# EvalAssist:
# Insights on Task-Specific Evaluations and AI-Assisted Judgment Strategy Preferences
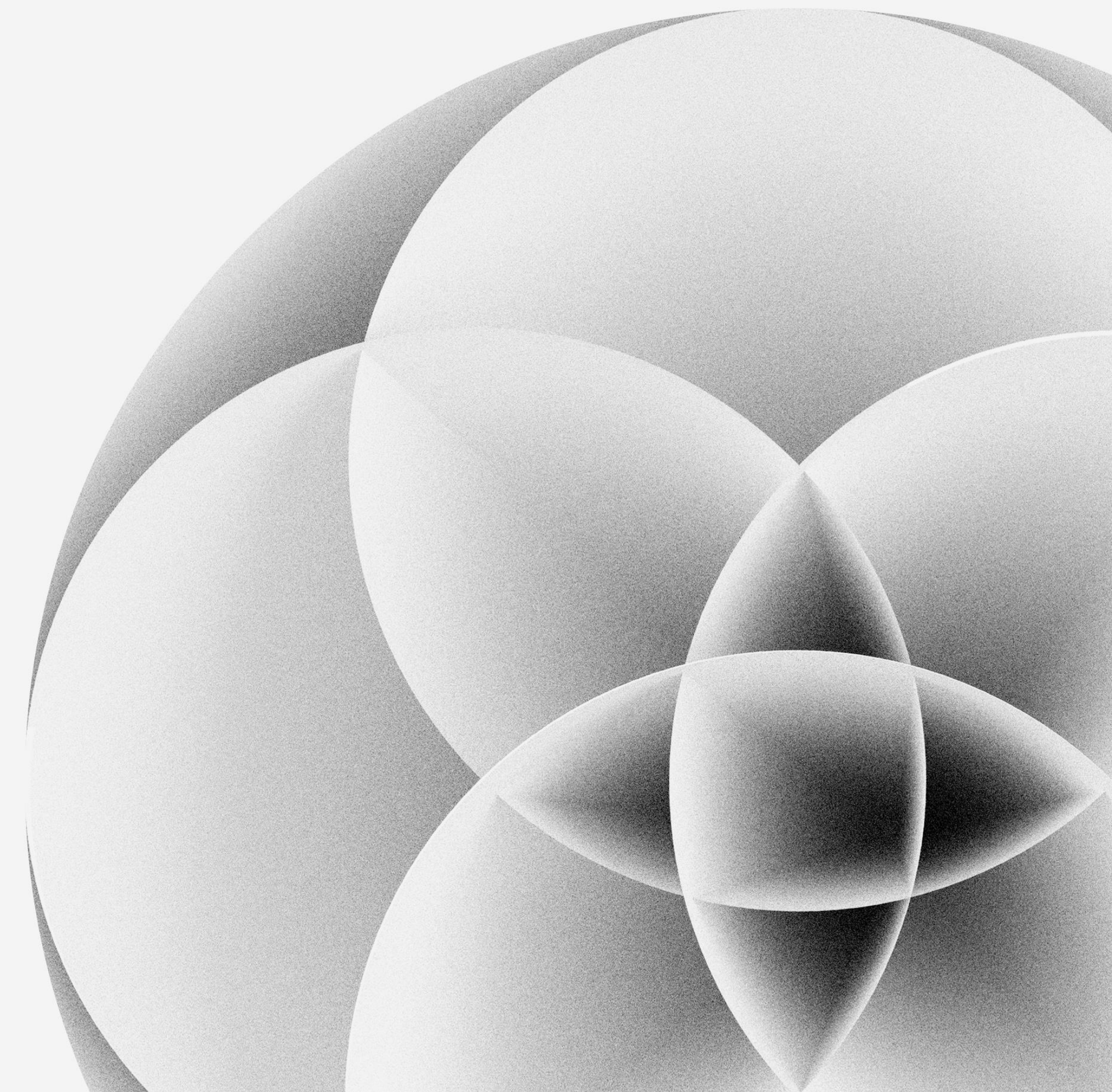
Presenter: Hyo Jin (Gina) Do (hjdo@ibm.com)

Authors: Zahra Ashktorab, Michael Desmond, Qian Pan, James M. Johnson, Martín Santillán Cooper, Elizabeth M. Daly, Rahul Nair, Tejaswini Pedapati, Hyo Jin Do, Werner Geyer

