



HILAK 2026 Workshop, Bergen, Norway

## When AI Questions, Human CLARIFY:

Using RE-LLM Coding Uncertainty to Resolve Codebook Ambiguities



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



VANDERBILT  
UNIVERSITY



VANDERBILT

Learning Innovation Incubator



# Rethinking AI's Roles in Codebook-Based Qual Coding

1

## Human + AI

code better,  
not just faster

2

## Uncertainty as Signal

AI uncertainty reveals  
ambiguity in the codebook

3

## CLARIFY Workflow

illustrated w/ nursing  
reflection analytics

UNCERTAINTY AS SIGNAL

CLARIFY WORKFLOW

PILOT STUDY

LIVE

# Human-AI Partnership in LLM-Assisted Content Analytics

Traditional Approach	Hybrid Intelligence Approach
<p>Code Faster</p> <p>AI for Automation</p> <p>Replicates Human Coding</p> <p>Obscures Ambiguity</p>	<p>Code Better</p> <p>AI as Thinking Partner</p> <p>Stress-Tests Human Coding</p> <p>Surfaces Uncertainty, Diagnoses Ambiguity</p>

UNCERTAINTY AS SIGNAL

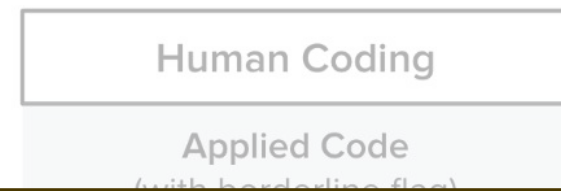
CLARIFY WORKFLOW

PILOT STUDY

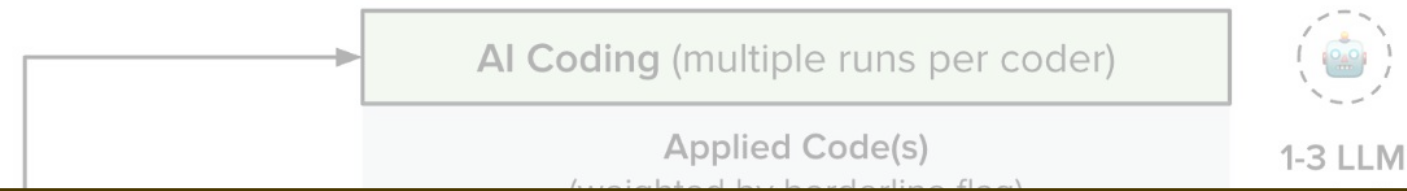


# CLARIFY Toolset and Workflow

Establish Initial Ground Truth

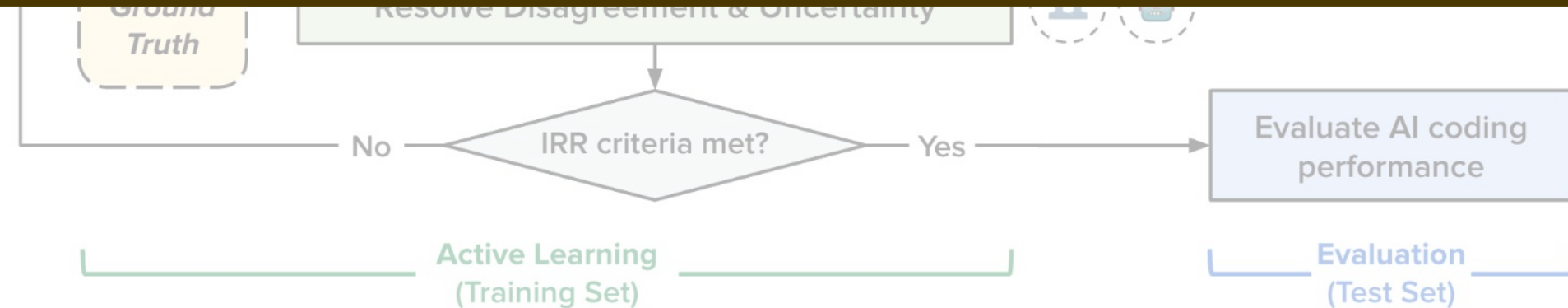


AI Coder Calibration with Human-in-the-Loop Validation



What if AI isn't just a tool for automation—  
but a partner that surfaces uncertainty  
and helps us reflect on ambiguity in our own thinking?

