

# AuraLearn: Evaluating Multi-Tier Explainability in a Hybrid RAG Framework for Educational Document Intelligence

Yashvanth Karunakaran<sup>1</sup>, Bhenedix Paul<sup>1</sup>, and Joshnavi Pokala<sup>1</sup>

School of Computer Science and Engineering,  
Vellore Institute of Technology Chennai, India  
{yashvanth.2023, bhenedix.paul2023, joshnavi.pokala2023}@vitstudent.ac.in

**Abstract.** Large Language Models (LLMs) and Retrieval-Augmented Generation (RAG) have transformed educational technology, yet their inherent black-box nature and propensity for hallucinations critically limit their trustworthiness in academic environments. Students require syllabus-bounded, mathematically verifiable answers rather than ungrounded probabilistic generation. To address this, we propose **AuraLearn**, a domain-specific educational assistant integrating a multi-tier Explainable AI (XAI) framework. AuraLearn employs a hybrid semantic-lexical retriever (FAISS combined with BM25/TF-IDF) coupled with a dual-stage summarizer utilizing an extractive Bidirectional Long Short-Term Memory (BiLSTM) network and an abstractive T5 transformer. Crucially, we introduce transparent interpretability across all architectural stages: post-hoc score decomposition via Transparent Approximations (BETA) for hybrid retrieval, intrinsic attention weight extraction for the BiLSTM, and leave-one-out sensitivity analysis with token-level confidence for the T5 generator. Experimental evaluations on the CNN/DailyMail 3.0.0 benchmark demonstrate that the extractive model achieves robust factual grounding (ROUGE-1: 0.3046) under severe class imbalance utilizing Focal Loss ( $\alpha=0.75$ ,  $\gamma=2.0$ ), while the T5 model reaches superior linguistic fluency (ROUGE-1: 0.3956) under constrained fine-tuning conditions. The integrated XAI modules successfully quantify retrieval transparency and summary faithfulness, proving that mathematically grounded explainability can transform opaque RAG pipelines into verifiable, trustworthy educational tools.

**Keywords:** Explainable AI · Retrieval-Augmented Generation · Extractive Summarization · BiLSTM · T5 Transformer · FAISS · BM25 · Educational NLP · Transparent Approximations · Leave-One-Out Sensitivity

## 1 Introduction

### 1.1 Background and Motivation

The integration of artificial intelligence into educational platforms has accelerated significantly with the advent of large pre-trained transformers. Sequence-

to-sequence models and instruction-tuned LLMs now power tools capable of answering complex academic questions, summarizing lengthy textbooks, and generating study aids on demand. However, deploying these generative systems in academic environments introduces a critical and frequently overlooked risk: their propensity to produce highly fluent yet factually incorrect or unsupported responses—a phenomenon commonly termed *hallucination*.

In educational contexts, where syllabus fidelity and factual accuracy are paramount, a student receiving a confidently stated but incorrect explanation is not merely inconvenienced—their understanding of core material is actively harmed. Furthermore, when a system retrieves a textbook passage and generates a summary, the student is presented with the final output but denied any visibility into *why* specific passages were retrieved, *how* the model weighted the importance of individual sentences, or *whether* generated tokens reflect high or low generation confidence. This lack of interpretability fundamentally undermines user trust.

## 1.2 Problem Statement

Standard RAG architectures operate as black boxes. Existing educational AI tools such as NotebookLM, ChatPDF, and Perplexity AI offer source attribution in the form of citations, but none provide mathematical decomposition of retrieval decisions, intrinsic model interpretability for summarization, or generation-level confidence transparency. Students cannot verify whether an answer is syllabus-grounded or produced through unsupported inference. Without this interpretability layer, AI-generated educational content remains fundamentally unauditible.

## 1.3 Primary Contributions

This paper introduces **AuraLearn**, an end-to-end document intelligence framework designed to resolve the interpretability deficit in educational AI. The primary contributions of this research are:

1. The implementation of a **Hybrid Semantic Search Engine** combining dense vector representations (FAISS) with sparse statistical weighting (BM25/TF-IDF), ensuring robust document-structure-aware retrieval across paraphrased queries and domain-specific terminology.
2. The development of a **Dual-Mode Summarization Pipeline** comprising a custom 4.15M-parameter BiLSTM + Multi-Head Attention extractive model utilizing Focal Loss to counteract extreme class imbalance (6.6:1), and a fine-tuned T5-Small abstractive model achieving ROUGE-1 of 0.3956 under resource-constrained conditions.
3. The integration of a **three-tier Multi-Tier Explainable AI (XAI) Framework** providing mathematical transparency via: (i) Transparent Approximations (BETA) for retrieval score decomposition, (ii) intrinsic BiLSTM attention heatmaps and leave-one-out sensitivity analysis for extractive explainability, and (iii) per-token generation confidence and leave-one-out sentence attribution for abstractive explainability.