IBM Research

UIST 25

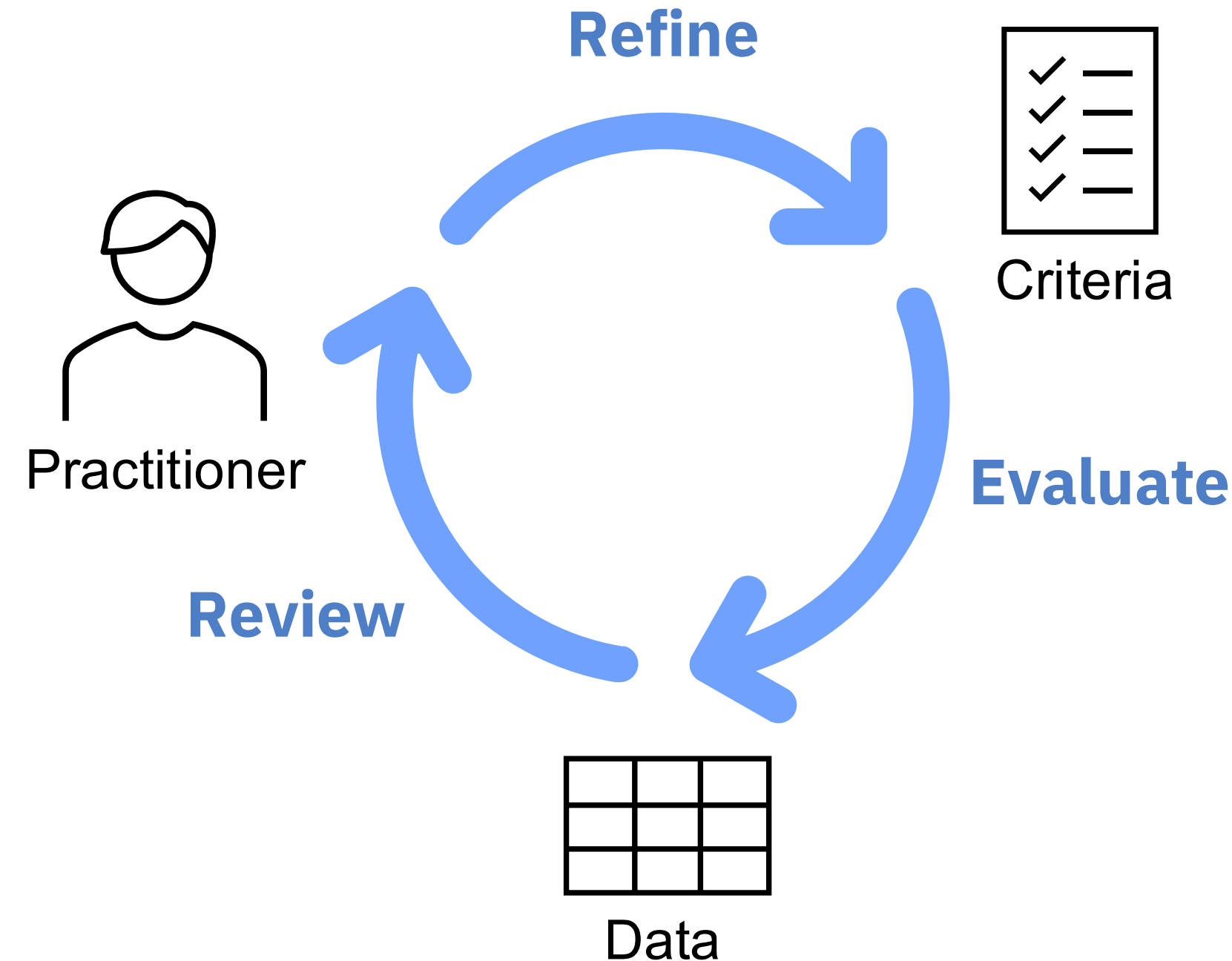BUSAN, KOREA | SEP. 28TH – OCT. 1ST 2025

# LLM-as-a-Judge

- **Manual evaluation** to assess LLM outputs is expensive, slow, and does not scale.
- LLMs are increasingly used as evaluators (i.e., **LLM-as-a-judge**).
- Offers flexibility through custom, task-specific **evaluation criteria.**



Credit: https://towardsdatascience.com/llm-as-a-judge-a-practical-guide/

# Human-in-the-Loop Evaluation

- "Criteria Drift" (Shankar et al., 2024): While predefined criteria may help users assess outputs, the act of grading also helps users to refine those criteria.
- Many LLM-as-a-judge systems support human-in-the-loop evaluation, integrating human feedback to improve human-AI alignment.
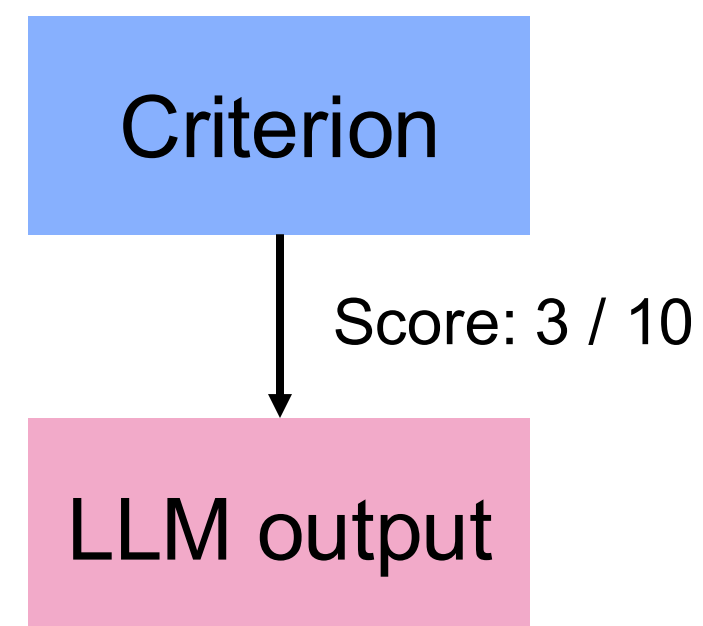
# Design Goals