(Naive) Ground Truth



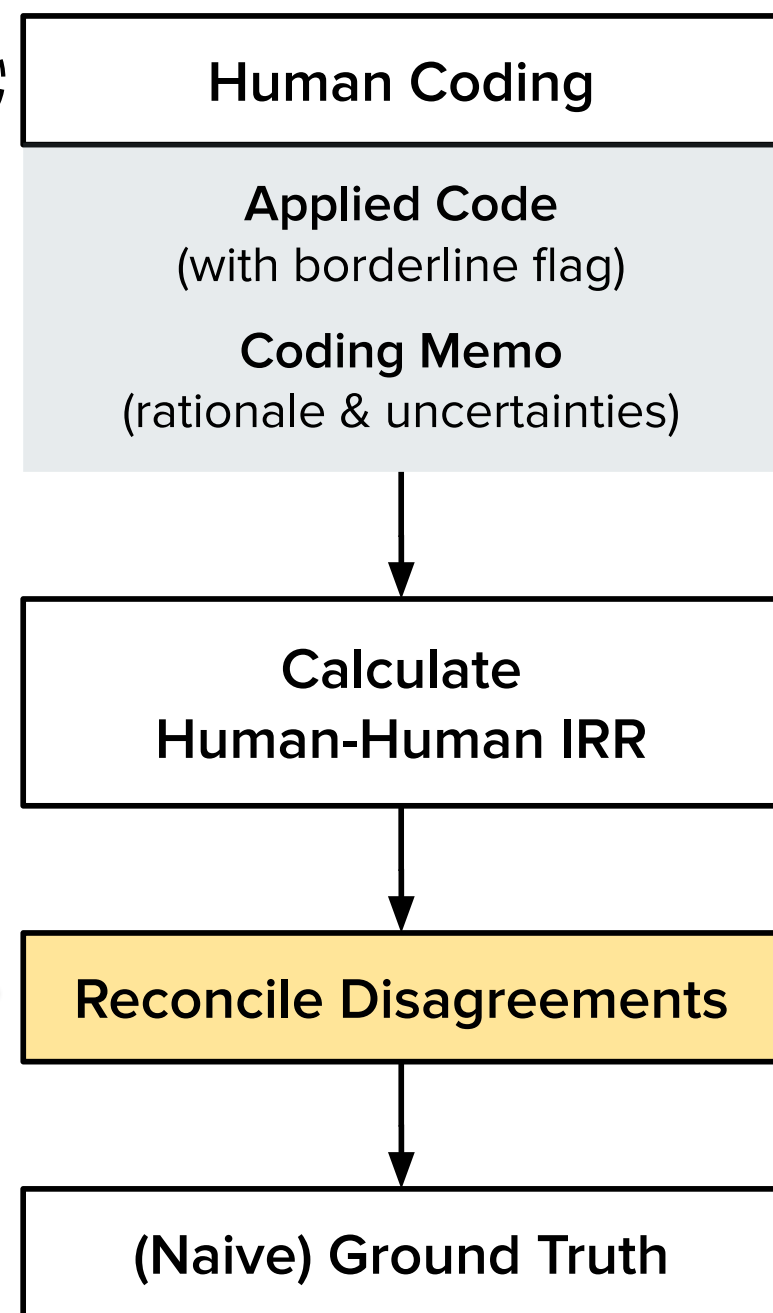
UNCERTAINTY AS SIGNAL

CLARIFY WORKFLOW

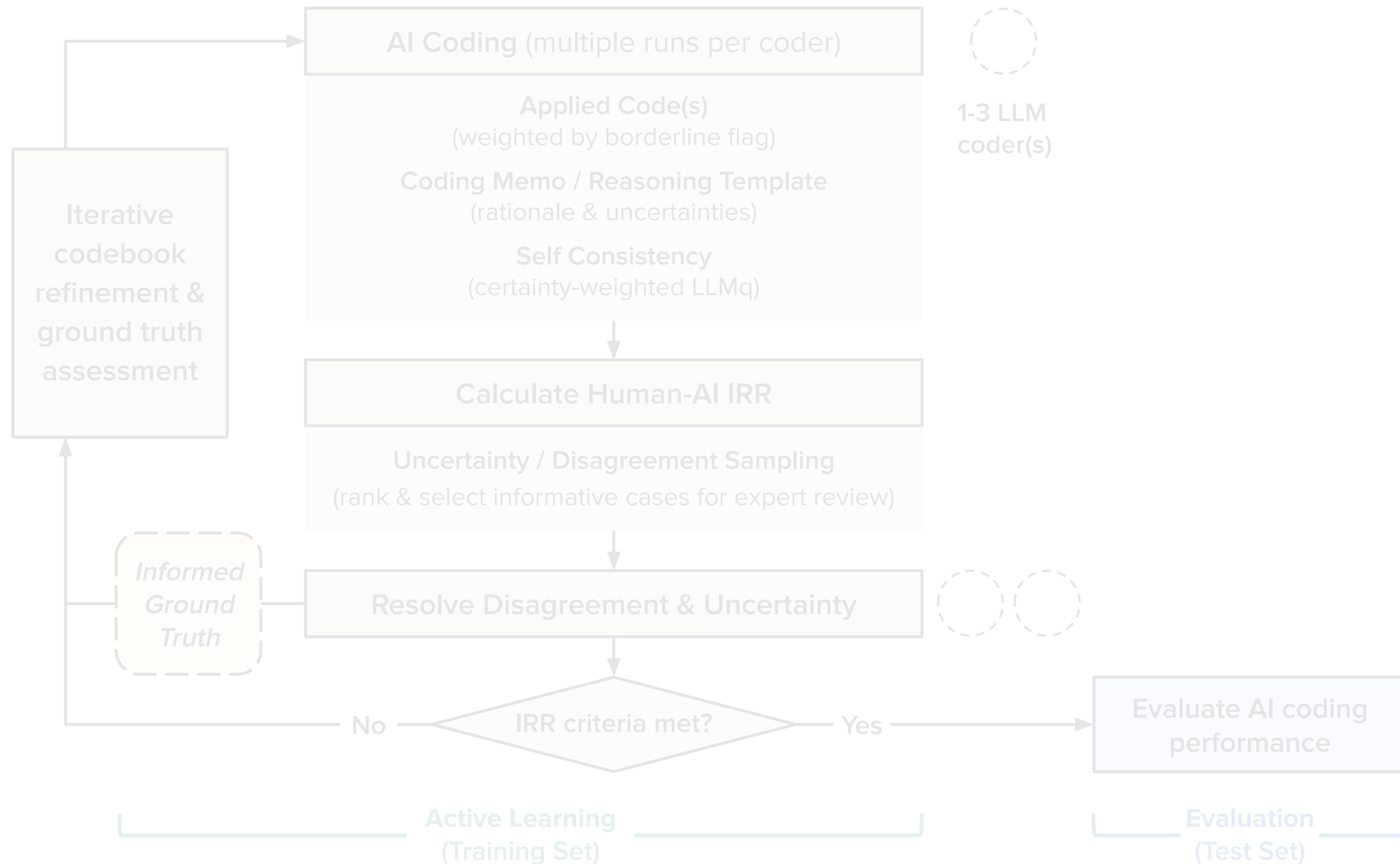
PILOT STUDY



