



Festival of Learning • AIED 2026, Seoul, Korea

## When Can We Trust AI Coding of Student-Generated Text?

A Committee-Based Approach to Diagnosing Agreement & Uncertainty at Scale

Fanjie Li, Madison Lee Mason, Daniel T. Levin, Alyssa Friend Wise



VANDERBILT  
UNIVERSITY



VANDERBILT  
Learning Innovation Incubator

# A New Challenge: Monitoring AI Assessment of Student-Generated Text

- **Deployment challenge:** As AI assessment of student work (e.g., essays, open-ended responses) moves into deployed educational systems, a critical question emerges:

*How do we know it continues to interpret student work as intended?*

- **Limits of standard test-set evaluation:** Initial offline evaluation establishes baseline performance, but cannot guarantee stable behavior as models encounter new learners, contexts, and forms of expression.
- **Need for ongoing monitoring:** This calls for interpretable diagnostics that support ongoing monitoring of AI coding on unlabeled data, once deployed.

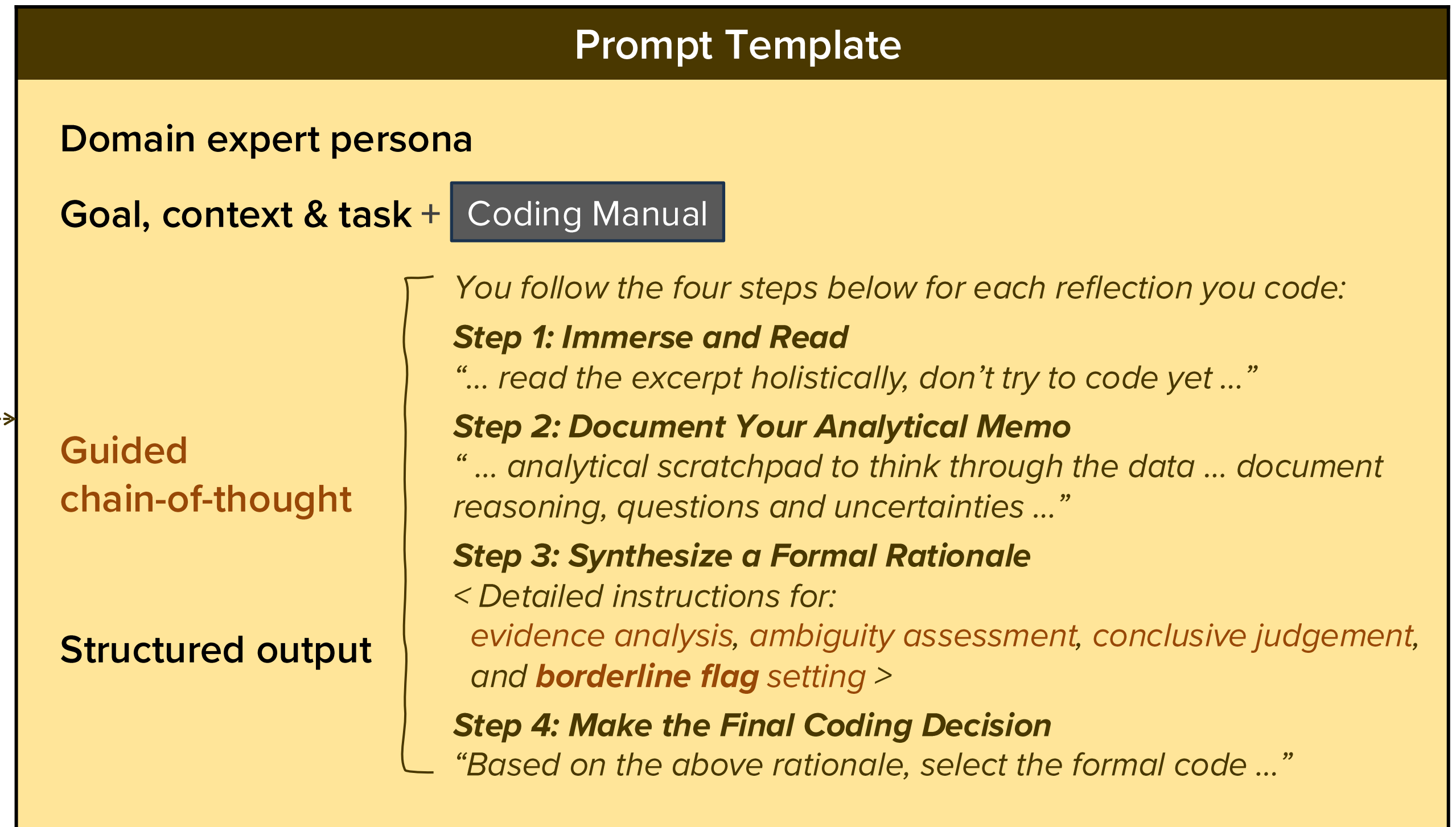
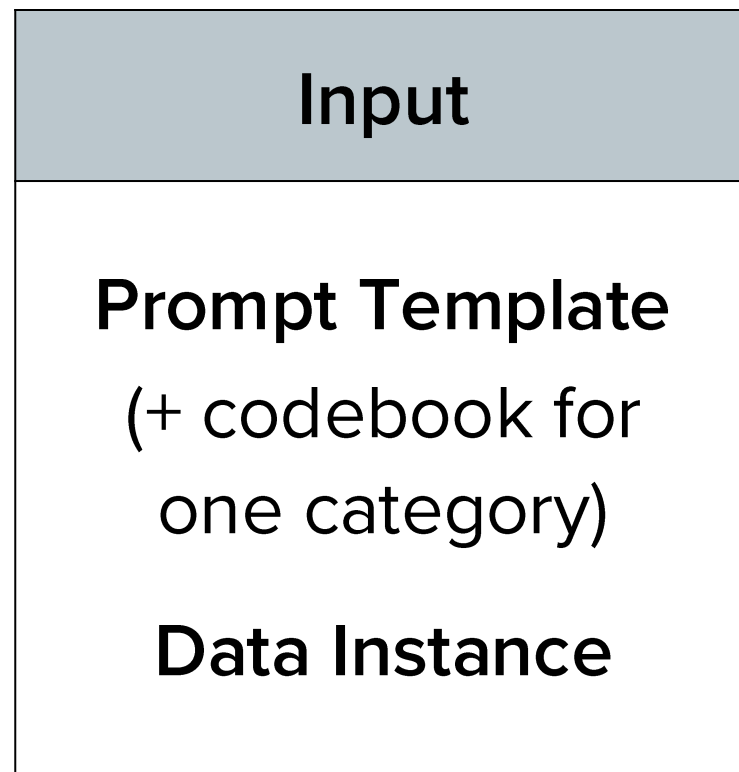
# A Committee-Based Diagnostic Framework to Support Interpretable Monitoring of Automated Coding

