

Delta-Gated Incremental Multi-Forward-Pass Modeling for Robust Multimodal Classroom Video Understanding

Anonymous CVPR submission

Paper ID ****

Abstract

001 *Multimodal transformers are attractive options for the*
002 *analysis of human activity in the classroom, but real-world*
003 *classroom data often contain missing or misaligned modal-*
004 *ities, making robust multimodal learning challenging. In*
005 *this work, we propose a **Delta-Gated Multi-Forward-Pass***
006 *(**DG-MFP**) **Longformer** for robust multimodal classroom*
007 *discourse understanding. The model treats the transcript*
008 *as the primary modality and models audio and video as in-*
009 *cremental sources of information. A shared encoder per-*
010 *forms multiple forward passes with different modality masks*
011 *(text, text+audio, text+video, text+audio+video), allow-*
012 *ing modality contributions to be isolated through differ-*
013 *ences between representations. These increments are fused*
014 *through class-specific delta gates that modulate modality*
015 *contributions relative to the text baseline. To evaluate ro-*
016 *bustness under realistic classroom conditions, we introduce*
017 *controlled missing-modality and cross-modal misalignment*
018 *tests. Experiments on the Artificial Intelligence for Advanc-*
019 *ing Instruction at Scale (AIAIS) dataset show that the pro-*
020 *posed method consistently improves overall F1 over stan-*
021 *dard multimodal fusion baselines and exhibits substantially*
022 *stronger robustness under missing or misaligned modal-*
023 *ities. Further analysis of the learned gates reveals inter-*
024 *pretable, task-specific patterns of modality activity, high-*
025 *lighting how multimodal signals provide complementary in-*
026 *formation when transcript evidence alone is insufficient.*

027 1. Introduction

028 Understanding classroom interactions requires reasoning
029 over multiple modalities such as spoken conversation and
030 visual behavior. Recent advances in transformer archi-
031 tectures have enabled powerful multimodal representation
032 learning across vision, language, and audio signals [5, 12,
033 14]. These models are increasingly used for educational
034 video analysis and classroom discourse understanding.

035 However, real classroom recordings often contain miss-

ing, noisy, or misaligned modalities due to recording con- 036
ditions. In this context, misalignment refers to imperfect 037
temporal alignment between the transcript segment and the 038
associated audio or video streams. Audio tracks may be cor- 039
rupted, video streams may be incomplete, and multimodal 040
signals are often imperfectly synchronized. Despite these 041
challenges, most multimodal learning systems are evaluated 042
under ideal conditions where all modalities are available 043
and perfectly aligned. As a result, it remains unclear how 044
robust these models are when applied to realistic classroom 045
environments where modalities may be unreliable. 046

A key challenge in multimodal learning is understand- 047
ing how different modalities influence the final predic- 048
tion. Standard multimodal fusion methods typically com- 049
bine modalities within a single forward pass. Although this 050
approach can improve predictive performance, it makes it 051
difficult to isolate modality contributions or analyze model 052
behavior when modalities are missing or misaligned. 053

To address this problem, we propose a **Delta-Gated** 054
Multi-Forward-Pass Longformer for robust multimodal 055
classroom discourse understanding, as illustrated in Fig- 056
ure 1. The central idea is to treat the transcript as the pri- 057
mary modality and to model audio and video as incremental 058
sources of information. 059

A shared encoder performs multiple forward passes with 060
different modality masks, producing representations for 061
text-only, text+audio, text+video, and text+audio+video in- 062
puts. Because the encoder parameters are shared across 063
passes, differences between these representations reflect the 064
additional information contributed by each modality. 065

This formulation explicitly separates the baseline tex- 066
tual signal from modality-specific increments, allowing the 067
model to estimate both the incremental contribution and the 068
reliability of each modality. The resulting increments are 069
fused through a delta-gated mechanism that models modal- 070
ity contributions relative to the text baseline while preserv- 071
ing interpretability through class-specific gating weights. 072

In addition to the model design, we introduce a diag- 073
nostic evaluation protocol for multimodal robustness under 074
realistic classroom conditions. The protocol includes con- 075

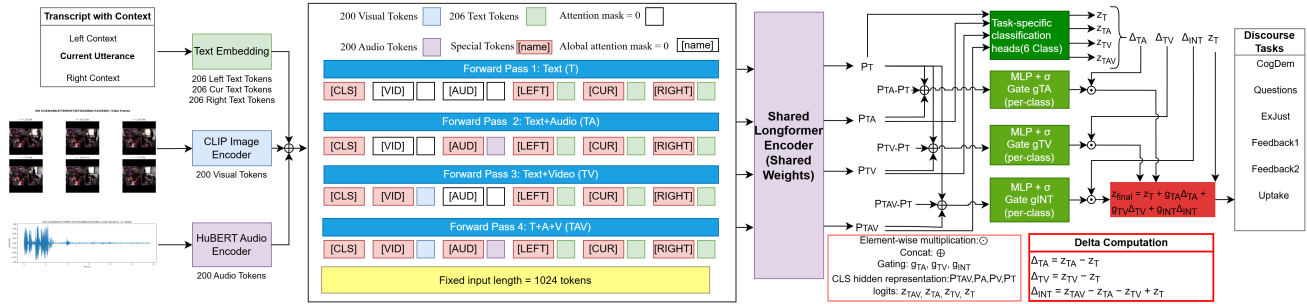


Figure 1. Overview of the proposed **DG-MFP** architecture. A unified multimodal token sequence is constructed from transcript context, video frames, and audio signals. A shared Longformer encoder is executed four times under different modality masks, producing representations for text-only (T), text+audio (TA), text+video (TV), and text+audio+video (TAV). Task-specific classification heads produce logits for each pass, while per-class delta gates modulate the incremental contributions of audio, video, and higher-order multimodal interaction terms.