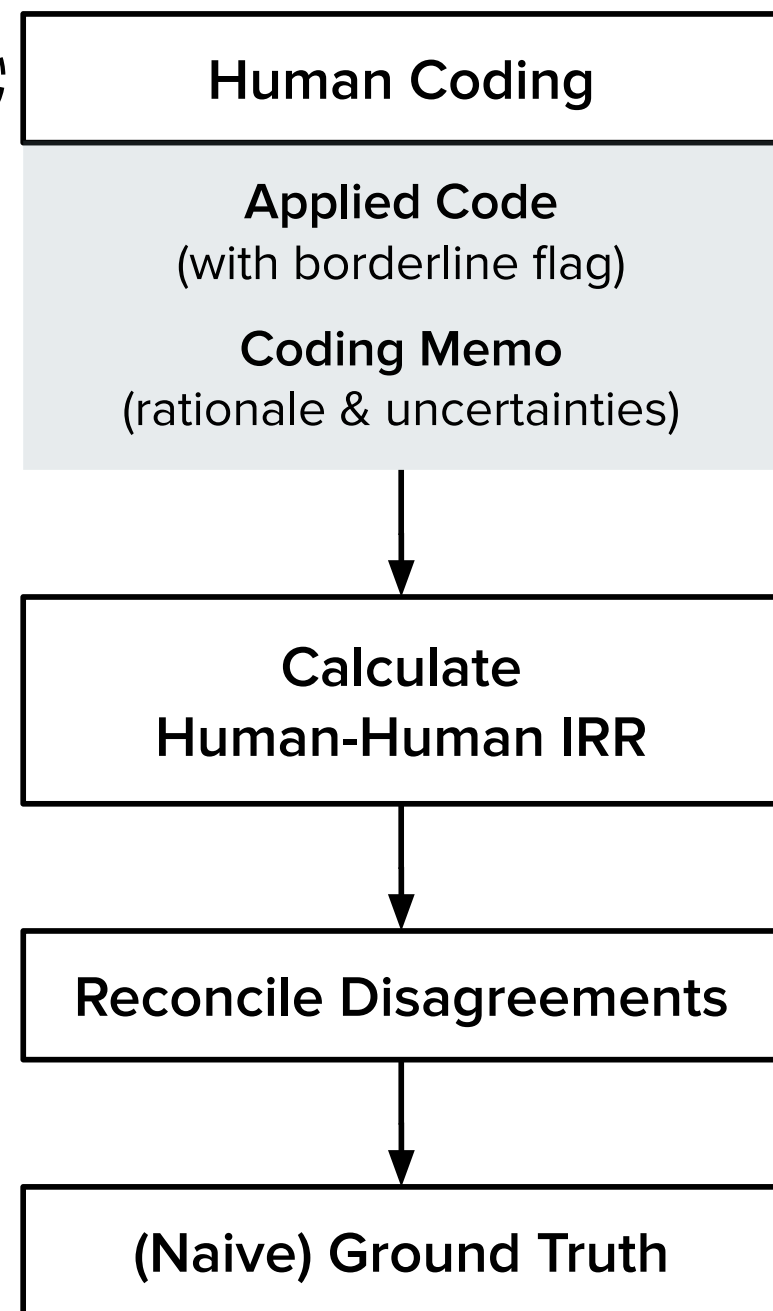
## Establish Initial Ground Truth



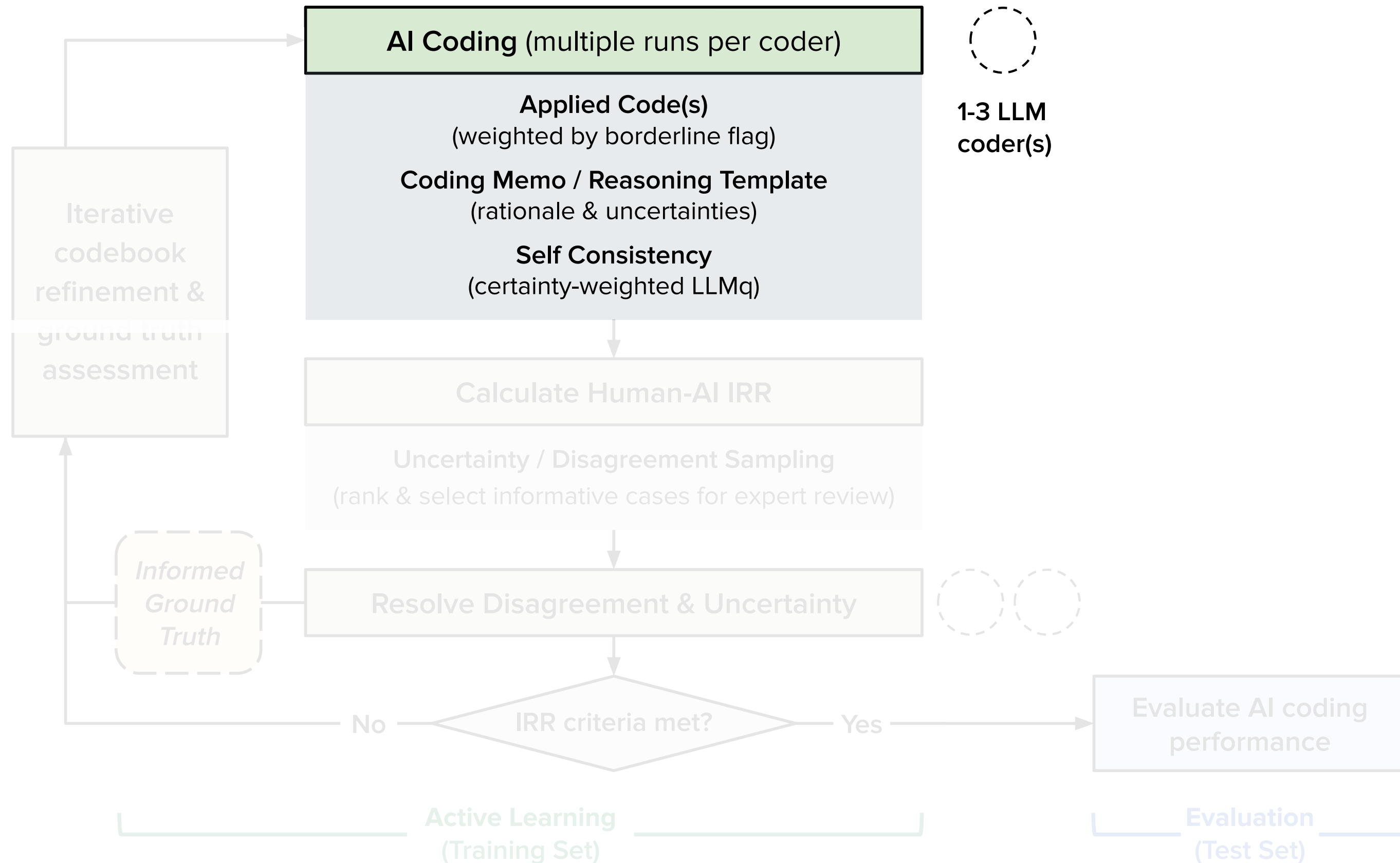
## AI Coder Calibration with Human-in-the-Loop Validation