Uncertainty  
Quantification

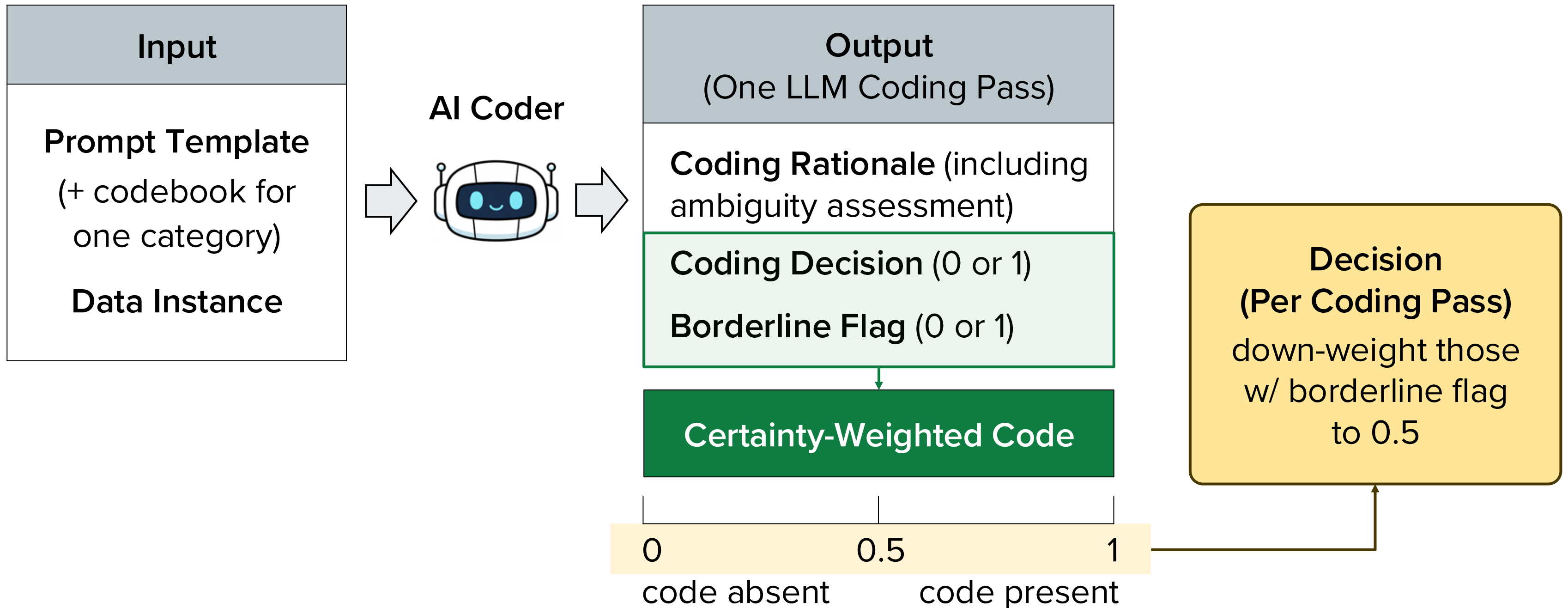
support human decisions about when to **trust**, **re-evaluate**, or **reject** automated categorization