1. Isolate generation from evaluation
2. Evaluate a subset of data first to reduce costs
3. Support multiple LLM evaluators
4. Include a positional bias indicator (i.e., a model favors one option based on its position)
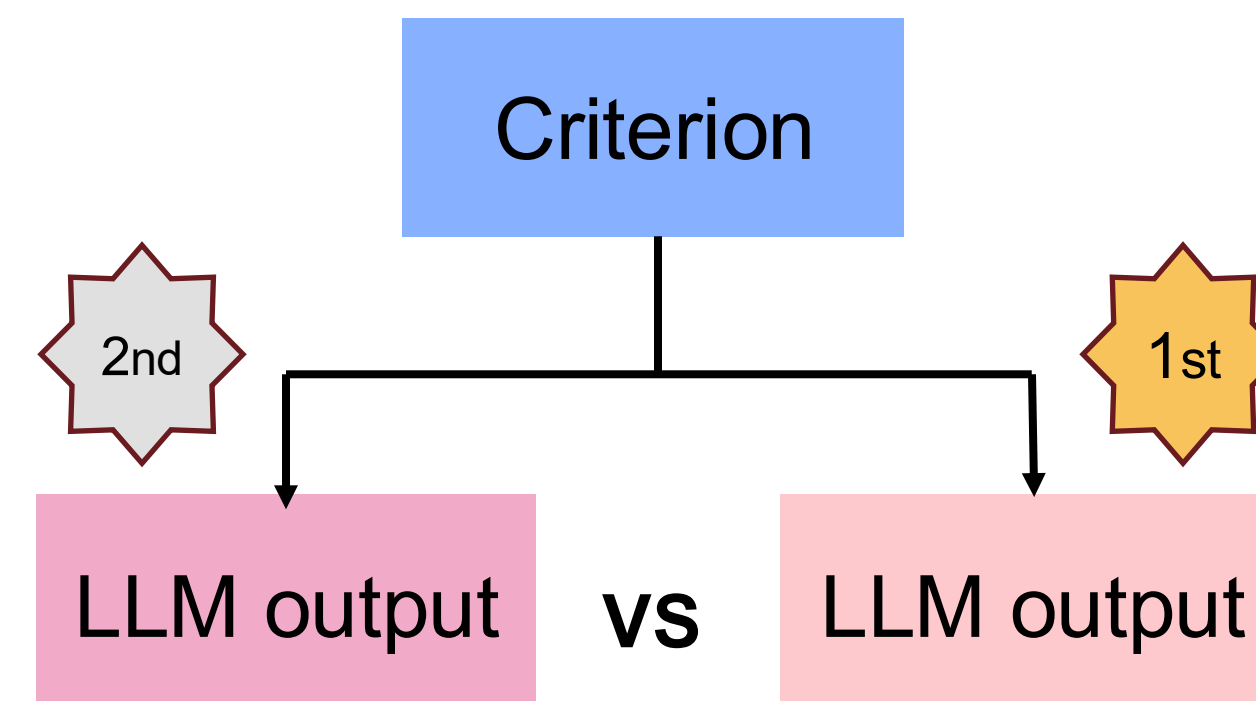5. Support two common LLM-as-a-judge strategies:

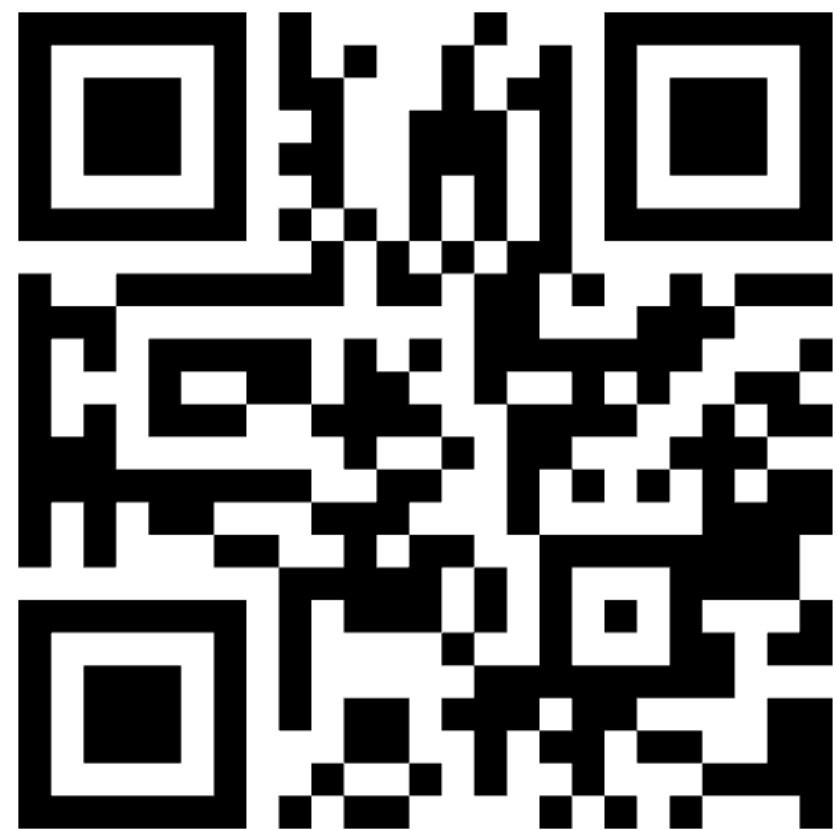**Direct Assessment**
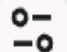Outputs are graded based on a rubric

**Pairwise Comparison**
Outputs are ranked by comparing outputs

# EvalAssist

Try out!



Note: EvalAssist has been updated since we completed this study, so some of the slides may reflect an older version of the system.

# Select Evaluation Strategy

# Define Criteria

## Direct Assessment:

**Evaluation Criteria**

| Form | JSON |

### Summary Quality ✎

Is the summary accurate, clear and coherent?

| Option | Option |
| --- | --- |
| Yes | The summary is accurate, clear and covers all the important elements in the original text. |
| Option | Option |
| No | The summary is inaccurate and includes unnecessary details or repetition. |

Add Option +

## Pairwise Comparison:

**Evaluation Criteria**

| Form | JSON |

### Summary Quality ✎

Criteria

The summary is accurate, clear and covers all the important elements in the original text.

# Set Task Context Variables

**Summary Quality** ✎   Direct Assessment   Save 🖫   Save as 🗐   New Test Case +   Delete Test Case 🗑

**Task context (optional)**

| Name | Value |
|---|---|
| instruction ✎ | Summarize this article: |
| article ✎ | A dress worn by Vivien Leigh when she played Scarlett O'Hara in the classic 1939 film Gone With the Wind has fetched $ 137,000 at auction . Heritage Auctions offered the gray jacket and skirt , featuring a black zigzag applique , plus more than 150 other items from the Academy Award-winning film at auction on Saturday in Beverly Hills , California . The dress - a jacket and full skirt ensemble - was worn in several key scenes in the 1939 movie , including when Scarlett O'Hara encounters Rhett Butler , played by Clark Gable , and when she gets attacked in the shanty town . Scroll down for video An outfit worn in several scenes of the 1939 film Gone With The Wind by Vivien Leigh as she played Scarlett O'Hara sold for $ 137,000 at auction on Saturday The dress - a jacket and full skirt ensemble - was worn in several key scenes in the 1939 movie but has suffered a little with age and has faded to light gray from original slate blue-gray color The outfit has suffered a little with age , however . When Leigh wore it in the movie , it was slate blue-gray but over the years it has faded to light gray . It was one of more than 150 items that were part of the private collection of James Tumblin , formerly in charge of the hair and makeup department at Universal Studios . Tumblin began collecting onscreen costumes , props and behind-the-scenes artifacts from the film in the 1960s , amassing a collection of more than 300,000 pieces of memorabilia . During a visit to the Western Costume Company he spotted the Scarlett O'Hara dress on the floor . He learned that the dress was about to be thrown away and negotiated a deal to buy it for $ 20 . Tumblin has 'devoted his life and efforts to promoting Hollywood and this film , touring his items throughout the United States , ' said Kathleen Guzman , managing director of Heritage Auctions . Gone With The Wind , which celebrated its 75th anniversary last year , was based on Margaret Mitchell 's 1936 best-selling book about a spoiled Old South socialite , played by Vivien Leigh , and co-starred Clark gable as Rhett Butler Hattie McDaniel ( left ) , Olivia DeHavilland ( middle ) , and Vivien Leigh : McDaniel famously became the first African-American actor to be nominated for and win an Academy Award Other top selling items from the auction were a straw hat worn by Leigh that sold for $ 52,500 ; the trousers and jacket from a suit worn by Clark Gable as Rhett Butler , selling for $ 55,000 ; and a black bonnet worn by both Leigh and Olivia De Havilland as Melanie Wilkes , which fetched $ 30,000 . Gone With The Wind , which celebrated its 75th anniversary last year , was based on Margaret Mitchell 's 1936 best-selling book about a spoiled Old South socialite . Actress Hattie McDaniel , who played Scarlett 's devoted nanny Mammy , a slave , famously became the first African-American actor to be nominated for and win an Academy Award . |