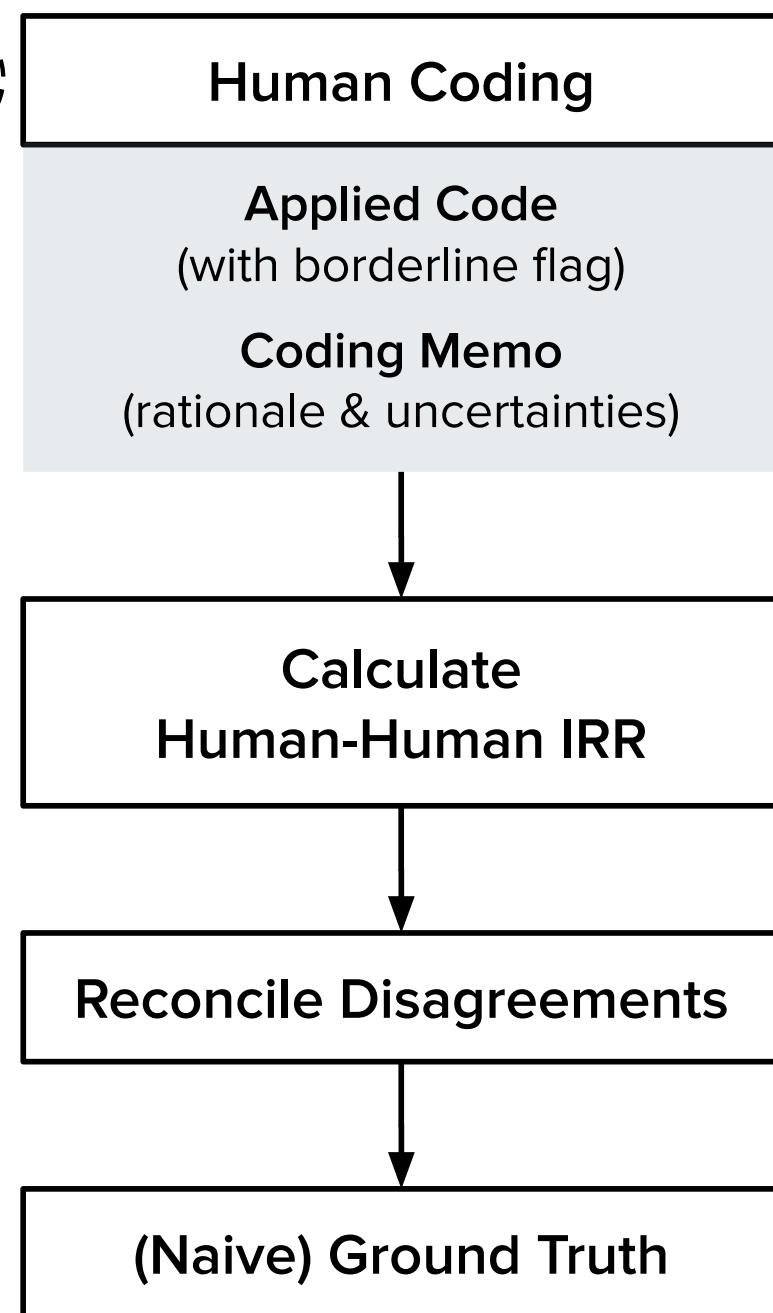
## Establish Initial Ground Truth



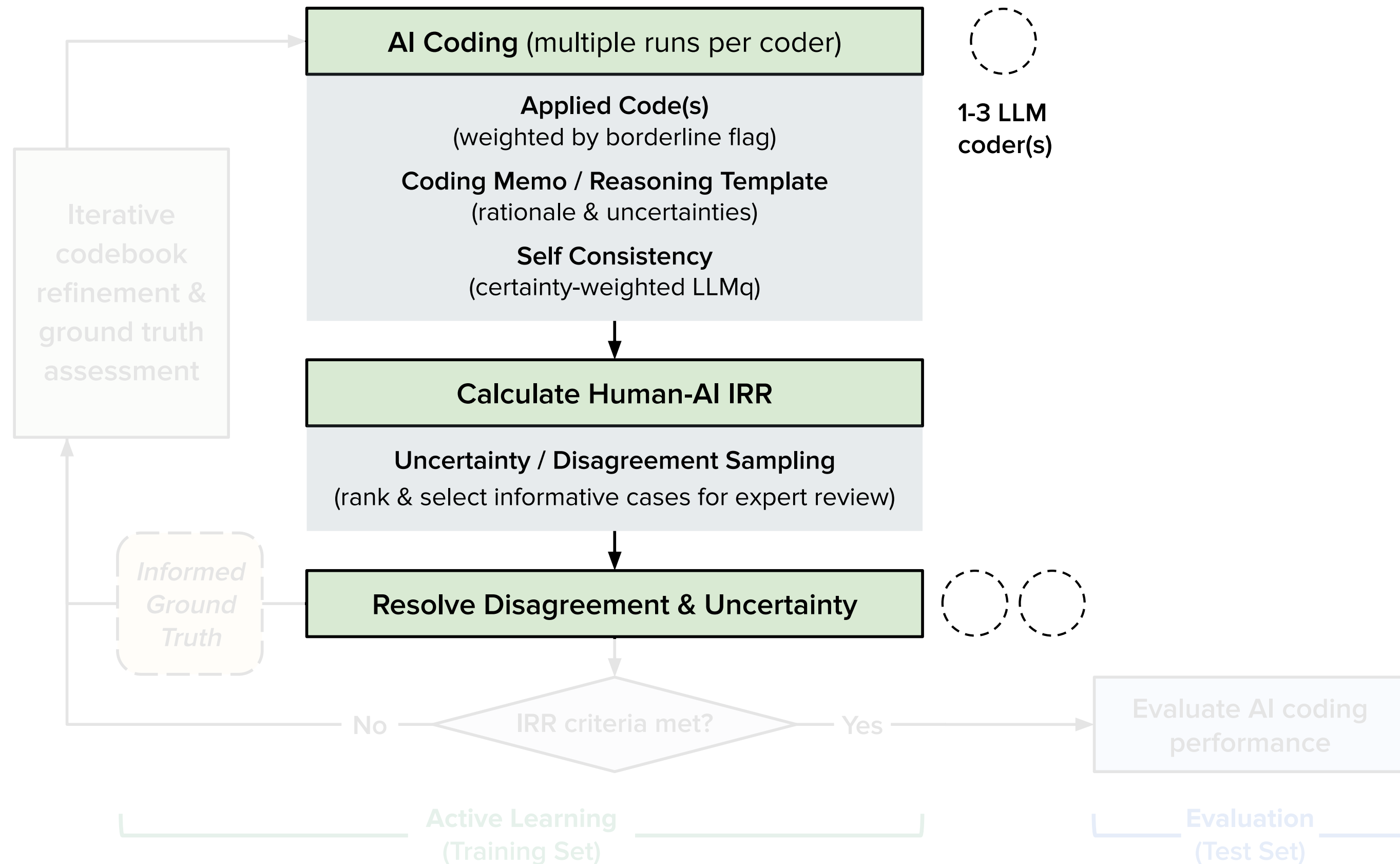
## AI Coder Calibration with Human-in-the-Loop Validation