Including this AI Coding System Element:	Allows us to Assess the <b>Uncertainty</b> :
<u>Guided Reasoning</u> with Uncertainty Flags	In Individual LLM Judgement
<u>Repeated Inference</u> for Decision Stability	Across Multiple Runs by an LLM
<u>Committee-Based</u> Coding with Multiple LLMs	Across Multiple LLMs

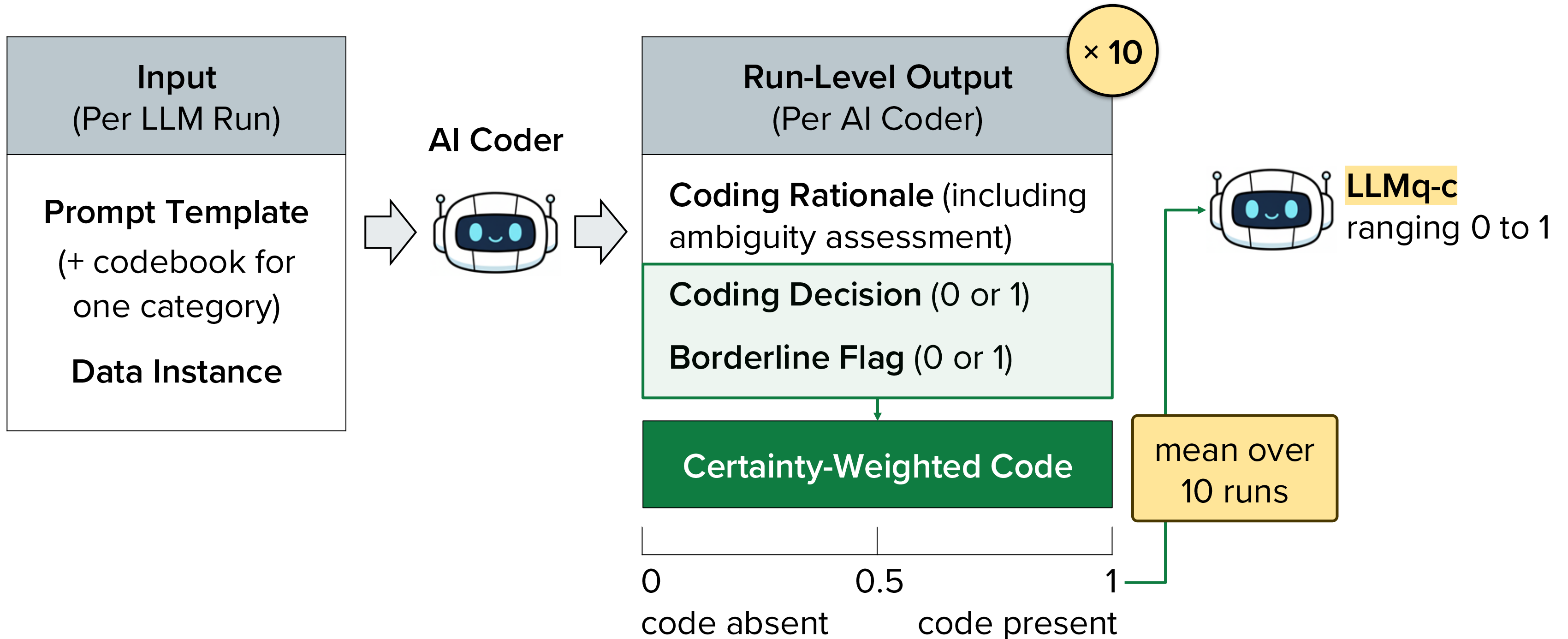
# Guided Reasoning with Uncertainty Flags