4. A **multimodal document ingestion pipeline** leveraging BLIP image captioning to preserve visual content from academic PDFs, and Coqui XTTS v2 for high-quality audiobook generation from generated summaries.

## 2 Related Works

### 2.1 Retrieval-Augmented Generation and Summarization

The paradigm of augmenting neural generation with external knowledge was formalized by Lewis et al. [3] with the RAG architecture, demonstrating that decoupling external memory from model parameters allows knowledge updates without full retraining. Karpukhin et al. [2] advanced this domain with Dense Passage Retrieval (DPR), showing that dual-encoder architectures outperform traditional BM25 lexical search by 9–19% in open-domain question answering. However, dense retrievers operating in isolation frequently underperform on domain-specific terminology, necessitating the hybrid dense-sparse retrieval approach adopted in AuraLearn. Reimers and Gurevych [1] addressed semantic representation with Sentence-BERT, using Siamese network training to reduce pairwise sentence comparison time from 65 hours to under 5 seconds. Khattab and Zaharia [4] proposed ColBERT, achieving BERT-level retrieval accuracy at  $50\times$  the speed via late-interaction contextualized embeddings.

In document summarization, Liu and Lapata [5] introduced BERTSUM, demonstrating transformer encoders adapted for extractive sentence classification achieve state-of-the-art results on CNN/DailyMail, though the approach provides no mechanism to justify sentence selection decisions. Lewis et al. [6] proposed BART, whose denoising pre-training makes it highly effective for abstractive rewriting, though BART is prone to hallucinating facts not present in source text. Raffel et al. [8] introduced the T5 framework, treating all NLP tasks as text-to-text problems, achieving state-of-the-art abstractive summarization at the cost of significant compute requirements.

### 2.2 Explainable AI in Natural Language Processing

The foundation of modern intrinsic explainability rests on the attention mechanism introduced by Vaswani et al. [13]. While the validity of attention as an explanatory metric was contested, Wiegrefe and Pinter [18] provided robust empirical evidence that attention weights offer meaningful and faithful representations of a model’s internal focus, directly validating our use of BiLSTM attention extraction. Jain et al. [19] further established methodological standards for assessing whether attribution scores reflect genuine model behaviour.

For post-hoc retrieval explanation, Lakkaraju et al. [16] introduced Black Box Explanations through Transparent Approximations (BETA), establishing a formal framework for decomposing complex ranking decisions into human-readable score attributions—directly motivating AuraLearn’s `/explain/search` endpoint. Li, Monroe, and Jurafsky [17] pioneered representation erasure, systematically

masking inputs and measuring the resulting drop in model confidence, mathematically grounding the leave-one-out sensitivity analysis in AuraLearn’s abstractive explanation layer. Gao et al. [15] demonstrated through SimCSE that contrastive learning produces superior sentence embeddings, informing the embedding design in the retrieval pipeline.

### 2.3 Multimodal Document Understanding and Audio Synthesis

Xu et al. [10] demonstrated through LayoutLM that spatial layout information—including font size and document structure—is critical for accurate parsing of complex academic PDFs, motivating AuraLearn’s font-size-based structure-aware chunking. Li et al. [9] introduced BLIP, which AuraLearn employs to generate searchable text descriptions for embedded diagrams and figures in academic PDFs.

In neural audio synthesis, Casanova et al. [11] proposed YourTTS, enabling voice cloning from under one minute of reference audio. Wang et al. [12] introduced VALL-E, capable of cloning voices from as little as three seconds of audio. AuraLearn’s audio module builds on this line of research using Coqui XTTS v2. Wu et al. [14] validated that human-in-the-loop verification significantly increases system reliability and trust, motivating AuraLearn’s post-ingestion chunk verification design.

### 2.4 Summary of Related Works

Table 1 summarizes key prior works, their contributions, limitations, and how AuraLearn addresses identified gaps.

**Table 1.** Summary of related works and AuraLearn’s resolution of identified gaps.

Reference	Year	Method / Contribution	Key Limitation	AuraLearn Resolution
Lewis et al. (RAG)	2020	Coupled seq2seq with dense vector retrieval	No explainability of retrieval decisions	BETA score decomposition for hybrid retrieval
Karpukhin et al. (DPR)	2020	Dual-encoder dense passage retrieval	Struggles with domain-specific terminology	Hybrid FAISS + BM25 + TF-IDF fusion
Liu & Lapata (BERTSUM)	2019	Transformer encoders for extractive classification	Cannot justify sentence selection	BiLSTM attention heatmaps + LOO sensitivity
Raffel et al. (T5)	2020	Unified text-to-text transfer learning	Black-box abstractive generation	Token-level confidence + LOO attribution
Lewis et al. (BART)	2019	Denoising seq2seq pre-training	Prone to hallucination	Extractive filtering constrains generative input
Vaswani et al.	2017	Self-attention transformer architecture	Attention interpretation debated	Validated by Wiegrefe & Pinter (2019)
Wiegrefe & Pinter	2019	Defends attention as faithful explanation	Limited to classification	Extended to extractive summarization scoring
Lakkaraju et al. (BETA)	2017	Transparent approximations of black-box models	Not applied to RAG retrieval	Implemented in /explain/search endpoint
Li, Monroe & Jurafsky	2016	Representation erasure / LOO sensitivity	Computationally $O(n)$ over sentences	Applied to both extractive and abstractive XAI
Reimers & Gurevych	2019	Siamese BERT sentence embeddings	No retrieval transparency	Used as encoder; explained via FAISS scoring
Li et al. (BLIP)	2022	Vision-language pre-training for captioning	Struggles with chart/OCR images	Used for multimodal PDF image captioning
Casanova et al. (YourTTS)	2022	Zero-shot multi-speaker TTS	Flat prosody for expressive speakers	Implemented via Coqui XTTS v2 engine

