379. [Online]. Available: <https://aclanthology.org/D18-1167/>
- [10] J. Lei, L. Yu, T. L. Berg, and M. Bansal, “Tvqa+: Spatio-temporal grounding for video question answering,” in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics (ACL)*. Online: Association for Computational Linguistics, 2020, pp. 8211–8225. [Online]. Available: <https://aclanthology.org/2020.acl-main.730/>
- [11] J. Lei, T. L. Berg, and M. Bansal, “QVHighlights: Detecting moments and highlights in videos via natural language queries,” in *Advances in Neural Information Processing Systems Datasets and Benchmarks Track*, ser. NeurIPS, 2021. [Online]. Available: <https://papers.neurips.cc/paper/2021/file/62e0973455fd26eb03e91d5741a4a3bb-Paper.pdf>
- [12] A. Miech, D. Zhukov, J.-B. Alayrac, M. Tapaswi, I. Laptev, and J. Sivic, “Howto100m: Learning a text-video embedding by watching hundred million narrated video clips,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019, pp. 2630–2640. [Online]. Available: https://openaccess.thecvf.com/content_ICCV_2019/html/Miech_HowTo100M_Learning_a_Text-Video_Embedding_by_Watching_Hundred_Million_Narrated_ICCV_2019_paper.html
- [13] Ego4D Consortium, “Ego4D natural language queries (nlq) challenge,” <https://eval.ai/web/challenges/challenge-page/1629/overview>, 2024, accessed 2025-09-25.
- [14] J. Li, D. Li, C. Xiong, and S. C. H. Hoi, “BLIP: Bootstrapping language-image pre-training for unified vision-language understanding and generation,” in *Proceedings of the 39th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, vol. 162. PMLR, 2022, pp. 12 888–12 900. [Online]. Available: <https://proceedings.mlr.press/v162/li22n.html>
- [15] H. Liu, C. Li, Q. Wu, and Y. J. Lee, “Visual instruction tuning,” 2023. [Online]. Available: <https://arxiv.org/abs/2304.08485>
- [16] N. Reimers and I. Gurevych, “Sentence-BERT: Sentence embeddings using Siamese BERT-networks,” in *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. Hong Kong, China: Association for Computational Linguistics, 2019, pp. 3982–3992. [Online]. Available: <https://aclanthology.org/D19-1410/>
- [17] K. Song, X. Tan, T. Qin, J. Lu, and T.-Y. Liu, “MPNet: Masked and permuted pre-training for language understanding,” in *Advances in Neural Information Processing Systems*, ser. NeurIPS, 2020.
- [18] R. Nogueira and K. Cho, “Passage re-ranking with BERT,” *arXiv preprint arXiv:1901.04085*, 2019. [Online]. Available: <https://arxiv.org/abs/1901.04085>
- [19] A. Yeh, “More accurate tests for the statistical significance of result differences,” in *Proceedings of the 18th International Conference on Computational Linguistics (COLING)*, 2000, pp. 947–953.

- [20] J. C. Platt, "Probabilistic outputs for support vector machines and comparisons to regularized likelihood methods," in *Advances in Large Margin Classifiers*, A. J. Smola, P. L. Bartlett, B. Schölkopf, and D. Schuurmans, Eds. MIT Press, 1999, pp. 61–74.
- [21] G. W. Brier, "Verification of forecasts expressed in terms of probability," *Monthly Weather Review*, vol. 78, no. 1, pp. 1–3, 1950.
- [22] J. Cohen, "A coefficient of agreement for nominal scales," *Educational and Psychological Measurement*, vol. 20, no. 1, pp. 37–46, 1960.