

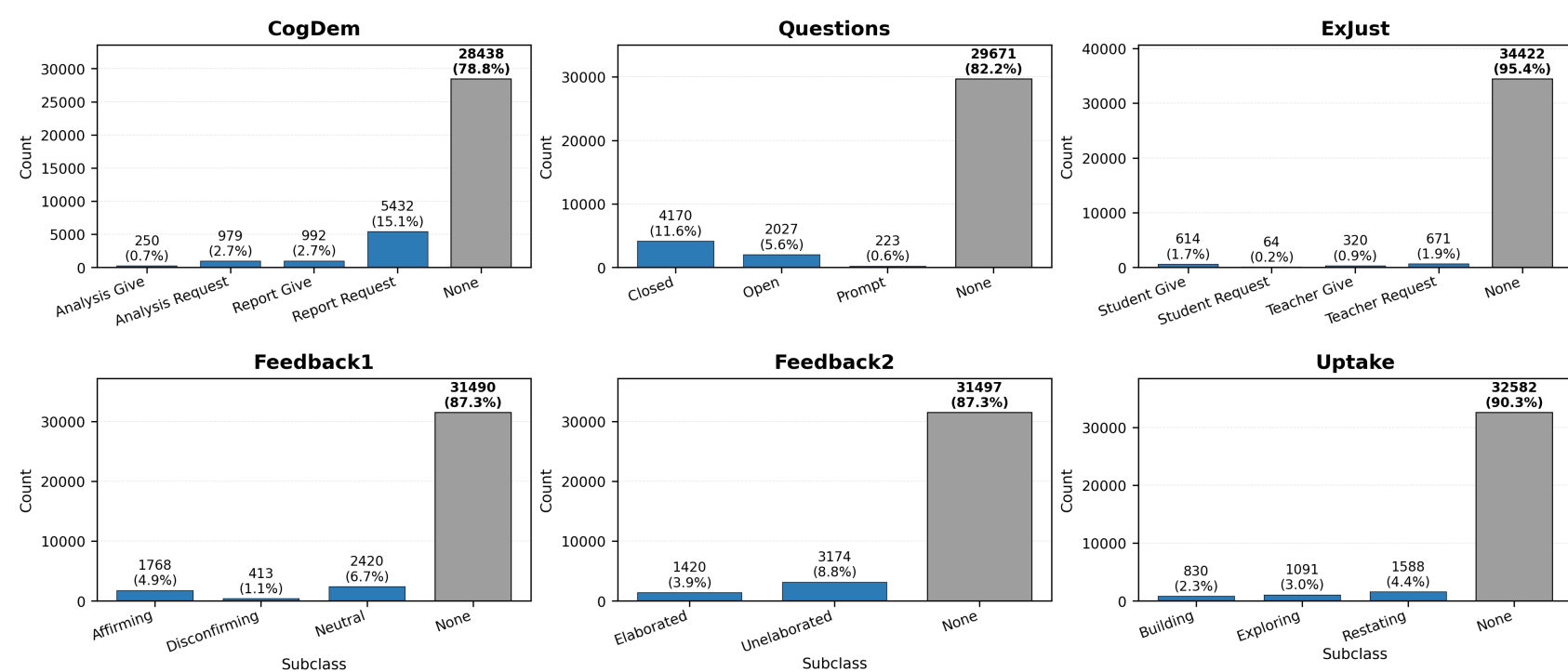
## Abstract

### In noisy classrooms, do audio/video help or mislead the transcript model?

- Classroom recordings are noisy, missing, and weakly aligned.
- We estimate contribution and reliability, not just fusion.
- DG-MFP isolates audio/video gains and gates trusted evidence.