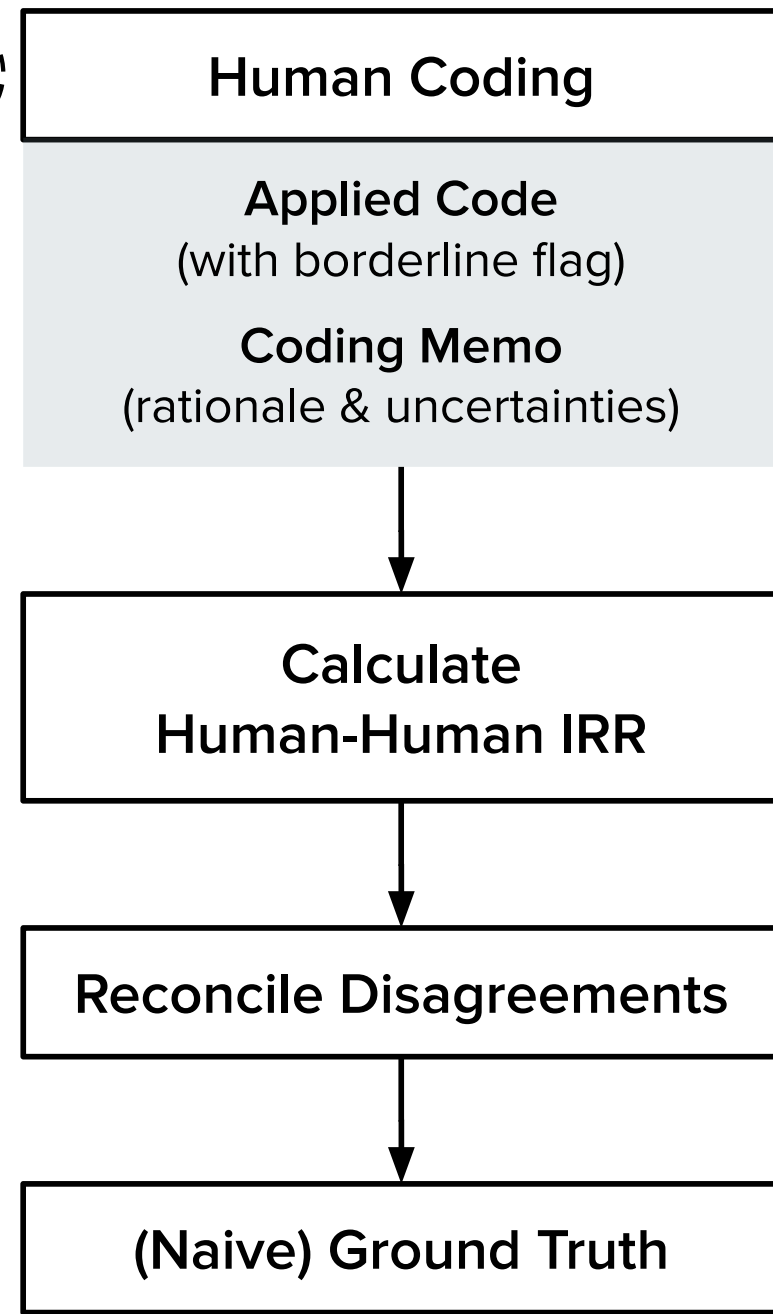
## Establish Initial Ground Truth



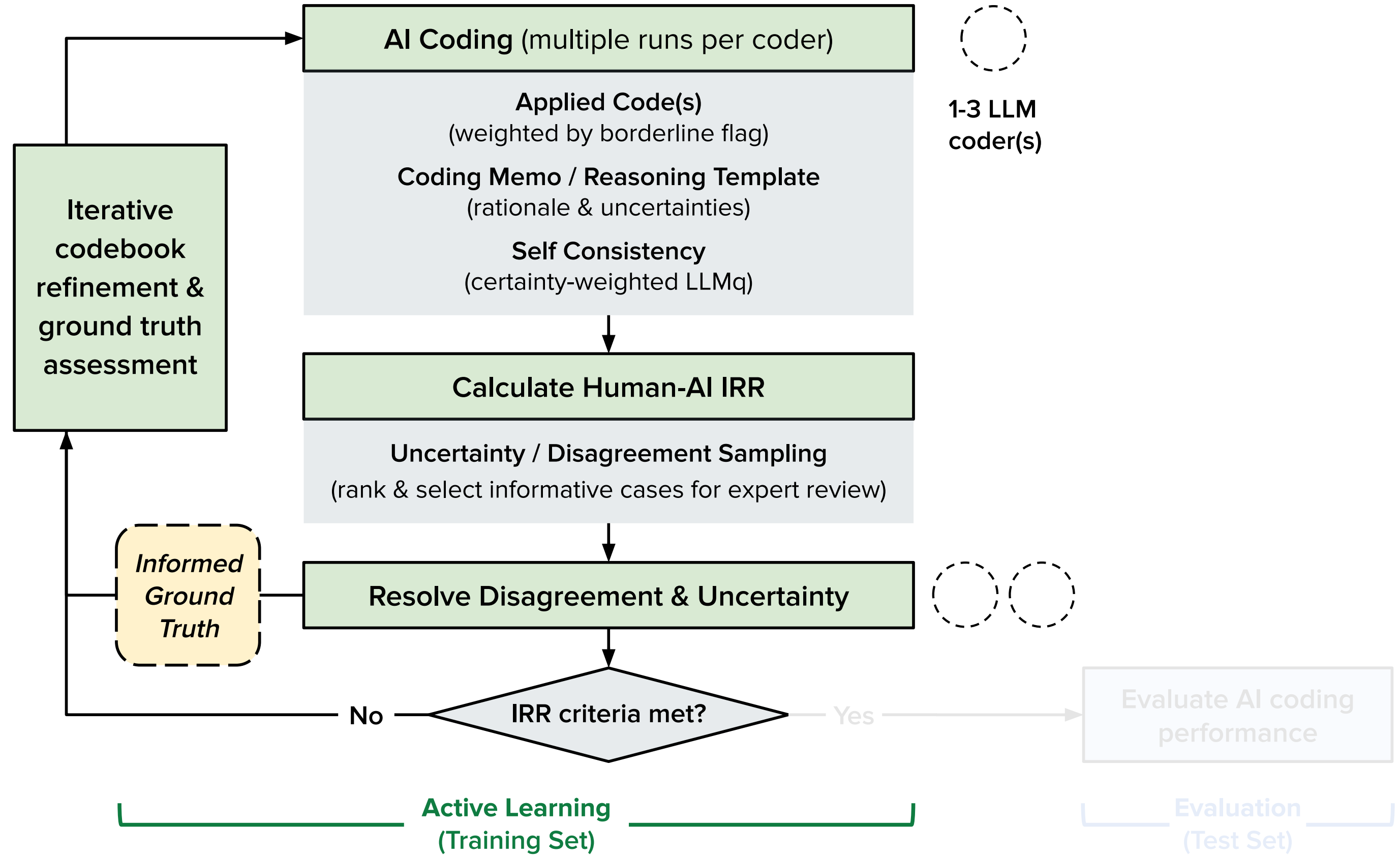
## AI Coder Calibration with Human-in-the-Loop Validation