076 trolled missing-modality scenarios and cross-modal mis- 110
 077 alignment tests that simulate common issues in classroom 111
 078 recordings. This design is particularly well suited to edu- 112
 079 cational video analysis, where transcripts often provide the 113
 080 most reliable signal while audio and visual streams may be 114
 081 noisy, partially missing, or imperfectly aligned in real class- 115
 082 room recordings. Furthermore, the explicit modeling of 116
 083 modality contributions provides interpretable insights into 117
 084 how different signals support discourse prediction, which is 118
 085 important for educational research and learning analytics. 119

086 **Our contributions are summarized as follows:**

- 087 1. We propose a **Delta-Gated Multi-Forward-Pass Long-** 120
 088 **former** that treats the transcript as the primary modal- 121
 089 ity and models audio and video as incremental contribu- 122
 090 tions through structured logit increments and class- 123
 091 specific gates. 124
 092 2. We introduce a **diagnostic evaluation protocol** for mul- 125
 093 timodal robustness that includes missing-modality and 126
 094 cross-modal misalignment stress tests. 127
 095 3. Experiments with classroom interaction data demon- 128
 096 strate improved robustness and provide interpretable in- 129
 097 sights into modality usage in discourse classification 130
 098 tasks. 131

099 2. Related Work

100 2.1. Multimodal Transformers

101 Transformer architectures have become a dominant frame- 136
 102 work for modeling long-range dependencies across modal- 137
 103 ities. The transformer architecture introduced in [12] has 138
 104 been widely adopted for language, vision, and multimodal 139
 105 learning, while pretrained models such as BERT [5] demon- 140
 106 strated the effectiveness of large-scale transformer-based 141
 107 representation learning. 142

108 Recent work extends transformers to multimodal set- 143
 109 tings. The Multimodal Transformer [11] introduced cross- 144

modal attention to jointly model language, visual, and 110
 acoustic signals in unaligned sequences. Other approaches 111
 such as VideoBERT [10] and vision–language models in- 112
 cluding CLIP [9] demonstrate the benefit of large-scale 113
 multimodal pretraining. More recent models such as 114
 Flamingo [1] and InternVideo [13] further explore large- 115
 scale vision–language and video representation learning. 116

Multimodal fusion strategies typically combine repre- 117
 sentations through early fusion, cross-modal attention, or 118
 gating mechanisms. Although these approaches improve 119
 predictive performance, they often make it difficult to ex- 120
 plicitly analyze how individual modalities contribute to the 121
 final prediction. 122

For long-sequence modeling, the Longformer [3] intro- 123
 duces sparse attention that enables efficient processing of 124
 long contexts. In this work, we build upon Longformer 125
 as the shared encoder for multimodal classroom discourse 126
 analysis. 127

128 2.2. Robust Multimodal Learning with Missing 129 Modalities

Many multimodal learning systems assume that all modal- 130
 ities are available during both training and inference. How- 131
 ever, in practice, real-world data often contain missing, 132
 corrupted, or unreliable modalities. Early work such as 133
 ModDrop [8] improves robustness by randomly removing 134
 modalities during training. 135

More recent studies analyze the behavior of multimodal 136
 transformers when one modality is absent. For example, Ma 137
 et al. [7] investigated the robustness of transformer-based 138
 multimodal models under missing-modality conditions and 139
 showed that model performance can depend strongly on the 140
 fusion strategy. 141

Most prior work focuses primarily on improving predic- 142
 tive performance. In contrast, our work explicitly models 143
 modality contributions and evaluates model behavior under 144