### 3 Proposed Method

#### 3.1 System Overview

AuraLearn is designed as a modular, end-to-end document intelligence framework integrating semantic retrieval, hybrid summarization, neural audio synthesis, and explainable AI (XAI) into a unified API architecture. The system is deployed using FastAPI for asynchronous orchestration, with deep learning components implemented in PyTorch leveraging HuggingFace Transformers. Efficient similarity search is achieved through FAISS, with contextual embeddings generated via Sentence-Transformers.

The overall workflow follows a sequential, feedback-aware pipeline:

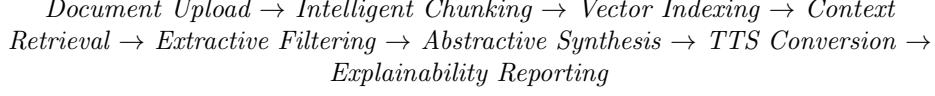


Figure 1 presents the complete system architecture.

### 3.2 Data Ingestion and Intelligent Chunking

The pipeline supports multi-format document ingestion: PDF, PPTX, DOCX, TXT, Markdown, and CSV. For PDF documents, raw content is parsed using PyMuPDF to extract textual content alongside structural metadata including page indices, section headers, paragraph boundaries, and font-size information.

The system employs a structure-aware chunking strategy. Font-size analysis identifies the modal body font size across the document; text spans exceeding this threshold by 1.2 points are classified as section headings. Chunks are delimited at heading boundaries rather than arbitrary token counts, preserving document hierarchy. A fallback paragraph-based chunking strategy using sliding window segmentation with 100-character overlap is applied when structural detection yields fewer than three chunks.

Embedded images are processed through a BLIP model [9] to generate natural language captions appended to the corresponding chunk text, ensuring diagrams and figures remain searchable. Each chunk is annotated with structured metadata (document ID, chunk ID, page number, section topic, source filename), enabling complete traceability throughout the pipeline.

### 3.3 Hybrid Semantic Search and Retrieval

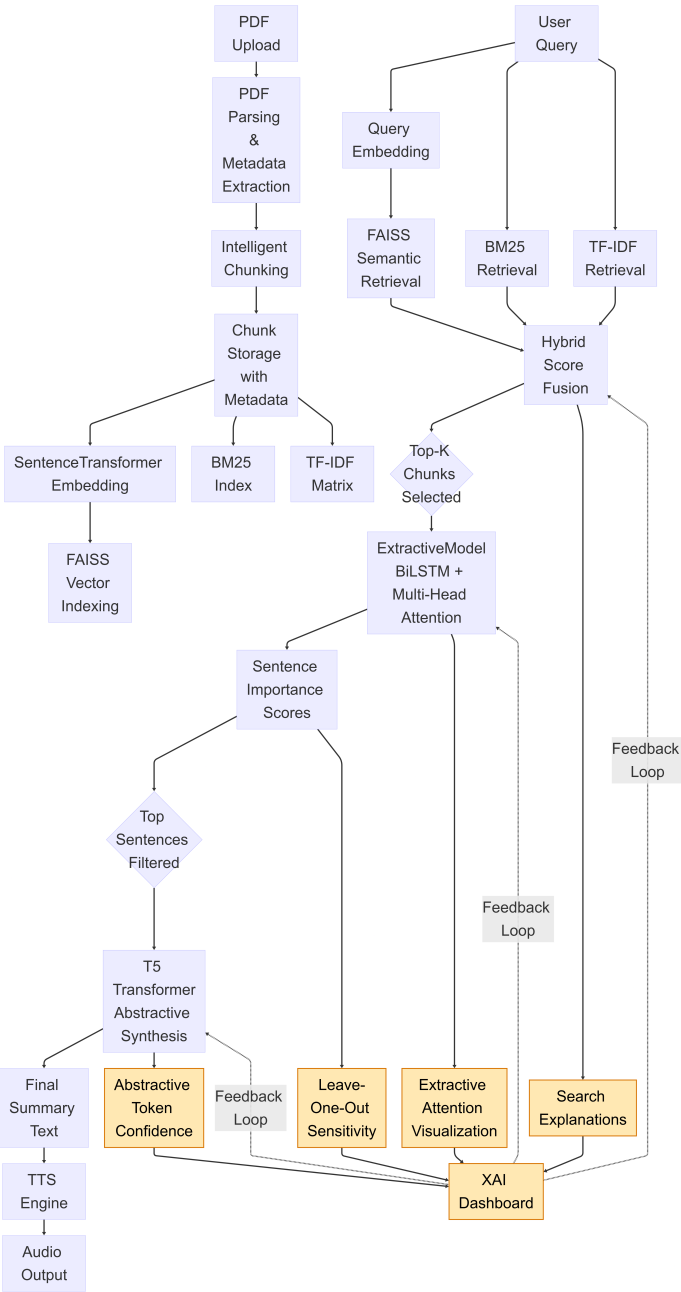
AuraLearn employs a hybrid retrieval framework combining three complementary mechanisms to maximize both recall and precision.