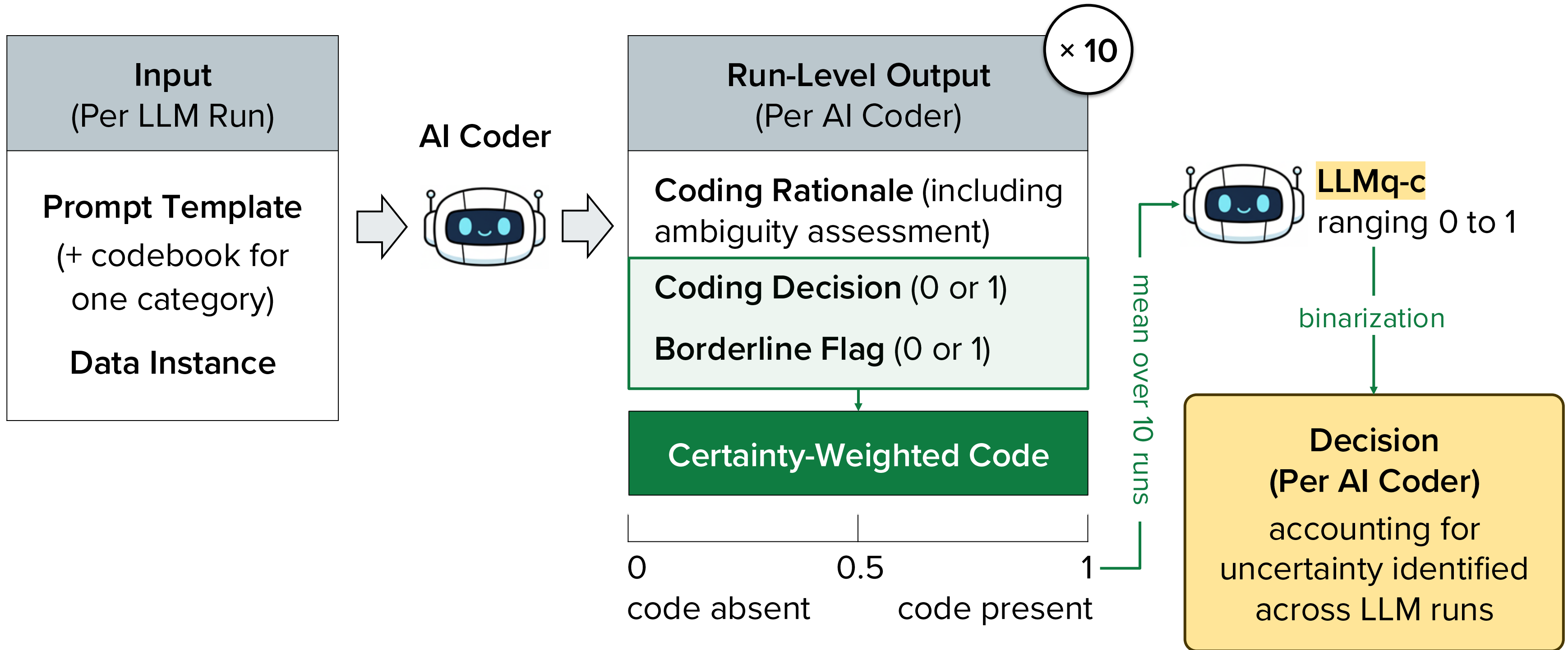
# Guided Reasoning with Uncertainty Flags