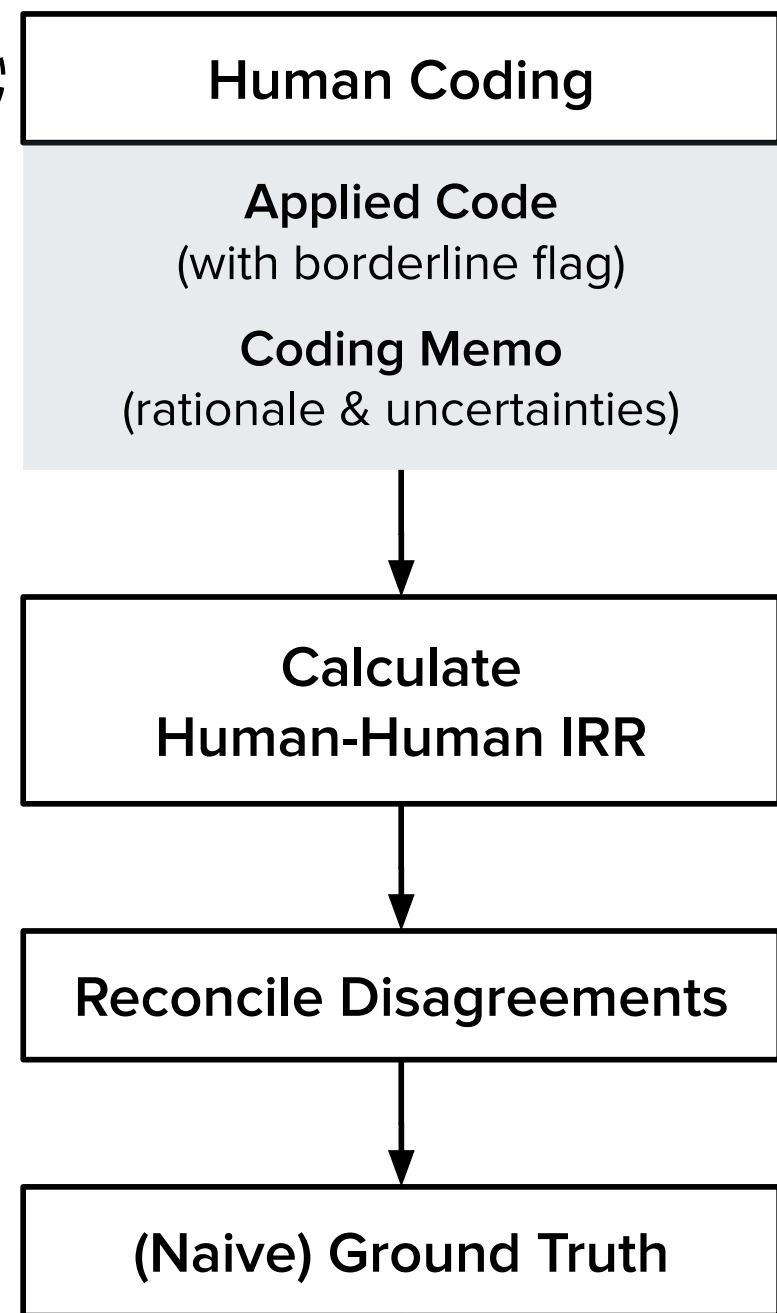
## Establish Initial Ground Truth



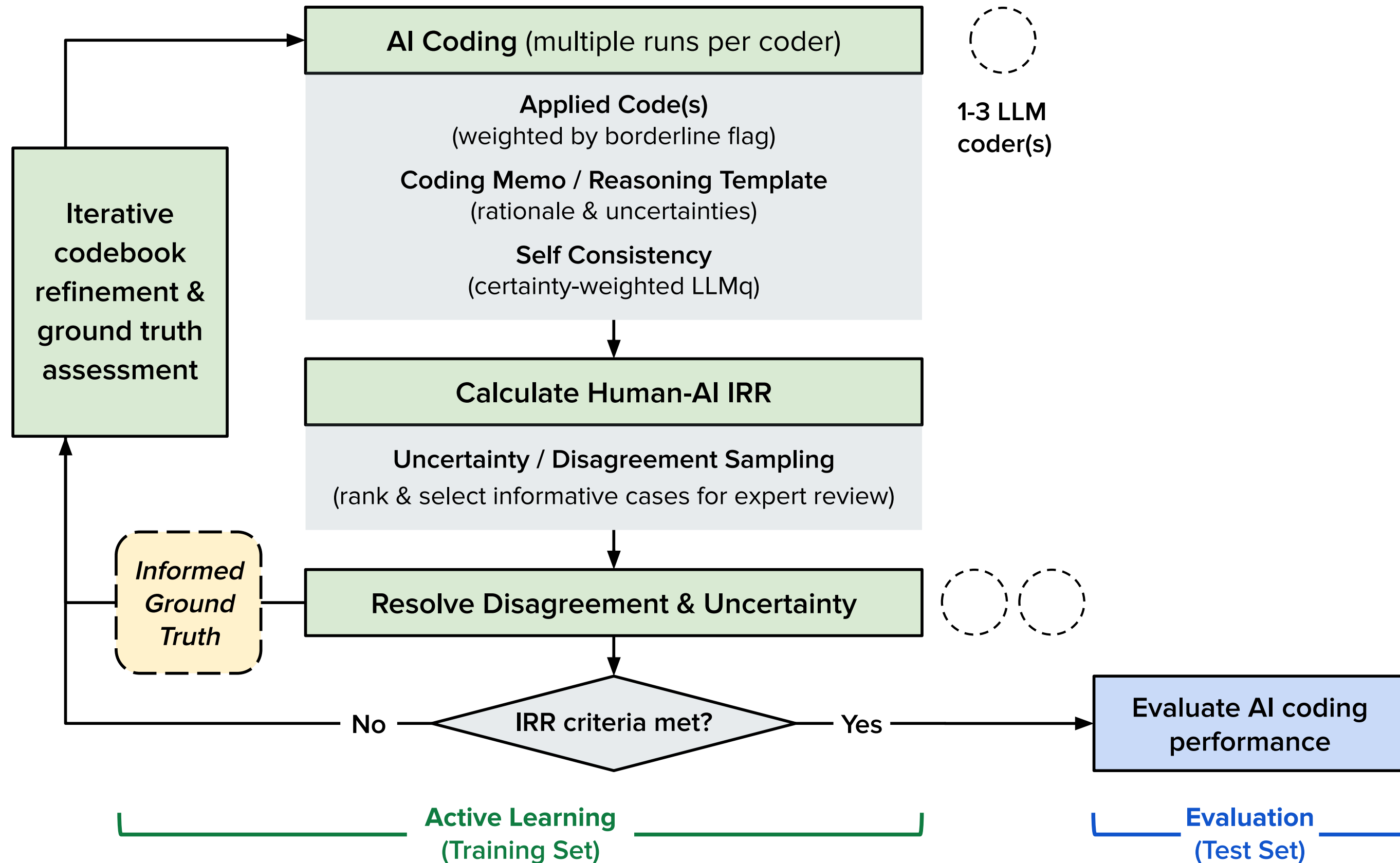
## AI Coder Calibration with Human-in-the-Loop Validation



### Establish Initial Ground Truth



### AI Coder Calibration with Human-in-the-Loop Validation



# Proof-of-Concept Pilot

## Context, Data and Coding Scheme

- **Nursing students' post-simulation written reflection** (213 event-level reflections, 16 students); Human coders double-coded 60 reflections (IRR:  $0.63 < \alpha < 0.91$ )
- **AI committee set up:**  
3 LLMs × 10 runs each
- **Focus code:**
  - **VALIDATE:** ensuring the correctness, safety, and readiness of a planned nursing intervention for implementation

### Codebook Overview (14 categories)

**Dimension 1: 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?"

**Dimension 2: How are they reflecting? (5 categories)**

- Are they evaluating how they performed and/or their knowledge or skill?
  - Positive evaluation / Negative evaluation
- Are they reporting challenging emotional experiences?
  - Cognitive-emotional discomfort (e.g., self-doubt, stress, frustration)
- Are they explaining what guided their reasoning — or how their reasoning (or gaps in reasoning) influenced care and patient outcomes?
  - Unpacking thinking
- Are they identifying what they will do differently or maintain in future practice?
  - Forward-looking reflection