Add variable +

8

# Select Judge Model

API credentials 🔑

**Model selection**

Evaluator

Llama3.3-70b ⌄

Granite3.0-8b
Granite3.2-8b
Granite3.3-8b
Llama3.1-8b
Llama3.1-70b
Llama3.1-405b
✓ Llama3.3-70b
Mixtral Large
Mixtral Medium
Mixtral Small
Mixtral8-7b

# Evaluate Test Data (Direct Assessment)

**Test data**

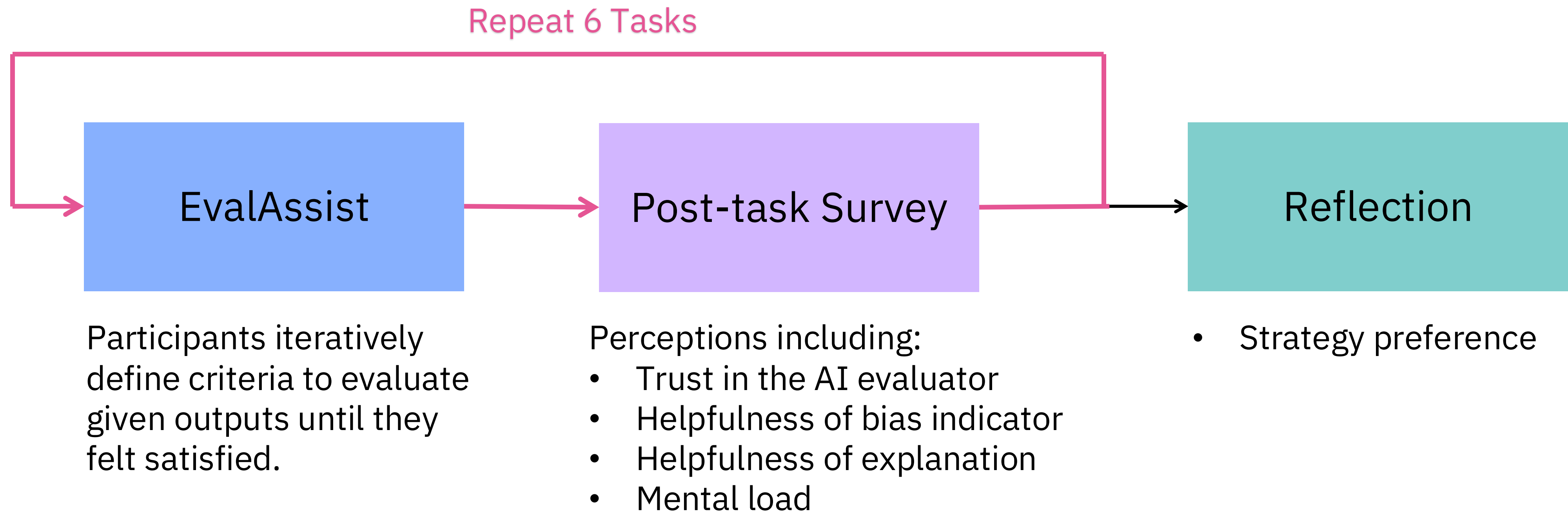| Responses to evaluate  summary ✎ | Expected result | Result | Explanation |
|---|---|---|---|
| Heritage Auctions offered the gray jacket and skirt , featuring a black zigzag applique , plus more than 150 other items from the Academy Award-winning film at auction on Saturday . The dress - a jacket and full skirt ensemble - was worn in several key scenes in the 1939 movie but has suffered a little with age and has faded to light gray . | No ⌄ | **No**<br>Agreement: Yes | The provided summary is inadequate, failing to accurately capture the essence of the original text. It omits crucial details, such as the auction price and the collector's background, and lacks clarity and coherence, making it unclear and disjointed. |
| A dress worn by vivien leigh when she played scarlett O'Hara in the classic 1939 film gone with the wind has fetched $ 137,000 at auction . Heritage auctions offered the gray jacket and skirt , featuring a black zigzag applique , plus more than 150 other items from the academy award-winning film at auction on saturday in beverly hills , california . The dress - a jacket and full skirt ensemble - was worn in several key scenes in the 1939 movie , including | Yes ⌄ | **No**<br>Agreement: No<br><br>Positional bias detected | The provided summary is inadequate, lacking accuracy, clarity, and coherence. It merely copies the first two paragraphs of the original text, failing to provide a concise overview of the main points. The response includes unnecessary details, repetition, and omits crucial elements from the |
| Heritage auctions offered the gray jacket and a black zigzag applique . The dress was worn in the classic 1939 movie . The wind has fetched $ 137,000 at auction . The heritage auctions offer the gray jackets and $ 150 more than 150 items . The academy award-winning film at auction on saturday in beverly hills , california . | No ⌄ | **No**<br>Agreement: Yes | The provided summary is inadequate, lacking clarity and coherence. It omits crucial information from the original text, includes unnecessary details, and repeats itself. Most importantly, it fails to convey the main idea of the article, which is about the auction of Vivien Leigh's iconic dress |
| Dress worn by vivien leigh when she played scarlett O'Hara in 1939 film gone with the wind has fetched $ 137,000 at auction . The dress - a jacket and full skirt ensemble - was worn in several key scenes in the 1939 movie , including when scarlett O'Hara encounters rhett butler , played by clark gable , and when she gets attacked in the shanty town . | No ⌄ | **No**<br>Agreement: Yes | The provided summary is inadequate, lacking clarity, coherence, and essential details. It fails to cover the majority of the original text, omitting crucial information such as the auction house, collection origin, dress condition, and other auction items, instead only repeating select |