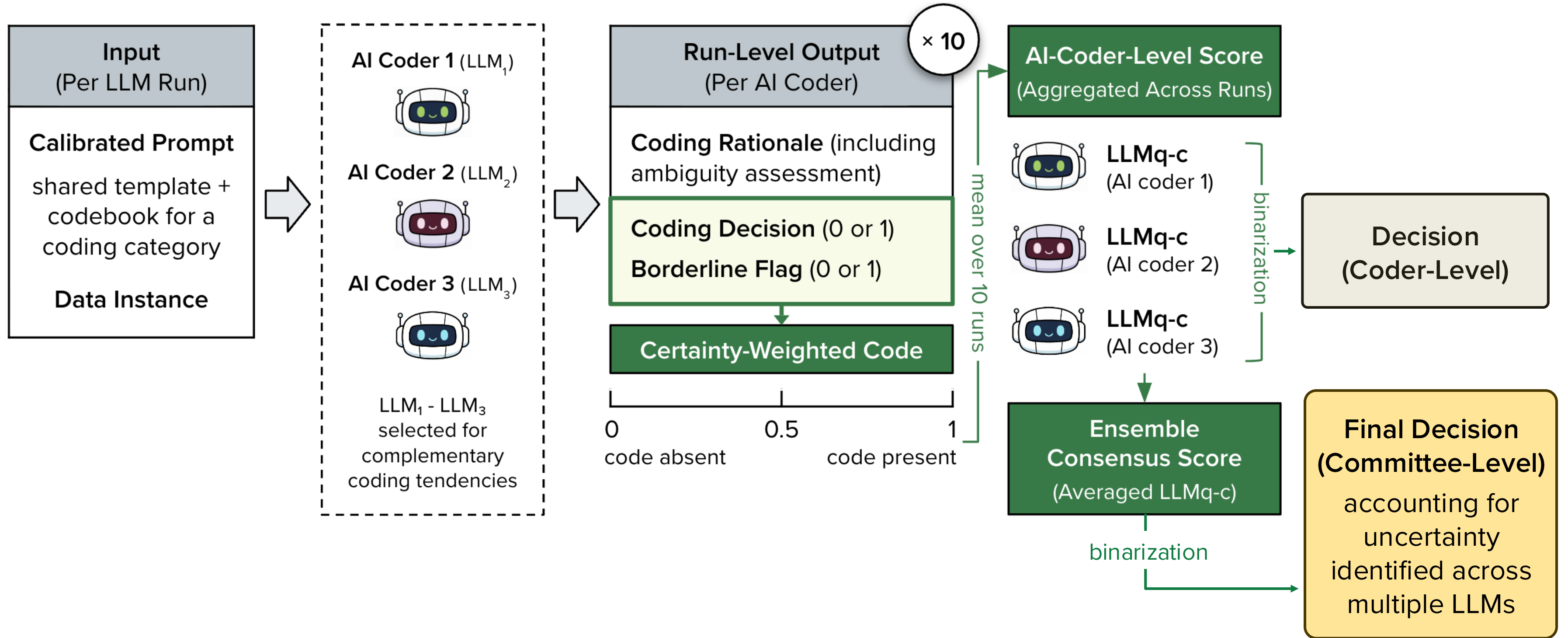
# Repeated Inference for Decision Stability



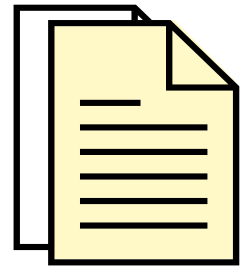
# Repeated Inference for Decision Stability





# Committee-Based Coding with Multiple LLMs



# Monitoring Metrics (Available on Unlabeled Data)



 (LLM <sub>1</sub> )	1 <sub>flag</sub>	1 <sub>flag</sub>	1 <sub>flag</sub>	1 <sub>flag</sub>	1 <sub>flag</sub>	1 <sub>flag</sub>	1	1 <sub>flag</sub>	1	1	→ LLMq-c = 0.79	} Ensemble Consensus Score <b>0.39</b>
 (LLM <sub>2</sub> )	0	0	0 <sub>flag</sub>	0	0	0 <sub>flag</sub>	0	0	0	0 <sub>flag</sub>	→ LLMq-c = 0.09	
 (LLM <sub>3</sub> )	0 <sub>flag</sub>	0 <sub>flag</sub>	0 <sub>flag</sub>	0 <sub>flag</sub>	0 <sub>flag</sub>	0 <sub>flag</sub>	0 <sub>flag</sub>	0 <sub>flag</sub>	0 <sub>flag</sub>	0 <sub>flag</sub>	→ LLMq-c = 0.3	

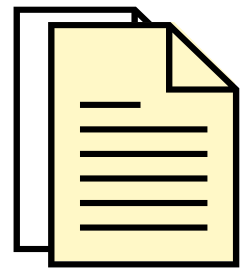
## Set-Level Monitoring

Inter-Model Agreement ( $\alpha_{AI}$ )

Consensus Entropy ( $H$ )

Borderline Rate ( $B$ )

# Monitoring Metrics (Available on Unlabeled Data)



 (LLM <sub>1</sub> )	1 <sub>flag</sub>	1 <sub>flag</sub>	1 <sub>flag</sub>	1 <sub>flag</sub>	1 <sub>flag</sub>	1 <sub>flag</sub>	1	1 <sub>flag</sub>	1	1	→ LLMq-c = 0.79
 (LLM <sub>2</sub> )	0	0	0 <sub>flag</sub>	0	0	0 <sub>flag</sub>	0	0	0	0 <sub>flag</sub>	→ LLMq-c = 0.09
 (LLM <sub>3</sub> )	0 <sub>flag</sub>	0 <sub>flag</sub>	0 <sub>flag</sub>	0 <sub>flag</sub>	0 <sub>flag</sub>	0 <sub>flag</sub>	0 <sub>flag</sub>	0 <sub>flag</sub>	0 <sub>flag</sub>	0 <sub>flag</sub>	→ LLMq-c = 0.3

} Ensemble Consensus Score **0.39**

## Set-Level Monitoring

Inter-Model Agreement ( $\alpha_{AI}$ )