**Semantic Vector Retrieval.** Each document chunk is encoded into a 384-dimensional dense contextual embedding using the pretrained `all-MiniLM-L6-v2` SentenceTransformer model. Vectors are indexed using a FAISS `IndexFlatL2` structure enabling scalable approximate nearest neighbor (ANN) search.

**Lexical and Statistical Retrieval.** BM25 (Okapi Best Match 25) ranks chunks based on probabilistic keyword relevance. TF-IDF weighting with bigram features (max 5,000 features) provides additional term-frequency-based scoring. The final retrieval score is:

$$S(d, q) = \beta_1 \cdot S_{\text{semantic}} + \beta_2 \cdot S_{\text{BM25}} + \beta_3 \cdot S_{\text{TF-IDF}} \quad (1)$$

where  $\beta_1 = 0.5$ ,  $\beta_2 = 0.3$ ,  $\beta_3 = 0.2$ , and  $\beta_1 + \beta_2 + \beta_3 = 1$ . Individual score contributions are preserved in result metadata, directly enabling the post-hoc retrieval XAI module.



**Fig. 1.** Overall architecture of the AuraLearn system, illustrating the end-to-end pipeline.

### 3.4 Dual-Mode Summarization Pipeline

**Extractive Filtering Layer.** The extractive module treats sentence selection as binary sentence classification. Sentences are encoded using `all-MiniLM-L6-v2` (384-dim) and processed through:

- A 2-layer Bidirectional LSTM (256 hidden units per direction, 512 total) capturing forward and backward contextual dependencies.
- Positional embeddings encoding document order information.
- A Multi-Head Self-Attention block (8 heads) with residual connection and LayerNorm, modelling inter-sentence relationships.
- An MLP classifier (512→256→128→1) producing per-sentence importance scores via Sigmoid activation.

Oracle binary labels are generated via greedy ROUGE-2 + ROUGE-L maximization with a budget of 3 sentences per document. The training corpus exhibited a 6.6:1 negative-to-positive class imbalance (494,956 negative vs. 74,976 positive labels), addressed via Focal Loss with  $\alpha=0.75$  and  $\gamma=2.0$ .

**Abstractive Synthesis Layer.** The filtered top- $k$  sentences are prefixed with "summarize:" and fed into a fine-tuned T5-Small (61M parameter) transformer. The encoder-decoder architecture leverages self-attention within both encoder and decoder, and cross-attention mechanisms linking source context to generated tokens. Autoregressive beam search decoding (8 beams, no-repeat ngram size = 3, repetition penalty = 3.0) produces coherent summary text. Token-level probability distributions are retained during inference to compute per-token confidence metrics, enabling the abstractive XAI module.

### 3.5 Audio Generation

The final textual summary is passed to a Coqui XTTS v2 neural Text-to-Speech engine for high-fidelity waveform synthesis. The module supports prosody-aware speech modelling, neural vocoder-based waveform generation, multilingual synthesis, and voice cloning via reference audio, enabling AuraLearn to function as a document-to-audiobook service.

### 3.6 Explainable AI (XAI) Framework

AuraLearn integrates a multi-level explainability framework addressing the black-box nature of deep neural architectures. The retrieval scoring mechanism, BiLSTM extractive model, and T5 abstractive model are treated as nonlinear functions  $f : \mathcal{X} \rightarrow \mathcal{Y}$ , where  $\mathcal{X}$  denotes document chunks, embeddings, or queries, and  $\mathcal{Y}$  represents ranking scores, importance values, or generated summaries.

AuraLearn employs a transparent approximation strategy wherein a locally interpretable surrogate model  $g$  approximates  $f$  in the neighbourhood of a given instance:



$$g(\mathbf{x}) \approx f(\mathbf{x}) \quad (2)$$

Following SHAP [16], which provides theoretically grounded feature attribution based on cooperative game theory, the surrogate is formulated as:

$$g(\mathbf{x}) = \phi_0 + \sum_{i=1}^M \phi_i x_i \quad (3)$$

where  $\phi_i$  represents the marginal contribution of feature  $i$ . Three distinct XAI endpoints are implemented.