<b>Data</b>	AIAIS; 107 videos; 36,091 samples
<b>Tasks</b>	6 tasks, 25 labels
<b>Signal</b>	Text + HuBERT + CLIP
<b>Metric</b>	F1 + macro-F1 robustness

Label Distribution Across Six Discourse Tasks (N = 36091)

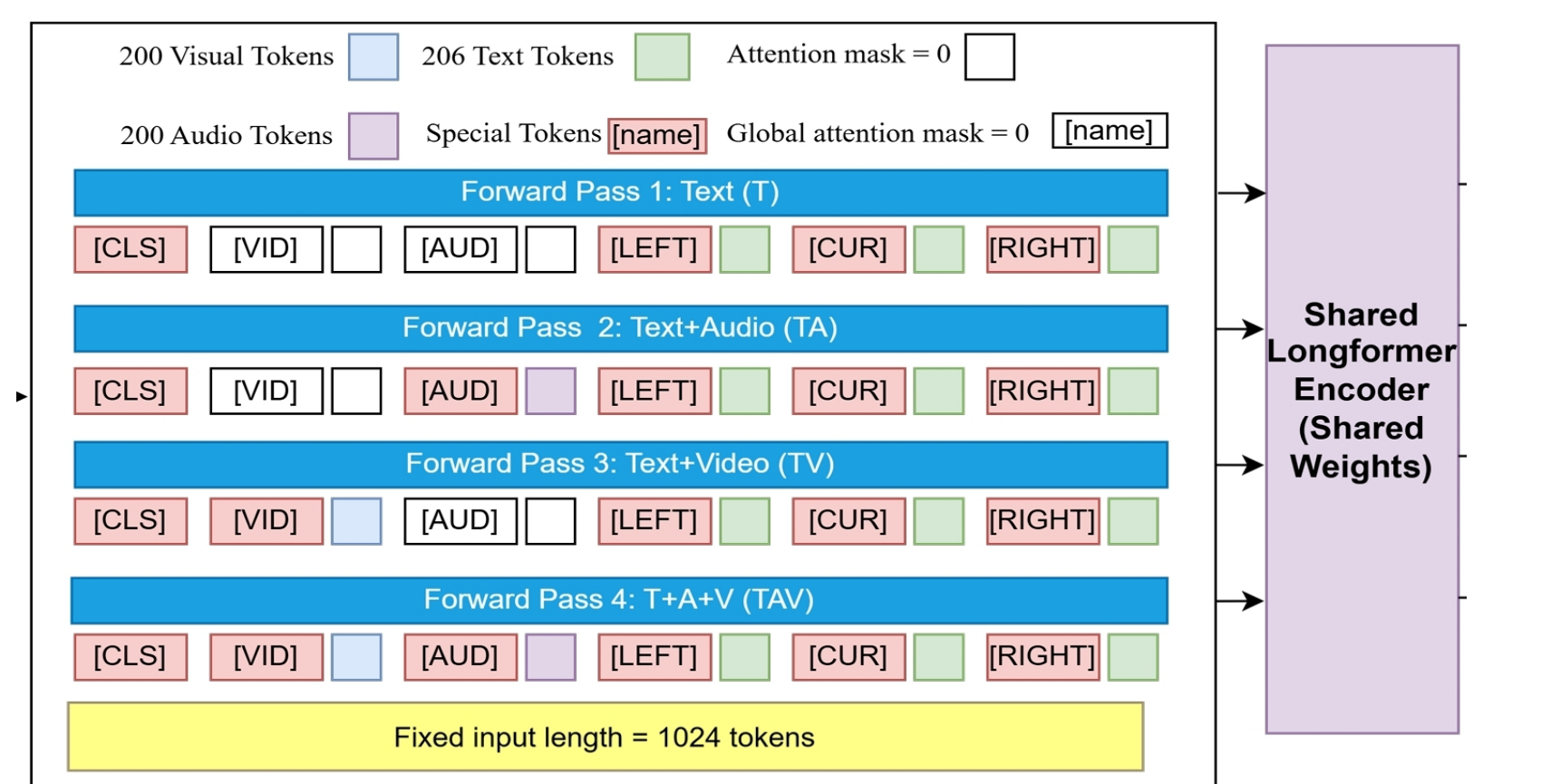


Long-tailed labels make rare-class robustness the real test.