145	controlled missing-modality and cross-modal misalignment	194
146	conditions.	195
147	2.3. Multimodal Analysis in Educational Settings	196
148	Multimodal signals play an important role in understanding	197
149	classroom interactions and learning processes. Educational	198
150	research increasingly uses multimodal learning analytics to	199
151	analyze speech, gestures, and visual behaviors in instruc-	200
152	tional environments [4].	201
153	More broadly, multimodal machine learning integrates	202
154	heterogeneous signals such as language, vision, and audio	
155	[2]. Transformer-based approaches further enable learning	
156	from partially aligned multimodal sequences by modeling	
157	cross-modal dependencies [11].	
158	However, educational datasets often contain substantial	
159	noise, imperfect alignment, and incomplete signals due	
160	to real classroom recording conditions. These challenges	
161	make robustness particularly important for the practical im-	
162	plementation of multimodal models in educational settings.	
163	Our work contributes to this area by evaluating multi-	
164	modal transformers under realistic classroom conditions	
165	with missing modalities and temporal misalignment, while	
166	also analyzing how individual modalities contribute to dis-	
167	course prediction tasks.	
168	3. Dataset and Data Processing	
169	3.1. Data Source	
170	Our dataset is derived from the Artificial Intelligence for	
171	Advancing Instruction at Scale (AIAIS) dataset released as	
172	part of the AIAI Challenge [6], which contains classroom	
173	recordings for studying teacher–student interactions.	
174	Rather than modeling entire videos, we focus on	
175	transcript-centered discourse segments , where each sam-	
176	ple corresponds to a short unit of classroom conversation	
177	during elementary mathematics and language arts instruc-	
178	tion.	
179	We use a subset of 107 classroom videos , producing	
180	36,091 transcript-centered samples aligned with the cor-	
181	responding audio and video streams.	
182	3.2. Multimodal Inputs and Context Construction	
183	Each sample contains three modalities: classroom video	
184	frames, audio signals, and transcript text produced by au-	
185	tomatic speech recognition (ASR). The transcript segment	
186	serves as the primary unit of analysis , while audio and	
187	video provide complementary contextual information.	
188	Modalities are aligned using timestamps from the origi-	
189	nal ASR output. Although the text of the transcript was	
190	manually corrected, the timing information was not reanno-	
191	tated. Alignment is therefore approximate and occurs at the	
192	sentence level rather than precise frame synchronization, re-	
193	fecting realistic classroom recording conditions.	
	To capture conversational context, each sample includes	
	a context window consisting of left context, current tran-	
	script segment, and right context. The model processes	
	these multimodal inputs within a maximum sequence length	
	of 1024 tokens .	
	Visual features are extracted using CLIP ViT-B/32 frame	
	embeddings, while audio features are obtained using a Hu-	
	BERT encoder. The transcript text is tokenized and pro-	
	cessed by the transformer model.	
	3.3. Discourse Annotation Tasks and Dataset Statis-	
	tics	
	The objective is to classify instructional discourse segments	
	within classroom interactions. Each transcript-centered	
	sample is annotated using a hierarchical scheme consisting	
	of six discourse tasks with 19 subclasses (25 labels in to-	
	tal including a <i>None</i> label for each task) .	
	The six tasks capture different aspects of classroom in-	
	structional interaction. <i>CogDem</i> (Cognitive Demand) re-	
	fects the level of reasoning involved in the discourse, dis-	
	tinguishing between simple reporting of facts and higher-	
	level analytical reasoning. <i>Questions</i> categorizes the types	
	of teacher questions posed to students, such as open-ended	
	or closed-ended questions. <i>ExJust</i> (Explanation and Justi-	
	fication) identifies when participants request or provide ex-	
	planations and reasoning. <i>Feedback1</i> and <i>Feedback2</i> re-	
	present different types of teacher feedback and evaluation	
	of student responses. Finally, <i>Uptake</i> captures how teach-	
	ers incorporate or build on students’ ideas during classroom	
	interaction. Together, these tasks characterize key instruc-	
	tional behaviors such as questioning, reasoning, feedback,	
	and the integration of student contributions.	
	The final dataset contains 36,091 samples from 107	
	classroom videos . We use an 80/20 split with 28,787 train-	
	ing samples and 7,304 test samples . Each sample may con-	
	tain labels for multiple discourse tasks, but within each task	
	exactly one subclass is assigned (including <i>None</i>).	
	The dataset exhibits substantial class imbalance, with	
	dominant <i>None</i> labels and long-tailed subclass distributions	
	across tasks. This creates a challenging multi-task long-	
	tailed classification problem . To better reflect perfor-	
	mance on minority classes we adopt macro-F1 as the pri-	
	mary evaluation metric. Figure 2 visualizes the label distri-	
	bution across the six discourse tasks.	
	4. Method	
	4.1. Overview	
	Figure 1 illustrates the proposed DG-MFP framework for	
	multimodal classroom discourse classification.	
	Given a transcript-centered classroom sample, we con-	
	struct a unified multimodal token sequence containing tran-	
	script context, video tokens, and audio tokens. These tokens	

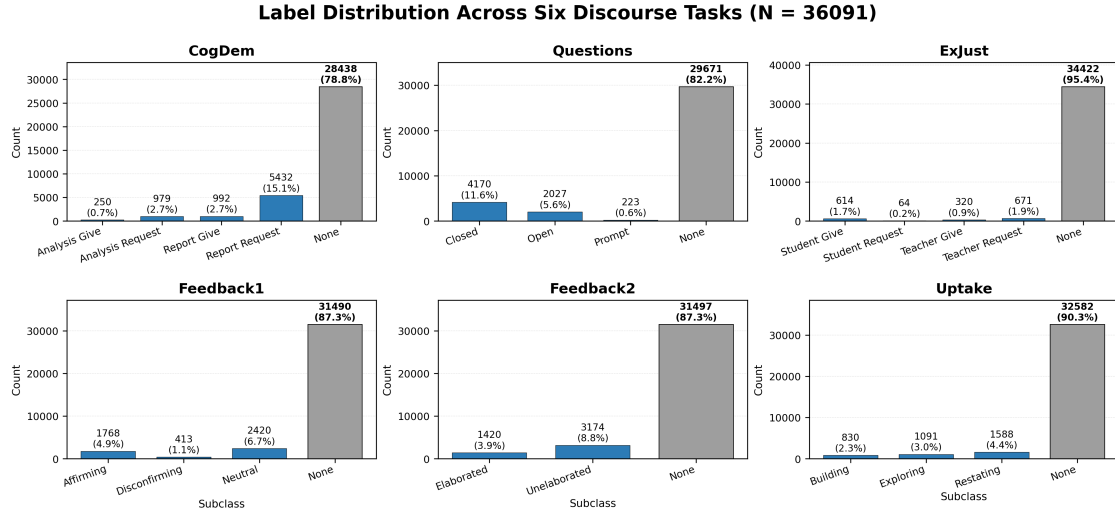


Figure 2. Label distribution across the six discourse tasks ($N = 36,091$ transcript-centered samples). Each task contains a dominant *None* class and several low-frequency subclasses, resulting in severe class imbalance and long-tailed label distributions. This motivates the use of macro-F1 for robustness experiments, as it better reflects performance on minority classes.