Evaluate

# Experimental Design

- 15 participants who had prior experience with model evaluation
- Within-subjects study design
- 6 Tasks: 2 Strategies (direct, pairwise) x 3 Types (Summarization, Email, Q&A)

| Type | Criteria | Example Output (shortened) |
|------|----------|----------------------------|
| **Summarization** (Fabbri et al., 2021) | **Preference** | Dress worn by Vivien Leigh when she played Scarlett O'Hara in the 1939 film *Gone with the Wind* has fetched $137,000 at auction. The dress—a jacket and full skirt ensemble—was worn in several key scenes in the movie, including when Scarlett encounters Rhett Butler, played by Clark Gable, and when she is attacked in the shanty town. |
| **Email** about an office Christmas party | **Inclusivity** | Subject: Join Us for a Multicultural Holiday Celebration! Dear Team, Our celebration will feature a variety of activities that honor different cultural traditions, including decorations from various holidays, inclusive gift exchange, international feast. We would love to hear your ideas and suggestions to make this event even more inclusive and enjoyable. |
| **Q&A** from an HR chatbot | **Faithfulness** | You can upgrade to business class by paying the business class portion of the airline ticket with a personal credit card to avoid the perception of a policy bypass. |

# Experimental Procedure

Repeat 6 Tasks

EvalAssist

Post-task Survey

Reflection

Participants iteratively
define criteria to evaluate
given outputs until they
felt satisfied.

Perceptions including:
- Trust in the AI evaluator
- Helpfulness of bias indicator
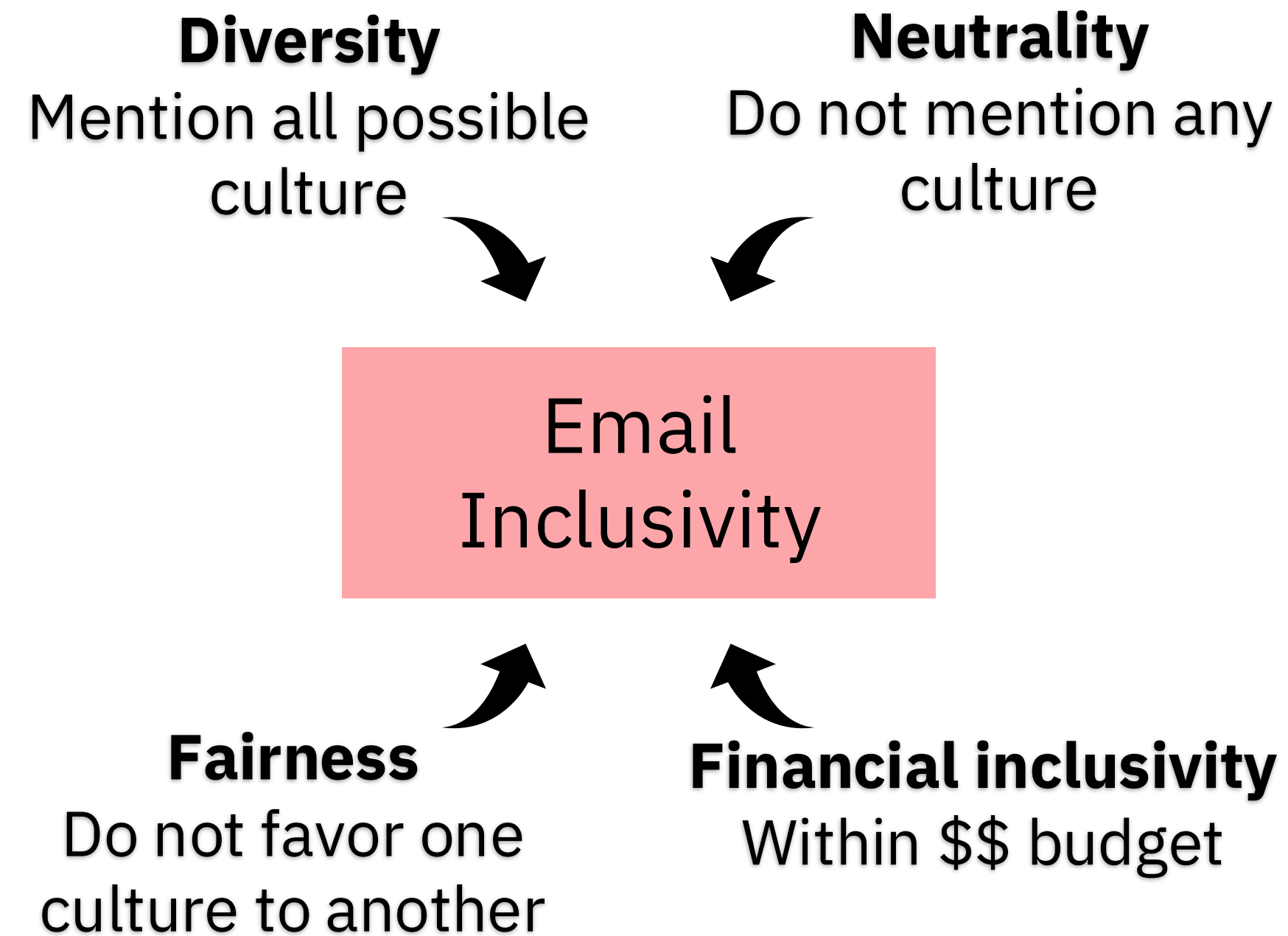- Helpfulness of explanation
- Mental load