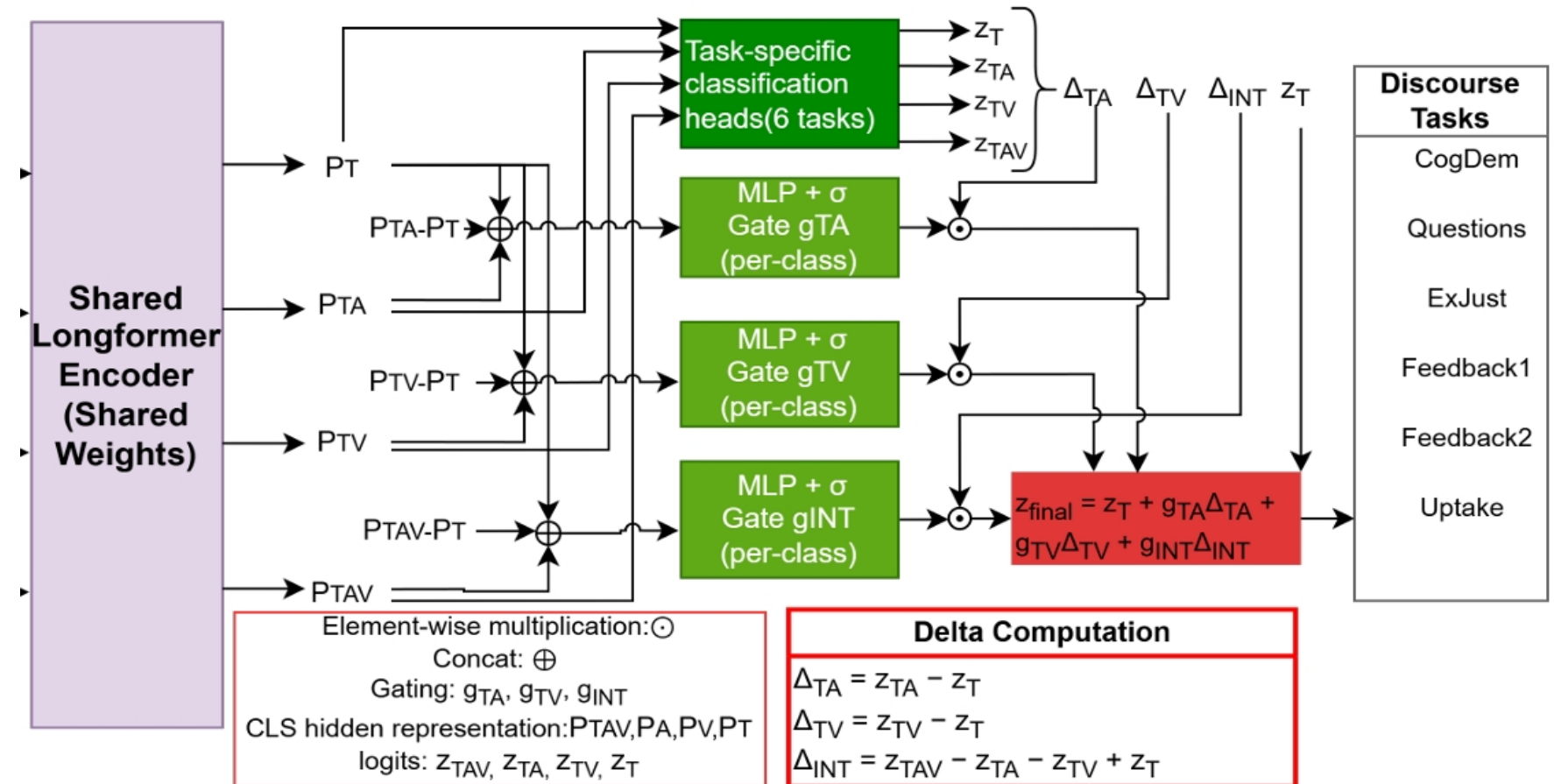
### Core idea

T = transcript text; A = audio; V = video.  
Use T as the baseline prediction.  
Compare TA, TV, and TAV passes to isolate added signal.  
Logit deltas estimate contribution; gates estimate reliability.

## Methods and Materials



### 1. Masked multi-pass inputs.



### 2. Delta-gated decomposition.

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

Notation: T=text, A=audio, V=video; TA/TV bimodal, TAV all modalities.