244 are processed by a **shared Longformer encoder** that is ex-
245 ecuted multiple times under different modality masks.

246 Unlike conventional multimodal models that fuse modal-
247 ities in a single forward pass, our method performs
248 four forward passes using the same encoder param-
249 eters: text-only (T), text+audio (TA), text+video (TV), and
250 text+audio+video (TAV). Each pass always includes tran-
251 script tokens, while additional modalities are introduced in-
252 crementally through modality masks. As a result, differ-
253 ences between passes reflect only the additional information
254 contributed by each modality.

255 Instead of directly predicting from the full multimodal
256 representation, we model how additional modalities mod-
257 ify a text-only baseline through delta increments. The final
258 logits are computed as

$$259 \quad z = z_T + g_{TA} \odot \Delta_{TA} + g_{TV} \odot \Delta_{TV} + g_{INT} \odot \Delta_{INT}, \quad (1)$$

260 where \odot denotes element-wise multiplication, Δ_{TA} and
261 Δ_{TV} represent modality-specific increments relative to the
262 text baseline, Δ_{INT} captures higher-order multimodal in-
263 teractions, and the gates g_{TA} , g_{TV} , g_{INT} control their con-
264 tributions. This formulation allows multimodal fusion to be
265 interpreted as a set of incremental information gains relative
266 to the baseline transcript prediction, with learned gates that
267 estimate the confidence of each information source.

268 4.2. Multimodal Representation and Encoder

269 Each sample contains three modalities: transcript text,
270 classroom video, and classroom audio. The text modality
271 includes the target transcript together with its left and right
272 discourse contexts (Section 3).

Visual and acoustic features are extracted using pre-
trained CLIP and HuBERT encoders and interpolated to
fixed-length token sequences. Each sample contains 200 vi-
sual tokens, 200 audio tokens, and three text segments (left
context, current transcript, and right context). The result-
ing multimodal sequence has a fixed length of **1024 tokens**
with modality markers [VID], [AUD], [LEFT], [CUR],
and [RIGHT].

Let the multimodal sequence be denoted by x . A shared
Longformer encoder $\text{Enc}(\cdot)$ processes the input under dif-
ferent modality masks:

$$284 \quad p_T = \text{Enc}(x; M_T), \quad (2)$$

$$285 \quad p_{TA} = \text{Enc}(x; M_{TA}), \quad (3)$$

$$286 \quad p_{TV} = \text{Enc}(x; M_{TV}), \quad (4)$$

$$287 \quad p_{TAV} = \text{Enc}(x; M_{TAV}), \quad (5)$$

288 where M_T enables transcript tokens only, M_{TA} enables
289 transcript and audio tokens, M_{TV} enables transcript and
290 video tokens, and M_{TAV} enables all modalities.

291 The encoder parameters are shared across forward
292 passes. The pooled representation is obtained from the final
293 hidden state of the [CLS] token, producing a text baseline
294 representation p_T , bimodal representations p_{TA} and p_{TV} ,
295 and the complete multimodal representation p_{TAV} .

296 4.3. Task Heads

297 The task is formulated as a **six-task discourse classifica-**
298 **tion** problem. For each discourse task t , we attach a task-
299 specific linear classifier $h^{(t)}(\cdot)$ to the pooled representation.

300 The same classifier is applied to all forward-pass repre-
301 sentations, producing logits

$$302 \quad z_T^{(t)} = h^{(t)}(p_T), \quad (6)$$

$$303 \quad z_{TA}^{(t)} = h^{(t)}(p_{TA}), \quad (7)$$

$$304 \quad z_{TV}^{(t)} = h^{(t)}(p_{TV}), \quad (8)$$

$$305 \quad z_{TAV}^{(t)} = h^{(t)}(p_{TAV}). \quad (9)$$

306 4.4. Delta-Gated Multimodal Fusion

307 Rather than directly using multimodal logits $z_{TAV}^{(t)}$, we de-
308 compose the prediction into modality-specific increments
309 relative to the text baseline.

Text–Audio Increment

$$310 \quad \Delta_{TA}^{(t)} = z_{TA}^{(t)} - z_T^{(t)} \quad (10)$$

$$311 \quad g_{TA}^{(t)} = \sigma\left(W_{TA}^{(t)}[p_T, p_{TA}, p_{TA} - p_T]\right) \quad (11)$$

$$312 \quad z^{(t)} = z_T^{(t)} + g_{TA}^{(t)} \odot \Delta_{TA}^{(t)} \quad (12)$$

Text–Video Increment

$$313 \quad \Delta_{TV}^{(t)} = z_{TV}^{(t)} - z_T^{(t)} \quad (13)$$

$$314 \quad g_{TV}^{(t)} = \sigma\left(W_{TV}^{(t)}[p_T, p_{TV}, p_{TV} - p_T]\right) \quad (14)$$

$$315 \quad z^{(t)} \leftarrow z^{(t)} + g_{TV}^{(t)} \odot \Delta_{TV}^{(t)} \quad (15)$$