- Strategy preference

# Variation in Criteria Development

- Significant variation in how participants defined criteria, including ineffective ones (e.g., providing examples and additional prompts)
- The most common strategy was to make criteria definitions **specific** (37%), leading to overly narrow criteria.

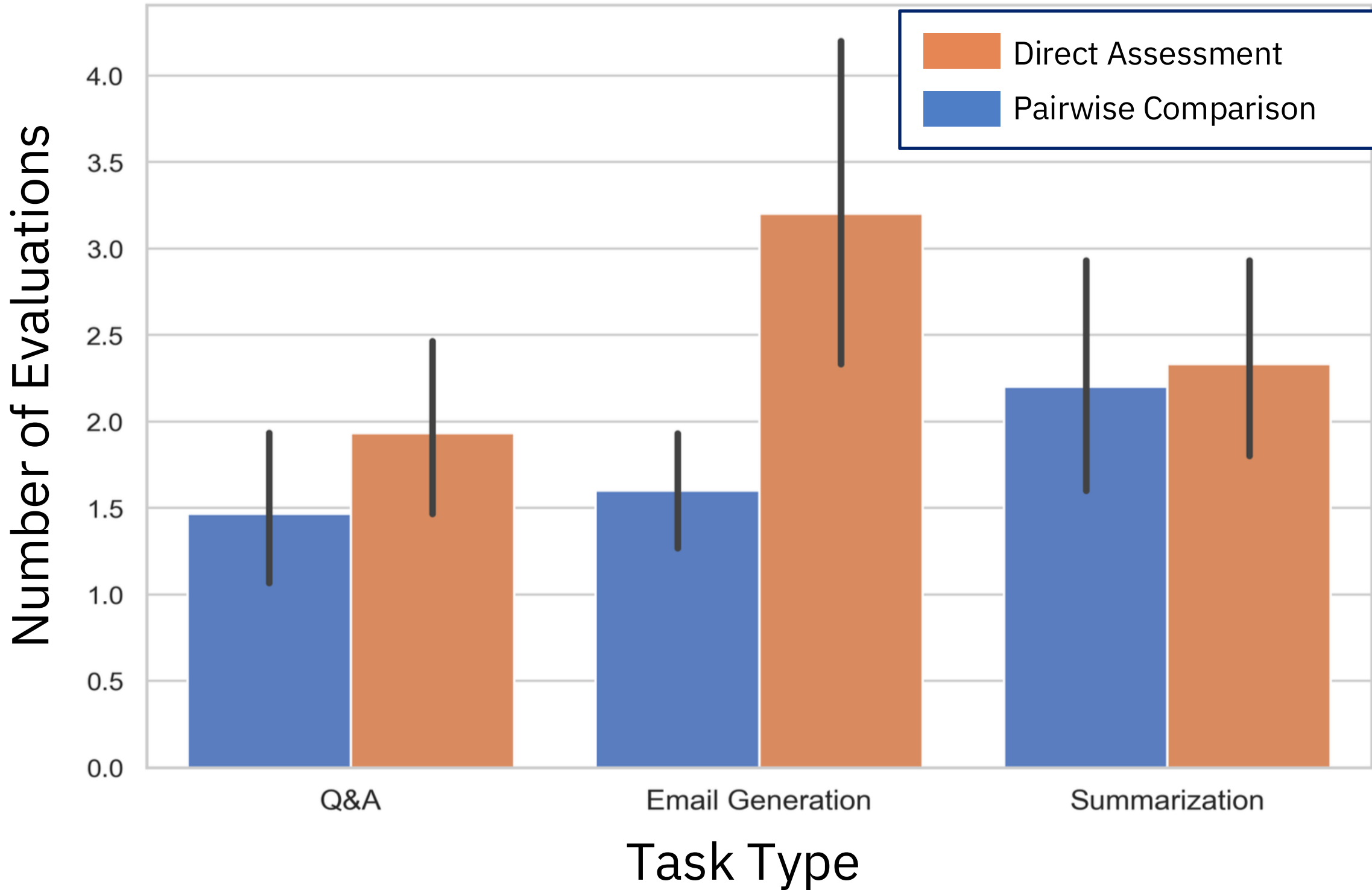*Inclusivity Criterion Example:*
*Please evaluate whether the following E-Mail is inclusive. This means that not only western traditions, such as Christmas, are celebrated, but employees are actively asked to contribute their customs and traditions to contribute to a diverse and inclusive company culture. Please also assess whether inclusive language is being used throughout the E-Mail.*