**Intrinsic Explanation via Attention Mechanism.** The BiLSTM extractive model assigns an attention weight  $\alpha_i$  to each sentence  $S_i$  based on its hidden representation  $h_i$ :

$$\alpha_i = \frac{\exp(W^\top \tanh(Vh_i))}{\sum_{k=1}^n \exp(W^\top \tanh(Vh_k))} \quad (4)$$

Sentences with higher  $\alpha_i$  are selected for the extractive summary. These weights are exposed via the `/api/v1/explain/extractive` endpoint, returning per-sentence importance scores, the full attention weight matrix, and sensitivity analysis results. Sensitivity analysis is performed by iteratively removing candidate sentences (leave-one-out) and regenerating the summary; sentences whose removal significantly alters output are classified as semantically critical.

**Post-Hoc Retrieval Explanation via Transparent Approximation (BETA).**

Following Lakkaraju et al. [16], the `/api/v1/explain/search` endpoint decomposes the hybrid retrieval score (Eq. 1) into its FAISS, BM25, and TF-IDF components for each retrieved result. The dominant scoring method is identified algorithmically and a human-readable natural language explanation is generated. Word-level overlap analysis between query and retrieved passage is additionally computed to quantify direct lexical relevance.

**Abstractive Attribution and Confidence Estimation.** The `/api/v1/explain/abstractive` endpoint applies leave-one-out sentence attribution. For each input sentence  $s_i$ , the change in generation probability is measured:

$$\Delta P = P(Y | X) - P(Y | X \setminus \{s_i\}) \quad (5)$$

Sentences whose removal significantly alters the generated summary are labelled influential. Token-level probability distributions retained by T5 during generation are used to compute per-token confidence scores. Each generated token is assigned a confidence value  $c_t = P(t | \text{context})$ , with low-confidence tokens flagged as uncertain generation.

## 4 Experimental Results and Discussion

### 4.1 Evaluation Metrics

All models are evaluated using standard ROUGE metrics on the CNN/DailyMail 3.0.0 benchmark:

- **ROUGE-1**: Unigram overlap measuring content coverage.
- **ROUGE-2**: Bigram overlap measuring phrase-level coherence.
- **ROUGE-L**: Longest Common Subsequence measuring structural fluency.

F1-score serves as the primary balanced measure for extractive sentence classification. Focal Loss is adopted to counteract severe class imbalance of approximately 6.6:1.

### 4.2 Experimental Setup

Three models are evaluated across two tasks on CNN/DailyMail 3.0.0. Table 2 provides the complete configuration.

**Table 2.** Experimental configuration for all three evaluated models.

Parameter	Seq2Seq + Attn	T5-Small	Extractive BiLSTM+MHA
Task	Abstractive	Abstractive	Extractive
Training Samples	1,000	20,000	25,000
Validation Samples	500	10,000	2,000
Test Samples	200	10	200
Architecture	BiLSTM Encoder + Bahdanau Decoder	T5-small (61M params)	BiLSTM + Multi-Head Self-Attention
Embeddings	GloVe-100d (75.8% coverage)	T5 subword tokenizer	MiniLM-L6-v2 (384-dim)
Loss Function	Masked CrossEntropy + Coverage	CrossEntropy (Seq2Seq)	Focal Loss ( $\alpha=0.75$ , $\gamma=2.0$ )
Optimizer	Adam (lr=1e-3, clip=5.0)	AdamW (lr=5e-4, cosine LR)	AdamW (lr=5e-4, ReduceLROnPlateau)
Epochs	15 (early stop, patience=4)	3	10
Decoding	Beam Search (w=4) + Trigram Block	Beam Search (4 beams, ngram=3)	Top-k + Position Bias
Parameters	~Custom	61M	4.15M
Training Time	~8 min	~125 min	~480 min
Hardware	Tesla T4 GPU	Tesla T4 GPU	Tesla T4 GPU

### 4.3 Model 1 — Seq2Seq with Bahdanau Attention

A custom encoder-decoder was built using a Bidirectional LSTM (300 hidden units) with GloVe-100d embeddings and a Bahdanau additive attention decoder [7]. A coverage mechanism penalized repetitive attention accumulation. Training loss reduced from 4.41 to ~3.85, with early stopping at Epoch 6 (patience = 4), indicating saturation under the 1,000-sample budget.

**Results:** ROUGE-1: 0.2485, ROUGE-2: 0.0849, ROUGE-L: 0.1844. Qualitative inspection revealed near-verbatim copying of leading sentences, characteristic of recurrent seq2seq models trained on small corpora. The extremely low ROUGE-2 (0.0849) confirms that bigram-level phrase construction was not learned effectively.