Higher-Order Interaction

$$316 \quad \Delta_{INT}^{(t)} = z_{TAV}^{(t)} - z_{TA}^{(t)} - z_{TV}^{(t)} + z_T^{(t)} \quad (16)$$

$$317 \quad g_{INT}^{(t)} = \sigma\left(W_{INT}^{(t)}[p_T, p_{TAV}, p_{TAV} - p_T]\right) \quad (17)$$

$$318 \quad z^{(t)} \leftarrow z^{(t)} + g_{INT}^{(t)} \odot \Delta_{INT}^{(t)} \quad (18)$$

319 All gates are **task-specific and per-class**, allowing
320 modality contributions to vary across discourse labels.

321 4.5. Training Objective

322 For each discourse task t , we apply cross-entropy loss

$$323 \quad \mathcal{L} = \sum_{t \in \mathcal{T}} \text{CE}(z^{(t)}, y^{(t)}) \quad (19)$$

324 where \mathcal{T} denotes the set of six discourse tasks.

Table 1. Overall F1 on the full test set. Each result corresponds to the best checkpoint within 30 training epochs.

Model	S42	S43	S44
Longformer (T)	0.6588	0.6599	0.6636
EF Longformer (TAV)	0.6711	0.6676	0.6678
DG-MFP (Ours)	0.6777	0.6712	0.6687

5. Experiments 325

5.1. Baselines 326

We compare the proposed method with two baseline archi-
327 tectures designed to isolate the effect of the fusion strategy.
328

Text-only baseline. This baseline uses only the transcript
329 together with the left and right context. It shares the same
330 Longformer backbone and text formatting as our method
331 but removes audio and video input entirely. This setting
332 measures how much predictive signal is available from the
333 transcript context alone.
334

Standard multimodal early-fusion baseline. This base-
335 line performs multimodal fusion in a single forward pass.
336 Visual and acoustic tokens are concatenated with the tran-
337 script sequence using the same prompt-token layout as
338 our method, and the pooled Longformer representation is
339 passed to the same task-specific classification heads. Unlike
340 our method, this baseline does not perform multiple masked
341 forward passes or delta-gated fusion.
342

5.2. Implementation Details 343

All models are trained for up to **30 epochs**, and we report
344 the best checkpoint selected within this training schedule.
345 To account for training variance, each main model is trained
346 with **three random seeds (42, 43, 44)**.
347

For the main comparison, we report **overall F1 on the**
348 **entire test set**. For robustness evaluation, we additionally
349 report **macro-F1** because perturbations may disproportion-
350 ately affect low-support classes.
351

5.3. Main Results 352

Table 1 reports overall F1 on the full test set. 353

The text-only baseline already achieves strong perfor-
354 mance, indicating that the transcript context provides sub-
355 stantial predictive information for the classification of class-
356 room discourse. Standard early fusion further improves per-
357 formance, demonstrating that audio and visual signals pro-
358 vide useful complementary cues when integrated into the
359 model.
360

The proposed DG-MFP model achieves the best over-
361 all F1 across all three random seeds used during training.
362 For example, with the random seed set to 42, our model
363

Table 2. Bucket-level F1 comparison for seed 42.

Bucket	# Labels	Text-only	Early Fusion	Ours
None	6	0.9685	0.9684	0.9681
Large	7	0.7071	0.7055	0.7107
Medium	6	0.5566	0.5607	0.5519
Small	6	0.3822	0.4034	0.4599

improves F1 from 0.6711 with early fusion to 0.6777. Although the absolute improvement is modest, the gain is consistent across seeds, indicating that explicitly modeling modality increments and higher-order interaction effects yields a more stable multimodal fusion strategy than single-pass early fusion.

Because the dataset is long-tailed and only approximately aligned across modalities, we next analyze performance by label frequency and under controlled robustness perturbations.

5.4. Performance Across Label-Frequency Buckets

As the dataset is strongly long-tailed, we further analyze which types of labels benefit the most from multimodal fusion. Using the seed-42 models, we group discourse labels into four frequency buckets: *None*, *Large*, *Medium*, and *Small*.

Figure 3 and Table 2 report the average F1 within each bucket. The *None* bucket contains the dominant background labels, whereas the *Small* bucket contains the rarest and most challenging discourse subclasses.

Several trends are apparent. First, all models perform similarly on the *None* bucket, indicating that dominant labels are already easy to predict from the transcript context.

Second, early fusion provides only limited and inconsistent gains across frequency buckets.

In contrast, the proposed DG-MFP model shows its clearest advantage on the *Small* bucket, which contains the rarest and most challenging discourse subclasses. Our method improves the average F1 from 0.3822 for the text-only baseline and 0.4034 for early fusion to 0.4599, corresponding to absolute gains of +0.0777 and +0.0565, respectively.

These results suggest that multimodal cues are particularly useful for low-support discourse labels, where the transcript context alone is often insufficient for reliable prediction.

5.5. Robustness Under Missing and Mismatched Modalities

To evaluate robustness beyond clean evaluation, we perform three controlled perturbation tests at inference time. These tests measure how multimodal fusion strategies behave when additional modalities are missing or semantically misaligned with the target transcript.

Cross-sample misalignment. We keep the transcript fixed, but replace the associated video and audio streams with mismatched signals drawn from other test samples. This produces a semantically incorrect multimodal context while preserving the text input. We denote this setting as **baseline_miss_alignment**.