💡 To code for *Dimension 2* (how are they reflecting?), focus on the student's (self-)evaluation, emotional reaction, explanation of reasoning, or the lessons learned for future practice, rather than the clinical task itself. Ask: "How is the student evaluating, explaining, or learning from their simulation experience?"

UNCERTAINTY AS SIGNAL

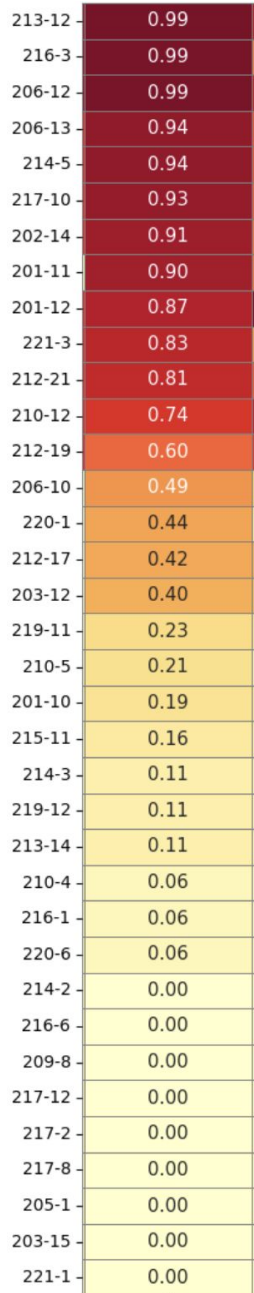
CLARIFY WORKFLOW

PILOT STUDY



# VALIDATE Code Trajectory

## Active learning round 1:



Human-AI agreement (Krippendorff's alpha) = 0.641

Number of cases with human-AI or AI-AI disagreement = 12

**Disagreement analysis:**

Review top 6 high-entropy cases (out of 36 training examples)

text id	consensus entropy	human	o3	claude	deepseek
213-12	0.99	1.00	0.79	0.09	0.30
216-3	0.99	0.70	0.94	0.31	0.27
206-13	0.94	0.70	1.00	0.31	0.56
214-5	0.94	0.70	1.00	0.57	0.30
217-10	0.93	0.00	0.77	0.44	0.18
206-12	0.93	0.70	1.00	0.63	0.30

**Note:** When human coders initially disagreed, reconciled labels were entered as .3 / .7 (rather than 0 / 1) to indicate the uncertainty identified.

# VALIDATE Code Trajectory

## Highest-entropy example in active learning round 1:

Reflection 213-12 | Coding Category: **VALIDATE** (Interpret Info to Validate Care Delivery Plan)

**Did this event segment contain a key moment that influenced the outcome of the simulation?**

"I was stressed with medication calculation."