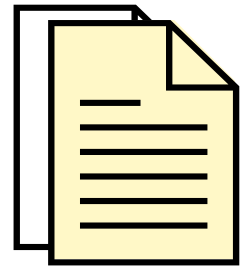
$\alpha_{AI} > 0.67$  ----- Moderate agreement

$\alpha_{AI} > 0.80$  ----- Strong agreement

Consensus Entropy ( $H$ )

Borderline Rate ( $B$ )

# Monitoring Metrics (Available on Unlabeled Data)



 ( LLM <sub>1</sub> )	1 <sub>flag</sub>	1 <sub>flag</sub>	1 <sub>flag</sub>	1 <sub>flag</sub>	1 <sub>flag</sub>	1 <sub>flag</sub>	1	1 <sub>flag</sub>	1	1	→ LLMq-c = 0.79
 ( LLM <sub>2</sub> )	0	0	0 <sub>flag</sub>	0	0	0 <sub>flag</sub>	0	0	0	0 <sub>flag</sub>	→ LLMq-c = 0.09
 ( LLM <sub>3</sub> )	0 <sub>flag</sub>	0 <sub>flag</sub>	0 <sub>flag</sub>	0 <sub>flag</sub>	0 <sub>flag</sub>	0 <sub>flag</sub>	0 <sub>flag</sub>	0 <sub>flag</sub>	0 <sub>flag</sub>	0 <sub>flag</sub>	→ LLMq-c = 0.3

$\hat{p}$

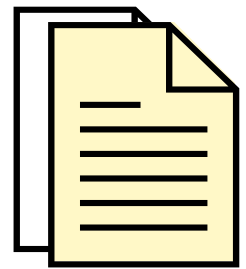
Ensemble  
Consensus Score  
**0.39**

## Set-Level Monitoring

- Inter-Model Agreement (  $\alpha_{AI}$  )
- Consensus Entropy (  $H$  )
- Borderline Rate (  $B$  )

**Lower  $H$**  --- more **decisive** judgments ( $\hat{p}$  near 0 or 1)  
**Higher  $H$**  --- greater **indecisiveness** ( $\hat{p}$  near 0.5)

# Monitoring Metrics (Available on Unlabeled Data)



 ( LLM <sub>1</sub> )	1 <sub>flag</sub>	1 <sub>flag</sub>	1 <sub>flag</sub>	1 <sub>flag</sub>	1 <sub>flag</sub>	1 <sub>flag</sub>	1	1 <sub>flag</sub>	1	1	→ LLMq-c = 0.79
 ( LLM <sub>2</sub> )	0	0	0 <sub>flag</sub>	0	0	0 <sub>flag</sub>	0	0	0	0 <sub>flag</sub>	→ LLMq-c = 0.09
 ( LLM <sub>3</sub> )	0 <sub>flag</sub>	0 <sub>flag</sub>	0 <sub>flag</sub>	0 <sub>flag</sub>	0 <sub>flag</sub>	0 <sub>flag</sub>	0 <sub>flag</sub>	0 <sub>flag</sub>	0 <sub>flag</sub>	0 <sub>flag</sub>	→ LLMq-c = 0.3

$\hat{p}$

Ensemble  
Consensus Score  
**0.39**

## Set-Level Monitoring

Inter-Model Agreement (  $\alpha_{AI}$  )

Consensus Entropy (  $H$  )

Borderline Rate (  $B$  )

**Lower  $H$**  --- more **decisive** judgments ( $\hat{p}$  near 0 or 1)

**Higher  $H$**  --- greater **indecisiveness** ( $\hat{p}$  near 0.5)

Percentage of items for which  $L \leq \hat{p} \leq U$

Distinguish between **pervasive uncertainty** vs **uncertainty concentrated on a few items**

## Proof-of-Concept Pilot | Context and Data

- **Nursing students' post-simulation written reflection** (213 event-level reflections, 16 students); Human coders double-coded 60 reflections (IRR:  $0.63 < \alpha < 0.91$ )
- **Coding scheme:**  
9 simulation-general nursing activity codes for student reflections
- **AI committee set up:**  
**3 LLMs × 10 runs each**  
(o3, claude-sonnet-4, deepseek-r1)



What are they reflecting about? (9 categories)

BEHAVIORAL MOVES	MENTAL MOVES
<ul style="list-style-type: none"><li>• <u>Gathering information</u></li><li>• <u>Implementation of nursing care</u></li><li>• <u>Communication</u><ul style="list-style-type: none"><li>◦ <u>w/ healthcare team members</u></li><li>◦ <u>w/ patient or patient family</u></li></ul></li></ul>	<ul style="list-style-type: none"><li>• <u>Interpreting information</u><ul style="list-style-type: none"><li>◦ <u>to understand the patient situation</u></li><li>◦ <u>to validate care delivery plan</u></li><li>◦ <u>to evaluate outcomes</u></li></ul></li><li>• <u>Establish goals &amp; generate solutions</u></li><li>• <u>Prioritization &amp; time management</u></li></ul>