Missing modality with prompt retention. We zero out all video and audio token values while retaining the [VID] and [AUD] prompt tokens. This tests whether the model remains stable when modality structure is present, but modality content is absent. We denote this setting as **baseline_missing**.

Missing modality with prompt removal. We remove both modality tokens and prompt tokens, leaving only the text portion of the input sequence. This tests robustness when the multimodal input structure itself changes. We denote this setting as **baseline_missing_removeprompt**.

All robustness tests are performed at evaluation time using the same trained checkpoints as in the main experiments, without additional training or finetuning. We report **macro-F1** under clean evaluation and the three perturbation settings. The text-only baseline is omitted because it does not use audio or video inputs.

Table 3 shows that the proposed model is substantially more robust than standard early fusion. Under cross-sample misalignment, early fusion drops by 0.0222, whereas our model drops only by 0.0006. Under the missing modality with prompt retention, the drop is 0.0109 for early fusion but only 0.0012 for our method. Even in the most severe setting, where both modality tokens and prompts are removed, our model retains a substantially higher macro-F1 (0.4304 vs. 0.3522).

These results indicate that modeling audio and video as incremental contributions relative to a baseline of transcripts makes the model less sensitive to misleading or absent multimodal signals.

5.6. Analysis of Learned Delta Gates

To understand how multimodal signals contribute to prediction, we analyze the learned gate values for the three modality pathways: text-audio (g_{TA}), text-video (g_{TV}), and higher-order interaction (g_{INT}). Each gate is a per-class scalar in the range $[0, 1]$ that modulates the contribution of the corresponding logit increment.

Relative importance of modality paths. Across all discourse tasks, the learned values g_{TA} and g_{TV} generally fall in the range of approximately 0.4–0.5, while the higher-order interaction gate g_{INT} remains substantially smaller, typically around 0.2. Figure 4 shows the average gate values

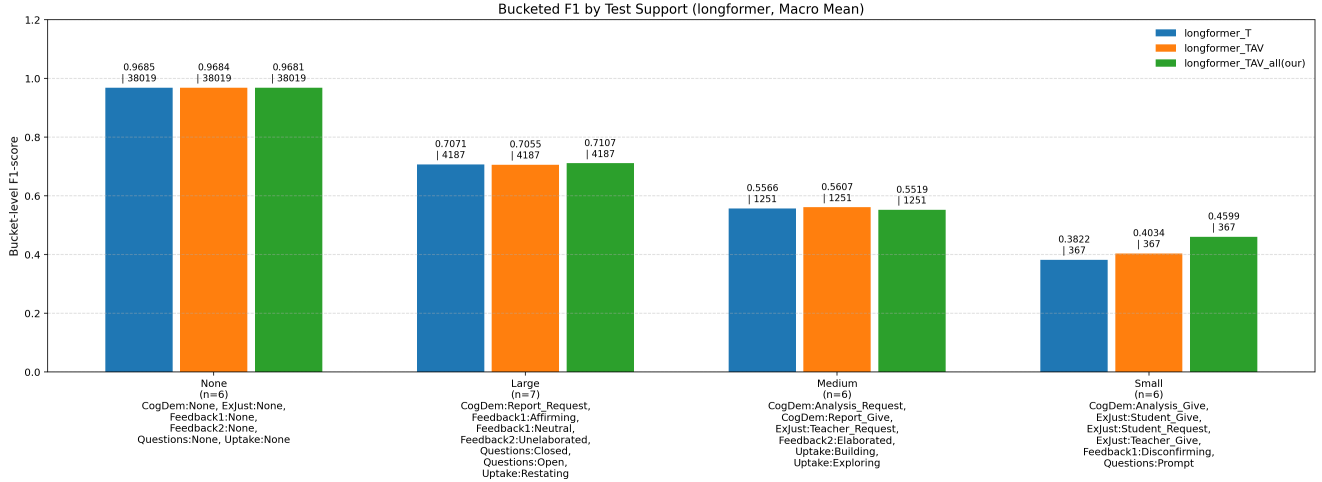


Figure 3. Bucketed F1 by test-set label frequency for the seed-42 models. Labels are grouped into *None*, *Large*, *Medium*, and *Small* buckets according to the number of training samples (support) for each discourse label. Label names follow the format *Task:Subclass* (e.g., *CogDem:Analysis_Give*, *Questions:Open*, *Feedback1:Affirming*), where the prefix denotes the discourse task and the suffix denotes the specific subclass. The proposed DG-MFP model achieves its clearest advantage on the *Small* bucket, indicating stronger benefits for low-support discourse subclasses.

Table 3. Robustness under missing and misaligned modalities. We report macro-F1, with the relative drop from the clean condition in parentheses.