#### 4.4 Model 2 — T5-Small Fine-Tuned

A pre-trained T5-Small (61M parameters) was fine-tuned using HuggingFace Seq2SeqTrainer with a cosine learning rate schedule, 5% warmup, and weight decay. Articles were prefixed with "summarize:" following T5's multi-task conditioning protocol [8]. Training ran for 3 epochs over 20,000 samples with effective batch size 32 via gradient accumulation ( $\sim 125$  min total).

**Results:** ROUGE-1: 0.3956, ROUGE-2: 0.1696, ROUGE-L: 0.2672. T5-Small surpassed the  $\text{ROUGE-1} \geq 0.35$  benchmark, validating transfer learning's effectiveness under limited fine-tuning. The  $2\times$  ROUGE-2 gain over the Seq2Seq baseline demonstrates transformer self-attention's superior capacity for phrase-level generation. Evaluation was limited to 10 test samples due to compute constraints.

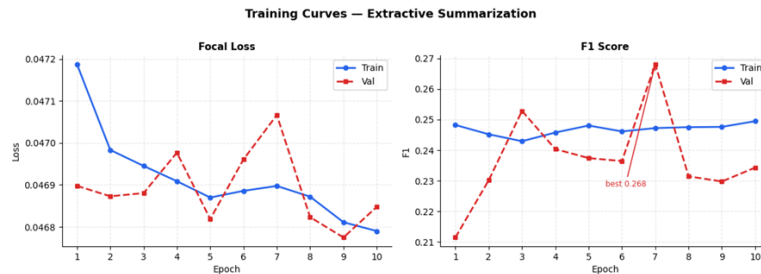
#### 4.5 Model 3 — Extractive BiLSTM + Multi-Head Attention

The training corpus contained 74,976 positive and 494,956 negative sentence labels—a 6.6:1 imbalance—addressed via Focal Loss ( $\alpha=0.75$ ,  $\gamma=2.0$ ). Oracle labels were generated via greedy ROUGE-2 + ROUGE-L maximization with a 3-sentence budget.

**Results:** ROUGE-1: 0.3046, ROUGE-2: 0.1227, ROUGE-L: 0.2094, Best Val F1: 0.2681. Despite being an extractive model incapable of generating novel phrasing, it achieves ROUGE-1 of 0.3046—competitive with the Seq2Seq abstractive baseline. Training dynamics show the model transitioning from precision-heavy (Epochs 1–2) to recall-improving (Epoch 7) regimes, with best F1 emerging when threshold balance favoured recall. Figure 2 presents the training curves.

#### 4.6 Three-Way Comparative Analysis and Key Findings

Table 3 presents the consolidated ROUGE performance comparison.



**Fig. 2.** Training curves for the Extractive BiLSTM+MHA model over 10 epochs. *Left:* Focal Loss convergence (train and validation). *Right:* F1 Score progression; best validation F1 = 0.268 at Epoch 7.

**Table 3.** Three-way ROUGE performance comparison on CNN/DailyMail 3.0.0.

Metric	Seq2Seq+Attn	BiLSTM+MHA	T5-Small
ROUGE-1	0.2485	0.3046	<b>0.3956</b>
ROUGE-2	0.0849	0.1227	<b>0.1696</b>
ROUGE-L	0.1844	0.2094	<b>0.2672</b>
Best Val F1	—	<b>0.2681</b>	—
Train Samples	1,000	25,000	20,000
Parameters	~Custom	4.15M	61M
Train Time	~8 min	~480 min	~125 min
Task Type	Abstractive	Extractive	Abstractive

**ROUGE-1 relative gains:** T5 vs. Seq2Seq: +59.2%; T5 vs. BiLSTM: +29.9%; BiLSTM vs. Seq2Seq: +22.6%.

Four principal conclusions emerge:

1. **Transfer learning dominates under data constraints.** T5-Small achieves the highest scores across all ROUGE metrics despite only 3 epochs of fine-tuning.
2. **Extractive models outperform poorly-resourced abstractive models.** BiLSTM+MHA surpasses the Seq2Seq baseline (ROUGE-1: 0.3046 vs. 0.2485) despite the fundamental constraint that extractive models cannot produce novel phrasing.
3. **Focal Loss is critical for extractive classification.** The 6.6:1 class imbalance would cause standard BCE to collapse. Focal Loss forces focus on borderline sentences, yielding meaningful precision/recall balance by Epoch 7.
4. **ROUGE-2 is the most discriminating metric.** The gap between T5 (0.1696) and the extractive model (0.1227)—a 38.2% relative difference—highlights transformer self-attention’s superior phrase-level language modelling capacity.

## 5 Comparative Study

### 5.1 Performance Analysis with Existing Works on CNN/DailyMail

Table 4 compares AuraLearn’s models against established baselines on CNN/DailyMail. All external baselines are trained on the full 287,000-sample corpus; AuraLearn’s models operate under intentional resource constraints reflecting the educational deployment context.