💡 To code for **Dimension 1** (what are they reflecting about?), focus on the specific behavioral moves (clinical actions) and/or mental moves (clinical reasoning) that the student describes, rather than the general phase of nursing process. **Ask: "Which type(s) of behavioral and/or mental moves is the student reflecting on?"**

# Proof-of-Concept Pilot | Results

Activity Code	Test Data				Unlabeled Data		
	F1	$\alpha_{AI}$	$H$	$B$	$\alpha_{AI}$	$H$	$B$
GATHER	0.966	0.942	0.295	4.17%	0.885	0.275	7.69%
IMPLMT	0.941	1.000	0.322	0.00%	0.951	0.218	5.13%
T-COMM	0.923	0.877	0.295	8.33%	0.946	0.230	2.56%
P-COMM	0.923	0.944	0.214	8.33%	0.917	0.131	3.42%
INTRPT-SIT	0.770	0.683	0.344	8.33%	0.756	0.276	9.40%
INTRPT-VLDT	0.947	1.000	0.321	0.00%	0.837	0.343	9.40%
INTRPT-EVAL	1.000	0.844	0.195	0.00%	0.532	0.163	1.71%
GOALS	0.824	0.698	0.506	25.00%	0.719	0.456	12.82%
TIME	1.000	0.735	0.379	4.17%	0.644	0.309	11.11%

**Behaviorally**  
explicit codes

Need inference  
about **cognitive**  
activity

# Proof-of-Concept Pilot | Results

Activity Code	Test Data				Unlabeled Data		
	F1	$\alpha_{AI}$	$H$	$B$	$\alpha_{AI}$	$H$	$B$
GATHER	0.966	0.942	0.295	4.17%	0.885	0.275	7.69%
IMPLMT	0.941	1.000	0.322	0.00%	0.951	0.218	5.13%
T-COMM	0.923	0.877	0.295	8.33%	0.946	0.230	2.56%
P-COMM	0.923	0.944	0.214	8.33%	0.917	0.131	3.42%
INTRPT-SIT	0.770	0.683	0.344	8.33%	0.756	0.276	9.40%
INTRPT-VLDT	0.947	1.000	0.321	0.00%	0.837	0.343	9.40%
INTRPT-EVAL	1.000	0.844	0.195	0.00%	0.532	0.163	1.71%
GOALS	0.824	0.698	0.506	25.00%	0.719	0.456	12.82%
TIME	1.000	0.735	0.379	4.17%	0.644	0.309	11.11%

**Behaviorally**  
explicit codes

**GATHER** ( gather info )  
**IMPLMT** ( implementation of nursing care )  
**T-COMM** ( communication w/ healthcare team )  
**P-COMM** ( communication w/ patient family )

## Stable consensus profile

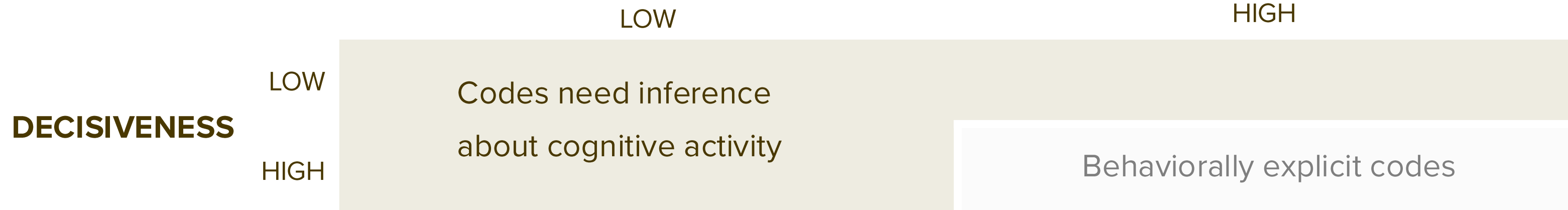
- High test-set performance (  $F1 > .9$  )
- Stable  $\alpha_{AI}$  across test and unlabeled data ( all  $\alpha_{AI} > .8$  )
- Consensus entropy (  $H$  ) remain relatively stable (low uncertainty)
- Borderline rates stayed below 10%