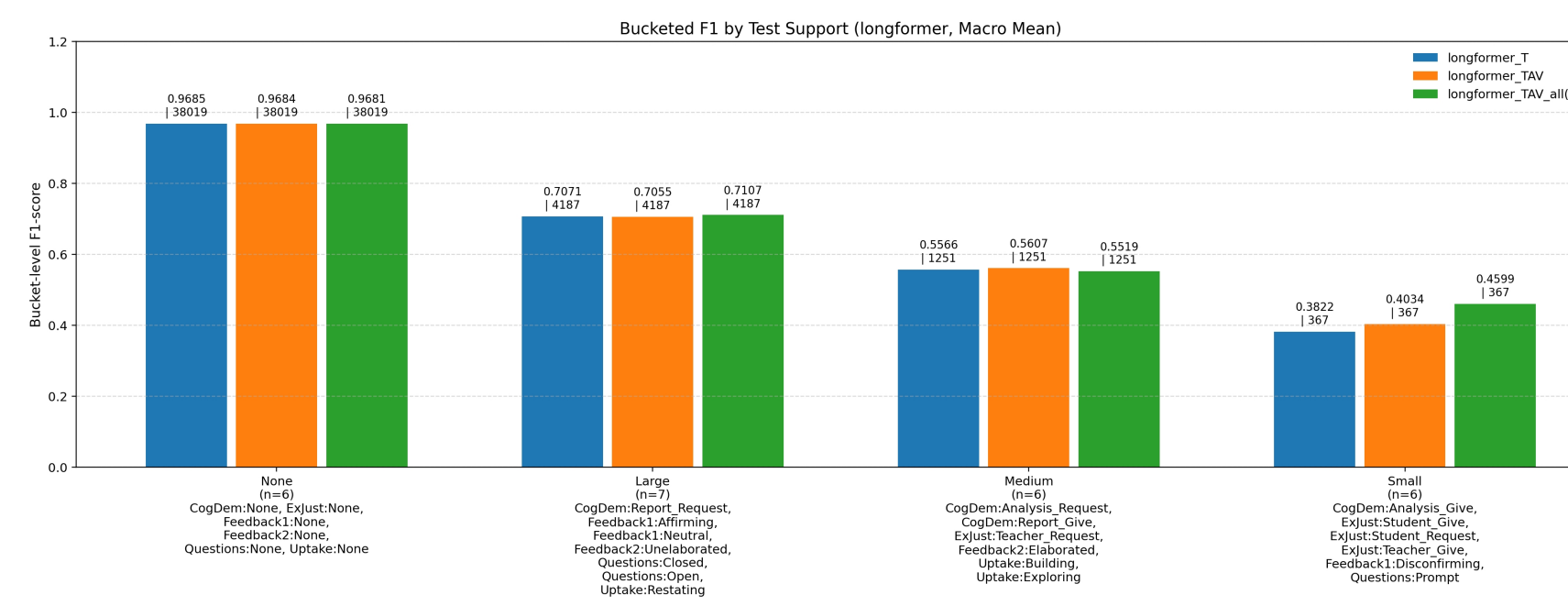
Question: do audio/video help, how much, or mislead?  
Deltas measure gain over the transcript baseline.  
Gates estimate reliability of each modality gain.  
Stress tests verify trust under missing/misaligned cues.

## Results

**0.6777** best F1  
**+0.0565** small-label gain  
**-0.0012** missing drop

### Overall F1 on full test set

Model	S42	S43	S44
Text	0.6588	0.6599	0.6636
Early	0.6711	0.6676	0.6678
DG-MFP	<b>0.6777</b>	<b>0.6712</b>	<b>0.6687</b>