Model	Clean	Misalignment	Missing	Missing w/o Prompt
Early-Fusion Longformer	0.6613	0.6392 (-0.0222)	0.6505 (-0.0109)	0.3522 (-0.3091)
DG-MFP (Ours)	0.6742	0.6736 (-0.0006)	0.6730 (-0.0012)	0.4304 (-0.2438)

455 for each discourse task. For example, the *CogDem* task ex- 455
 456 hibits average gate values of $g_{TA} = 0.44$, $g_{TV} = 0.50$, and 456
 457 $g_{INT} = 0.23$, while the *Questions* task shows $g_{TA} = 0.47$, 457
 458 $g_{TV} = 0.49$, and $g_{INT} = 0.24$. This consistent pattern in- 458
 459 dicates that bimodal text–audio and text–video signals pro- 459
 460 vide most of the complementary information beyond the 460
 461 text-only baseline, whereas higher-order three-way modal- 461
 462 ity interactions appear less dominant. 462

463 **Task-specific modality patterns.** Although the overall 463
 464 trend is consistent, different discourse tasks exhibit subtle 464
 465 differences in modality usage. For example, the *Feedback1* 465
 466 task shows slightly stronger text–audio influence ($g_{TA} =$ 466
 467 0.49) than text–video ($g_{TV} = 0.46$), suggesting that acous- 467
 468 tic cues may help distinguish different forms of teacher 468
 469 feedback. In contrast, the *CogDem* and *Feedback2* tasks 469
 470 exhibit slightly larger g_{TV} values (around 0.48–0.50), in- 470
 471 dicated that visual classroom cues can provide complemen- 471
 472 tary signals for identifying cognitive demand and feedback 472
 473 structures. 473

474 **Behavior on dominant *None* classes.** Many *None* cat- 474
 475 egories also exhibit relatively strong bimodal gates. This 475
 476 suggests that multimodal cues may help the model confirm 476

the absence of specific discourse functions. However, the 477
 interaction gate g_{INT} remains consistently smaller for these 478
 dominant classes, indicating that detecting the absence of 479
 discourse signals typically does not require complex three- 480
 way cross-modal interactions. 481

482 **Multimodal reliance for rare labels.** A closer inspec- 482
 483 tion of the label-level gate values (Figure 5) reveals that 483
 484 several low-frequency discourse labels rely more strongly 484
 485 on multimodal signals. For example, rare labels such as 485
 486 *Questions:Prompt* (test-set support = 44) and *Cog-* 486
 487 *Dem:Analysis_Give* (test-set support = 43) exhibit rela- 487
 488 tively large text–audio or text–video gate values (g_{TA} or 488
 489 g_{TV} around 0.43–0.48). Similarly, *CogDem:Report_Give* 489
 490 shows a particularly strong visual contribution with $g_{TV} \approx$ 490
 491 0.60. 491

492 In contrast, dominant *None* categories, which have thou- 492
 493 sands of training examples, often show more moderate gate 493
 494 magnitudes. This suggests that the model can rely primar- 494
 495 ily on transcript context for frequent labels, whereas multi- 495
 496 modal cues become more helpful when predicting less fre- 496
 497 quent discourse categories. 497

498 This observation is consistent with the bucket-level re- 498
 499 sults in Table 2, where the proposed model achieves the 499

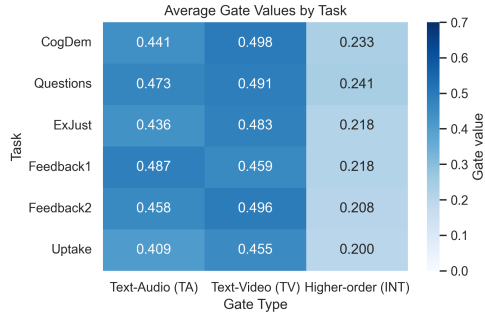


Figure 4. Average delta gate values across discourse tasks. Columns correspond (from left to right) to the text–audio (g_{TA}), text–video (g_{TV}), and higher-order interaction (g_{INT}) gates. The bimodal gates (g_{TA} and g_{TV}) consistently exhibit larger values than the interaction gate (g_{INT}), suggesting that bimodal modality signals contribute most of the complementary information beyond the text-only baseline, while three-way modality interactions are comparatively less dominant.

500 largest performance improvement on the small-label bucket.
 501 Together, these findings suggest that delta-gated fusion
 502 helps the model selectively leverage multimodal signals
 503 when transcript evidence alone is insufficient.

504 5.7. MoE Comparison as a Qualitative Reference

505 To further examine interpretability, we compare DG-MFP
 506 with a pooled-level Mixture-of-Experts (MoE) variant used
 507 as a qualitative reference. In this model, the pooled [CLS]
 508 representation is routed to multiple MLP experts through a
 509 soft router, and expert outputs are combined before classifica-
 510 tion. Unlike DG-MFP, this MoE variant does not explic-
 511 itly separate text–audio, text–video, and higher-order inter-
 512 action effects.

513 During training, routing statistics reveal limited expert
 514 specialization. The router often assigns most selections to
 515 a single expert in early epochs and gradually converges to
 516 nearly uniform expert usage ($\approx 25\%$ per expert, entropy
 517 $\approx \log 4$), suggesting weak semantic specialization.

518 In contrast, DG-MFP produces gates directly tied to
 519 structured modality contributions. These gates are task-
 520 specific and per-class, making them more interpretable than
 521 latent expert assignments in MoE.

522 6. Conclusion

523 In this work, we studied robust multimodal classroom dis-
 524 course classification under realistic conditions where audio
 525 and video signals may be missing, noisy, or imperfectly
 526 aligned with the transcript. To address this problem, we pro-
 527 posed a **Delta-Gated Multi-Forward-Pass Longformer**
 528 that treats the transcript as the primary modality and models
 529 audio and video as incremental sources of information.