# Proof-of-Concept Pilot | Results

Activity Code	Test Data				Unlabeled Data		
	F1	$\alpha_{AI}$	$H$	$B$	$\alpha_{AI}$	$H$	$B$
GATHER	0.966	0.942	0.295	4.17%	0.885	0.275	7.69%
IMPLMT	0.941	1.000	0.322	0.00%	0.951	0.218	5.13%
T-COMM	0.923	0.877	0.295	8.33%	0.946	0.230	2.56%
P-COMM	0.923	0.944	0.214	8.33%	0.917	0.131	3.42%
INTRPT-SIT	0.770	0.683	0.344	8.33%	0.756	0.276	9.40%
INTRPT-VLDT	0.947	1.000	0.321	0.00%	0.837	0.343	9.40%
INTRPT-EVAL	1.000	0.844	0.195	0.00%	0.532	0.163	1.71%
GOALS	0.824	0.698	0.506	25.00%	0.719	0.456	12.82%
TIME	1.000	0.735	0.379	4.17%	0.644	0.309	11.11%

Need inference about **cognitive** activity

## INTER-MODEL AGREEMENT



# Proof-of-Concept Pilot | Results

Activity Code	Test Data				Unlabeled Data		
	F1	$\alpha_{AI}$	$H$	$B$	$\alpha_{AI}$	$H$	$B$
INTRPT-SIT	0.770	0.683	0.344	8.33%	0.756	0.276	9.40%



Example Reflection | Coding Category: **INTRPT-SIT** (Interpret Info to Understand Patient Situation)

**What aspects of the simulation or your own actions led to you achieving your goals during this event?**

The mother informing me that the patient has dysphagia influenced my nursing actions before administering med.

 (o3)	1	1	1	1	1	1	1	1	1	1	1	1	→ LLMq-c <sub>(o3)</sub> = <b>0.95</b>	} Ensemble Consensus Score <b>0.58</b>
 (claude)	1	0	1	0	0	0	1	1	0	1	1	0	→ LLMq-c <sub>(claude)</sub> = <b>0.6</b>	
 (deepseek)	0	0	0	0	0	0	0	0	0	0	0	0	→ LLMq-c <sub>(deepseek)</sub> = <b>0.2</b>	

### Sample Rationale (Code = 0):

"While the reflection could be interpreted as implicitly containing interpretation work (understanding that dysphagia requires modified medication administration), they do not explicitly describe interpreting this information to understand the patient's condition."

### Sample Rationale (Code = 1):

"Some uncertainty: the reflection is brief and most of the emphasis is on the subsequent action (elevating the patient). However, the student does explicitly note the significance of the dysphagia information for understanding the patient's current condition, providing enough evidence for the code."

## Proof-of-Concept Pilot | Results

Activity Code	Test Data				Unlabeled Data		
	F1	$\alpha_{AI}$	$H$	$B$	$\alpha_{AI}$	$H$	$B$
GATHER	0.966	0.942	0.295	4.17%	0.885	0.275	7.69%
IMPLMT	0.941	1.000	0.322	0.00%	0.951	0.218	5.13%
T-COMM	0.923	0.877	0.295	8.33%	0.946	0.230	2.56%
P-COMM	0.923	0.944	0.214	8.33%	0.917	0.131	3.42%
INTRPT-SIT	0.770	0.683	0.344	8.33%	0.756	0.276	9.40%
INTRPT-VLDT	0.947	1.000	0.321	0.00%	0.837	0.343	9.40%
INTRPT-EVAL	1.000	0.844	0.195	0.00%	0.532	0.163	1.71%
GOALS	0.824	0.698	0.506	25.00%	0.719	0.456	12.82%
TIME	1.000	0.735	0.379	4.17%	0.644	0.309	11.11%

Illustrate how  $\alpha_{AI}$ ,  $H$ , and  $B$  can be jointly interpreted to distinguish **stable vs problematic ensemble behavior** and guide actionable next steps, such as **targeted data review** or **construct refinement**

# Limitations & Summary of Contribution

## Contribution

- **Interpretable uncertainty for LLM coding** that move beyond token-level confidence to diagnostics grounded in coding task semantics
- **Complementary diagnostics on agreement, decisiveness, and certainty** to characterize AI coding behavior on unlabeled data
- **Continuous deployment-time monitoring** that support ongoing human oversight of AI coding beyond standard test-set evaluation.

## Limitations and Future Work

- **Pilot-scale validation** (one coding scheme, one educational domain)
- **Computational cost of committee-based inference** is non-trivial; future work need to explore sparse and adaptive monitoring strategies.
- Diagnostics indicate **when further investigation is needed**, not yet **what intervention should follow**.

# Thank You! Questions?

Calibrate AI coder committee via  
iterative codebook refinement

[bit.ly/lak26-clarify](https://bit.ly/lak26-clarify)



Support human oversight of  
AI coding after deployment

[bit.ly/aied26-monitor](https://bit.ly/aied26-monitor)



**Fanjie Li**, Vanderbilt University, LIVE Learning Innovation Incubator