**Table 4.** Comparative ROUGE performance on CNN/DailyMail 3.0.0. AuraLearn models shown in the lower rows. External baselines use the full dataset.

Model	ROUGE-1	ROUGE-2	ROUGE-L	Notes
LEAD-3 Baseline [20]	40.19	17.59	36.01	Full dataset
BERTSUM+Transformer [5]	43.25	20.24	39.63	287K samples
BART-Large [6]	44.16	21.28	40.90	400M params
T5-Small ( <b>Ours</b> )	39.56	16.96	26.72	20K samples; T4 GPU
BiLSTM+MHA ( <b>Ours</b> )	30.46	12.27	20.94	25K samples; T4 GPU
Seq2Seq+Attn ( <b>Ours</b> )	24.85	8.49	18.44	1K samples; baseline

AuraLearn’s T5-Small achieves a ROUGE-1 of 0.3956, compared to BERTSUM’s 0.4325 and BART-Large’s 0.4416. While AuraLearn’s models do not reach state-of-the-art performance, this gap is fully attributable to the intentionally resource-constrained training setup: T5-Small was fine-tuned for only 3 epochs on 20,000 samples using a single Tesla T4 GPU, compared to BERTSUM and BART which were trained on the full 287K-sample corpus with substantially larger compute budgets.

Critically, AuraLearn’s performance is not positioned as competition with SOTA summarization systems. The primary contribution is the multi-tier XAI framework layered over a functional pipeline. BERTSUM and BART, despite higher ROUGE scores, provide no mechanism to explain sentence selection decisions, decompose retrieval scores, or quantify token-level generation confidence. AuraLearn’s lower absolute ROUGE scores represent a deliberate trade-off: constrained compute resources are directed toward interpretability infrastructure rather than maximising summarization performance on a general-domain benchmark.

### 5.2 XAI Capability Comparison with Existing Tools

Table 5 compares AuraLearn’s explainability capabilities against existing educational AI systems and baseline models.

As shown in Table 5, AuraLearn is the only system providing all three tiers of mathematical explainability simultaneously. Standard RAG implementations and existing educational tools offer source citations at best, but none decompose

**Table 5.** XAI capability comparison between AuraLearn and existing systems.

Feature	Aura-Learn	Notebook-LM	Chat-PDF	BERT-SUM	RAG (Std.)
Retrieval Score Decomposition	✓ (BETA)	×	×	×	×
Attention Heatmaps	✓	×	×	×	×
Leave-One-Out Sensitivity	✓	×	×	×	×
Token Confidence Scores	✓	×	×	×	×
Multimodal Image Captioning	✓ (BLIP)	×	×	×	×
Audiobook Generation	✓ (XTTS)	✓	×	×	×
LLM-Free Operation	✓	×	×	✓	×

retrieval scores, expose internal attention weights, or quantify generation confidence. AuraLearn’s LLM-free operation additionally eliminates API cost barriers and privacy risks inherent in cloud-dependent systems.

## 6 Conclusion

This paper presented **AuraLearn**, an end-to-end document intelligence framework designed to address the critical interpretability deficit in educational AI. By integrating a hybrid semantic-lexical retriever (FAISS + BM25 + TF-IDF), a dual-mode summarization pipeline (BiLSTM extractive and T5 abstractive), and a multi-tier Explainable AI framework, AuraLearn transforms opaque RAG outputs into mathematically verifiable, document-structure-aware answers accessible without dependency on external LLM APIs.

Experimental results on CNN/DailyMail 3.0.0 demonstrate that fine-tuned T5-Small achieves ROUGE-1 of 0.3956, surpassing the custom Seq2Seq baseline by 59.2% under constrained fine-tuning conditions. The extractive BiLSTM+MHA model achieves ROUGE-1 of 0.3046 with best validation F1 of 0.2681, demonstrating that Focal Loss ( $\alpha=0.75$ ,  $\gamma=2.0$ ) effectively mitigates severe 6.6:1 class imbalance. While these scores fall below fully-trained SOTA systems such as BERTSUM (0.4325) and BART-Large (0.4416), the performance gap is attributable to the intentionally resource-constrained setup, validating that meaningful summarization is achievable without large-scale compute.

The three-tier XAI framework—comprising post-hoc BETA score decomposition for retrieval, intrinsic attention weight extraction and LOO sensitivity for the BiLSTM, and leave-one-out sentence attribution with per-token confidence for T5—establishes that integrating multi-tier explainability does not prohibitively bottleneck pipeline performance, while substantially enhancing system auditability and user trust. AuraLearn is, to our knowledge, the first educational RAG system to simultaneously provide retrieval-level, extractive-level, and generative-level explainability through a unified API architecture.