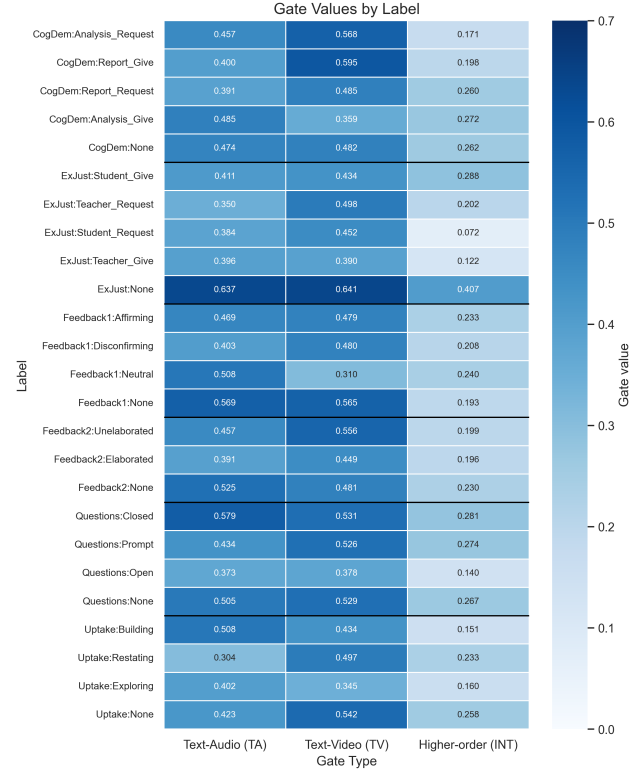


Figure 5. Delta gate values for each discourse label. Rows correspond to task-specific labels, and columns correspond (from left to right) to the text–audio (g_{TA}), text–video (g_{TV}), and higher-order interaction (g_{INT}) modality pathways. The visualization shows that multimodal contributions vary across labels, with several categories exhibiting relatively stronger bimodal modality gates.

530 The proposed framework uses a shared encoder in mul-
 531 tiple masked forward passes and decomposes the prediction
 532 into a baseline transcript, modality-specific increments, and
 533 a higher-order interaction term. This formulation enables
 534 explicit modeling of modality contributions while maintain-
 535 ing a consistent representation space across forward passes.

536 Experiments on classroom discourse data show that
 537 the proposed method improves performance over standard
 538 early-fusion baselines, achieves strong gains on rare dis-
 539 course labels, and demonstrates substantially stronger ro-
 540 bustness under missing-modality and misalignment condi-
 541 tions. Analysis of the learned gates further reveals inter-
 542 pretable task-specific patterns of modality use, indicating
 543 that multimodal signals are most beneficial when transcript
 544 evidence alone is insufficient for reliable prediction.

545 In general, our results suggest that multimodal fusion
 546 can be effectively modeled as incremental contributions
 547 to a transcript-based prediction. We hope that this per-
 548 spective will encourage further work on robust and inter-
 549 pretable multimodal learning in real-world educational en-
 550 vironments.

551

References

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

[1] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katie Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford, Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob Menick, Sebastian Borgeaud, Andrew Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikolaj Binkowski, Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, and Karen Simonyan. Flamingo: a visual language model for few-shot learning, 2022. 2

[2] Tadas Baltrušaitis, Chaitanya Ahuja, and Louis-Philippe Morency. Multimodal machine learning: A survey and taxonomy. *IEEE TPAMI*, 41(2):423–443, 2019. 3

[3] Iz Beltagy, Matthew E. Peters, and Arman Cohan. Longformer: The long-document transformer, 2020. arXiv:2004.05150. 2

[4] Paulo Blikstein. Multimodal learning analytics. In *Proceedings of the Third International Conference on Learning Analytics and Knowledge*, pages 102–106, 2013. 3

[5] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding, 2019. arXiv:1810.04805. 1, 2

[6] DrivenData and University of Virginia. Artificial intelligence for advancing instruction (aiai) challenge. <https://github.com/drivendataorg/ai-ai-challenge>, 2025. Accessed: 2026. 3

[7] Mengmeng Ma, Jian Ren, Long Zhao, Davide Testuggine, and Xi Peng. Are multimodal transformers robust to missing modality? In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 18177–18186, 2022. 2

[8] Natalia Neverova, Christian Wolf, Graham Taylor, and Florian Nebout. Moddrop: Adaptive multi-modal gesture recognition. *IEEE TPAMI*, 38(8):1692–1706, 2016. 2

[9] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision, 2021. arXiv:2103.00020. 2

[10] Chen Sun, Austin Myers, Carl Vondrick, Kevin Murphy, and Cordelia Schmid. Videobert: A joint model for video and language representation learning. In *ICCV*, 2019. 2

[11] Yao-Hung Hubert Tsai, Shaojie Bai, Paul Pu Liang, J. Zico Kolter, Louis-Philippe Morency, and Ruslan Salakhutdinov. Multimodal transformer for unaligned multimodal language sequences, 2019. arXiv:1906.00295. 2, 3

[12] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need, 2017. arXiv:1706.03762. 1, 2

[13] Yi Wang, Kunchang Li, Yizhuo Li, Yanan He, Bingkun Huang, Zhiyu Zhao, Hongjie Zhang, Jilan Xu, Yi Liu, Zun Wang, Sen Xing, Guo Chen, Junting Pan, Jiashuo Yu, Yali Wang, Limin Wang, and Yu Qiao. Internvideo: General

video foundation models via generative and discriminative learning, 2022. 2 608
609

[14] Peng Xu, Xiatian Zhu, and David A. Clifton. Multimodal learning with transformers: A survey. *IEEE TPAMI*, 45(10): 12113–12132, 2023. 1 610
611
612