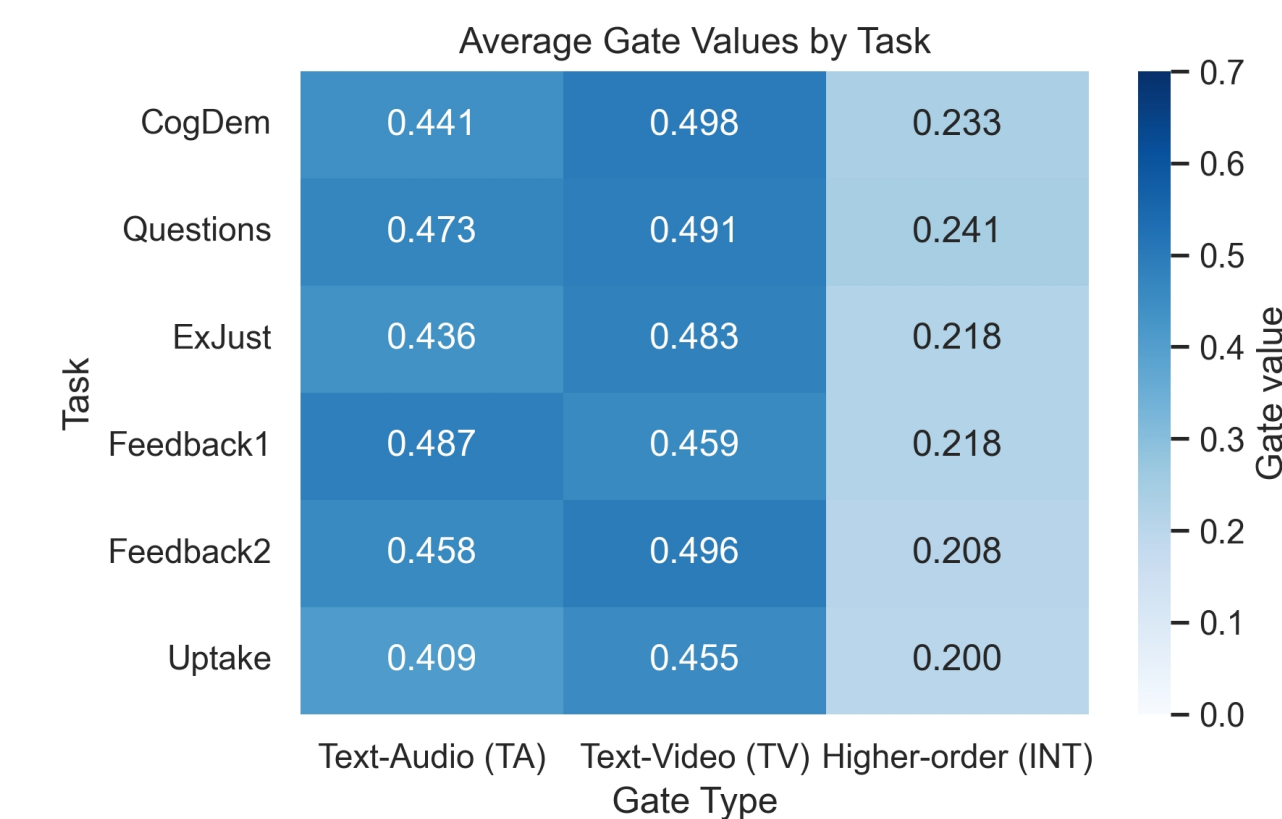
### Bucketed F1 by label support (seed 42)

Bucket	#	Text	Early	Ours
None	6	0.9685	0.9684	0.9681
Large	7	0.7071	0.7055	<b>0.7107</b>
Medium	6	0.5566	<b>0.5607</b>	0.5519
Small	6	0.3822	0.4034	<b>0.4599</b>