 (o3)	1	1	1	1	1	1	1	1	1	1	→ LLMq-c <sub>(o3)</sub> = 0.79	} Ensemble Consensus Score <b>0.39</b>
 (claude)	0	0	0	0	0	0	0	0	0	0	→ LLMq-c <sub>(claude)</sub> = 0.09	
 (deepseek)	0	0	0	0	0	0	0	0	0	0	→ LLMq-c <sub>(deepseek)</sub> = 0.3	

**Sample Rationale (Code = 0):**

"While 'medication calculation' nominally connects to an indicator, **the lack of detail about validation behaviors creates uncertainty**. Codebook requires evidence of cognitive work to validate care plans, not just acknowledgement that calculations occurred. **However, one could argue the mere act of doing calculations implies validation attempt.**"

**Sample Rationale (Code = 1):**

"The statement is extremely brief and does not explicitly describe performing or checking the calculation, merely being stressed. **One could argue it does not provide enough detail about the mental process.** Nonetheless, the **explicit reference to medication calculation—as a key moment—provides minimal but sufficient evidence** under a generous reading of the codebook. Borderline but leaning present."

# VALIDATE Code Trajectory

**Added heuristics:** Apply this code when students are reflecting on interpreting information to validate care delivery plan, regardless of how they are reflecting on it. That is, **reflection on interpreting information to validate care delivery plan can take various forms, including** but not limited to describing what they did or failed to do during the simulation, evaluating their performance, **expressing feelings or challenges about the experience**, identifying lessons learned, and planning for future practice. All of these reflection types qualify as evidence for coding.

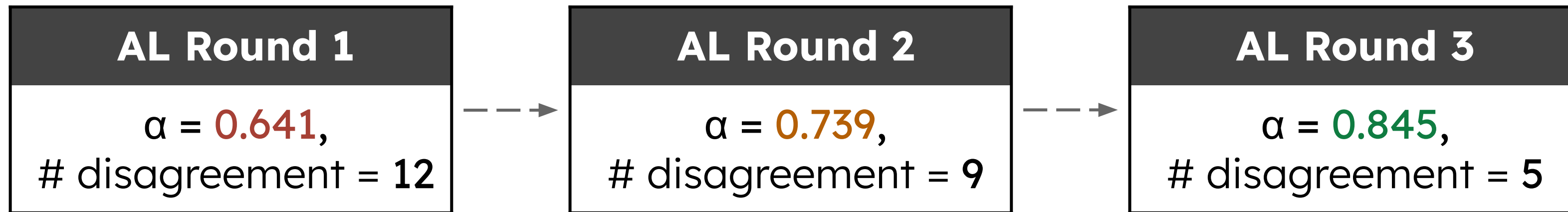
 (o3)	1	1	1	1	1	1	1	1	1	1	1	→ LLMq-c <sub>(o3)</sub> = 1.00	} Ensemble Consensus Score <b>0.97</b>
 (claude)	1	1	1	1	1	1	1	1	1	1	1	→ LLMq-c <sub>(claude)</sub> = 1.00	
 (deepseek)	1	1 <i>flag</i>	1	1	1	1	1 <i>flag</i>	1	1 <i>flag</i>	1	1	→ LLMq-c <sub>(deepseek)</sub> = 0.91	

**Sample Rationale BEFORE Codebook Clarification:** "While 'medication calculation' nominally connects to an indicator, the lack of detail about validation behaviors creates uncertainty. Codebook requires evidence of cognitive work to validate care plans, not just acknowledgement that calculations occurred. However, one could argue the mere act of doing calculations implies validation attempt."

**Sample Rationale AFTER Codebook Clarification:** "Minor uncertainty exists because the student does not explicitly describe the validation process (e.g., verifying dosage accuracy) but only references stress associated with the task. **However, the codebook's heuristics emphasize that outcomes or emotional reactions related to validation work still qualify as evidence.**"

# VALIDATE Code Trajectory

Tracking performance across three **active learning (AL)** rounds:



Performance on **held-out test set**:



UNCERTAINTY AS SIGNAL

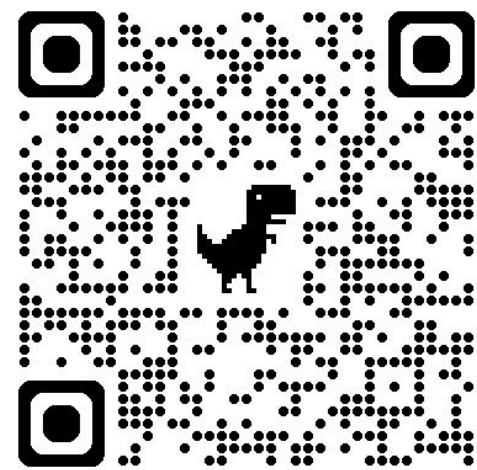
CLARIFY WORKFLOW

PILOT STUDY



# Concluding Thoughts

Scan for full paper  
[bit.ly/lak26-clarify](https://bit.ly/lak26-clarify)



Objectives	AI roles	Human roles
AI coding as a supplementary verification layer	Detect drift in code application across the entire dataset (Zambrano et al., 2023)	Resolve inconsistencies observed through human-in-the-loop validation (Cohn et al., 2024; Ramanathan et al., 2025; Zambrano et al., 2023)
	Assess intra-coder consistency through recursive model queries (Ramanathan et al., 2025; Tai et al., 2024) Use convergence of codes to quantify decision confidence (Wang et al., 2023)	Evaluate trends in code stabilization (Tai et al., 2024; Ramanathan et al., 2025) Review low-confidence coding decisions (Li et al., 2026)
AI coder(s) as collaborative partners for identifying and resolving gaps in the current coding scheme	Articulate reasoning behind coding decisions, which either prompts human reflection on their rationales and assumptions or exposes unclear code definitions that cause coder confusions (Ramanathan et al., 2025; Zambrano et al., 2023; Cohn et al., 2024)	Perform disagreement analysis and resolve AI coder confusions through codebook improvements and/or few-shot examples with human-corrected reasoning (Ramanathan et al., 2025; Zambrano et al., 2023; Cohn et al., 2024)
	Provide suggestions for updating the code definition upon human request (Zambrano et al., 2023)	Decide whether to accept/reject changes suggested by LLM (Zambrano et al., 2023; Ramanathan et al., 2025)