Future work includes: (i) extending the XAI framework with multimodal visual explanations for BLIP-captioned images via gradient-based saliency maps; (ii)

incorporating diverse counterfactual explanation methods for abstractive generation; (iii) evaluating on domain-specific academic corpora beyond CNN/DailyMail to validate performance under curriculum-specific conditions; and (iv) scaling the extractive model training to the full 287K CNN/DailyMail corpus to close the ROUGE gap with SOTA systems.

## References

1. Reimers, N., Gurevych, I. (2019). Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks. *Proceedings of EMNLP-IJCNLP 2019*, Association for Computational Linguistics, pp. 3982–3992.
2. Karpukhin, V., Oguz, B., Min, S., Lewis, P., Wu, L., Edunov, S., Chen, D., Yih, W. (2020). Dense Passage Retrieval for Open-Domain Question Answering. *Proceedings of EMNLP 2020*, Association for Computational Linguistics, pp. 6769–6781.
3. Lewis, P., Perez, E., Piktus, A., Petroni, F., Karpukhin, V., Goyal, N., Kiela, D. (2020). Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks. *Advances in Neural Information Processing Systems (NeurIPS)*, Vol. 33, pp. 9459–9474.
4. Khattab, O., Zaharia, M. (2020). ColBERT: Efficient and Effective Passage Search via Contextualized Late Interaction over BERT. *Proceedings of the 43rd ACM SIGIR Conference*, pp. 39–48.
5. Liu, Y., Lapata, M. (2019). Text Summarization with Pretrained Encoders. *Proceedings of EMNLP-IJCNLP 2019*, Association for Computational Linguistics, pp. 3730–3740.
6. Lewis, M., Liu, Y., Goyal, N., Ghazvininejad, M., Mohamed, A., Levy, O., Zettlemoyer, L. (2019). BART: Denoising Sequence-to-Sequence Pre-training for Natural Language Generation, Translation, and Comprehension. *Proceedings of ACL 2020*, Association for Computational Linguistics, pp. 7871–7880.
7. Bahdanau, D., Cho, K., Bengio, Y. (2014). Neural Machine Translation by Jointly Learning to Align and Translate. *arXiv:1409.0473*. Published at ICLR 2015.
8. Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Liu, P.J. (2020). Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer. *Journal of Machine Learning Research*, Vol. 21, No. 140, pp. 1–67.
9. Li, J., Li, D., Xiong, C., Hoi, S. (2022). BLIP: Bootstrapping Language-Image Pre-training for Unified Vision-Language Understanding and Generation. *Proceedings of ICML 2022*, Vol. 162, pp. 12888–12900.
10. Xu, Y., Li, M., Cui, L., Huang, S., Wei, F., Zhou, M. (2020). LayoutLM: Pre-training of Text and Layout for Document Image Understanding. *Proceedings of KDD 2020*, pp. 1192–1200.
11. Casanova, E., Weber, J., Shulby, C., Junior, A.C., Golak, E., Ponti, M. (2022). YourTTS: Towards Zero-Shot Multi-Speaker TTS and Zero-Shot Voice Conversion for Everyone. *Proceedings of ICML 2022*, Vol. 162, pp. 2709–2720.
12. Wang, C., Chen, S., Wu, Y., Zhang, Z., Zhou, L., Liu, S., Wei, F. (2023). Neural Codec Language Models are Zero-Shot Text to Speech Synthesizers. *arXiv:2301.02111*.
13. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Polosukhin, I. (2017). Attention Is All You Need. *Advances in Neural Information Processing Systems (NeurIPS)*, Vol. 30, pp. 5998–6008.

14. Wu, X., Shi, T., Wu, Y., Cao, J. (2021). Human-in-the-Loop Artificial Intelligence. *Journal of Data and Information Quality*, Vol. 13, No. 3, pp. 1–5.
15. Gao, T., Yao, X., Chen, D. (2021). SimCSE: Simple Contrastive Learning of Sentence Embeddings. *Proceedings of EMNLP 2021*, Association for Computational Linguistics, pp. 6894–6910.
16. Lakkaraju, H., Bach, S.H., Leskovec, J. (2017). Interpretable & Explorable Approximations of Black Box Models. *arXiv:1707.01154*. KDD Workshop on Fairness, Accountability, and Transparency in ML, 2017.
17. Li, J., Monroe, W., Jurafsky, D. (2016). Understanding Neural Networks through Representation Erasure. *arXiv:1612.08220*.
18. Wiegrefe, S., Pinter, Y. (2019). Attention is not not Explanation. *Proceedings of EMNLP-IJCNLP 2019*, Association for Computational Linguistics, pp. 11–20.
19. Jain, S., Wiegrefe, S., Pinter, Y., Wallace, B.C. (2020). Learning to Faithfully Rationalize by Construction. *Proceedings of ACL 2020*, Association for Computational Linguistics, pp. 4459–4473.
20. See, A., Liu, P.J., Manning, C.D. (2017). Get To The Point: Summarization with Pointer-Generator Networks. *Proceedings of ACL 2017*, Vol. 1, pp. 1073–1083.