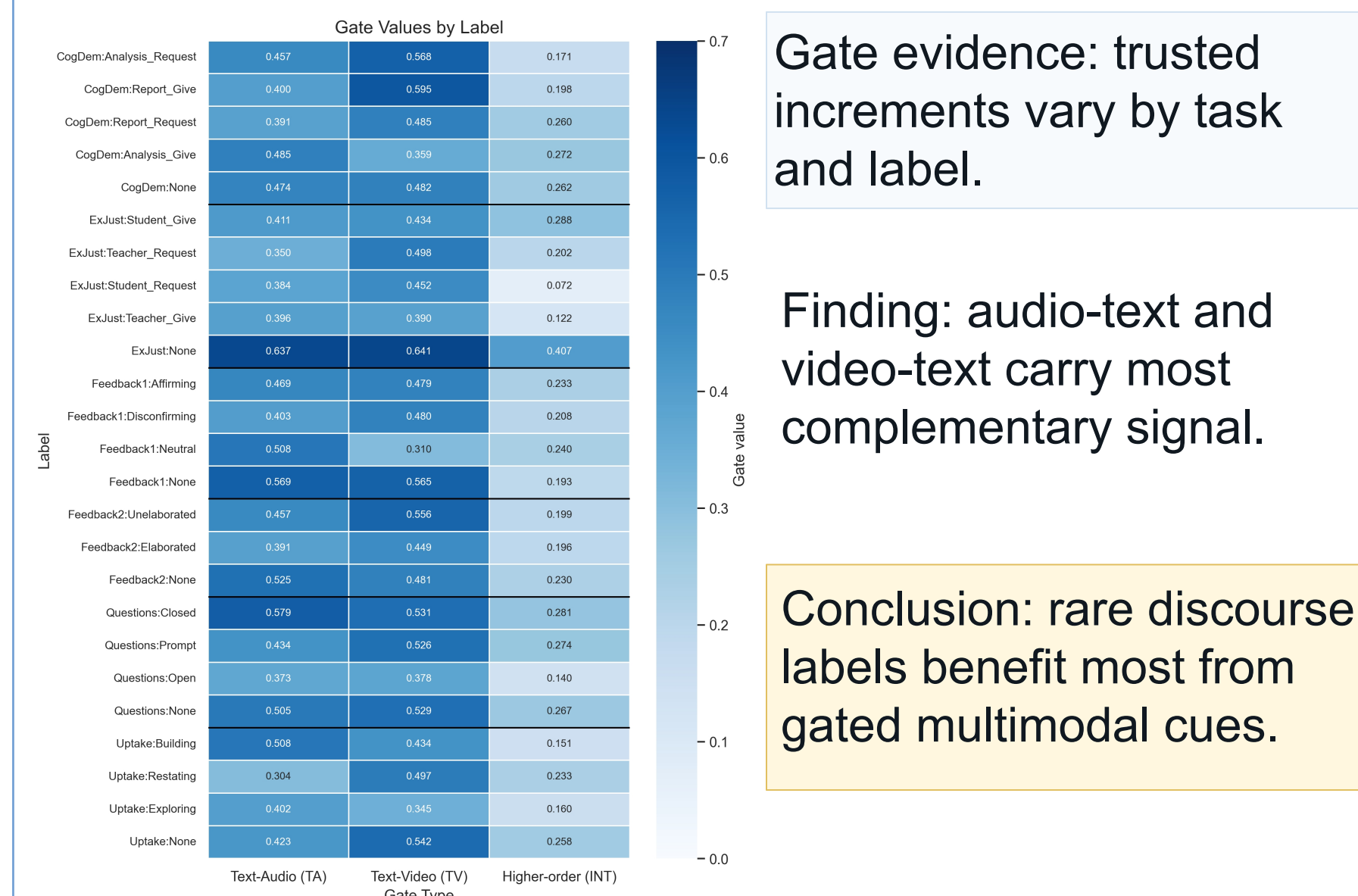
### Stress robustness: macro-F1 (drop from clean)

Model	Clean	Misalign	Missing	NoPrompt
Early	0.6613	0.6392 (-0.0222)	0.6505 (-0.0109)	0.3522 (-0.3091)
DG-MFP	<b>0.6742</b>	<b>0.6736 (-0.0006)</b>	<b>0.6730 (-0.0012)</b>	<b>0.4304 (-0.2438)</b>

## Discussion and Conclusions



Discussion: gate patterns explain model behavior.



Gate evidence: trusted increments vary by task and label.

Finding: audio-text and video-text carry most complementary signal.

Conclusion: rare discourse labels benefit most from gated multimodal cues.

DG-MFP turns fusion into measurable contribution and reliability.

Contact: Chongyu He | ekn8kz@virginia.edu | CV4Edu-17  
Data reference: AIAIS dataset, AIAI Challenge 2025