# Varied Prioritization of Evaluative Dimensions

**Diversity**
Mention all possible culture

**Neutrality**
Do not mention any culture

Email Inclusivity

**Fairness**
Do not favor one culture to another

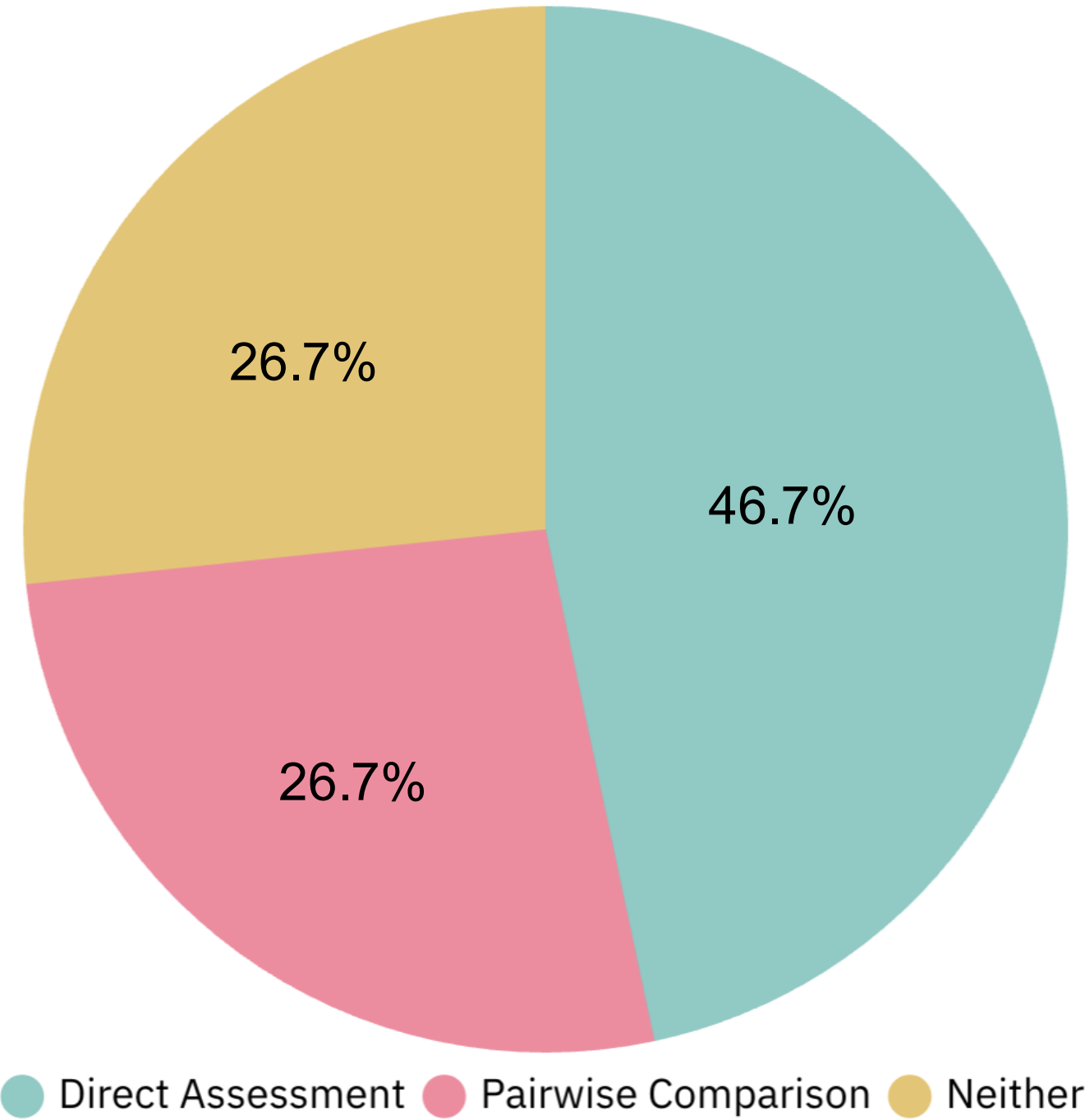**Financial inclusivity**
Within $$ budget

# Direct Assessment vs. Pairwise Comparison

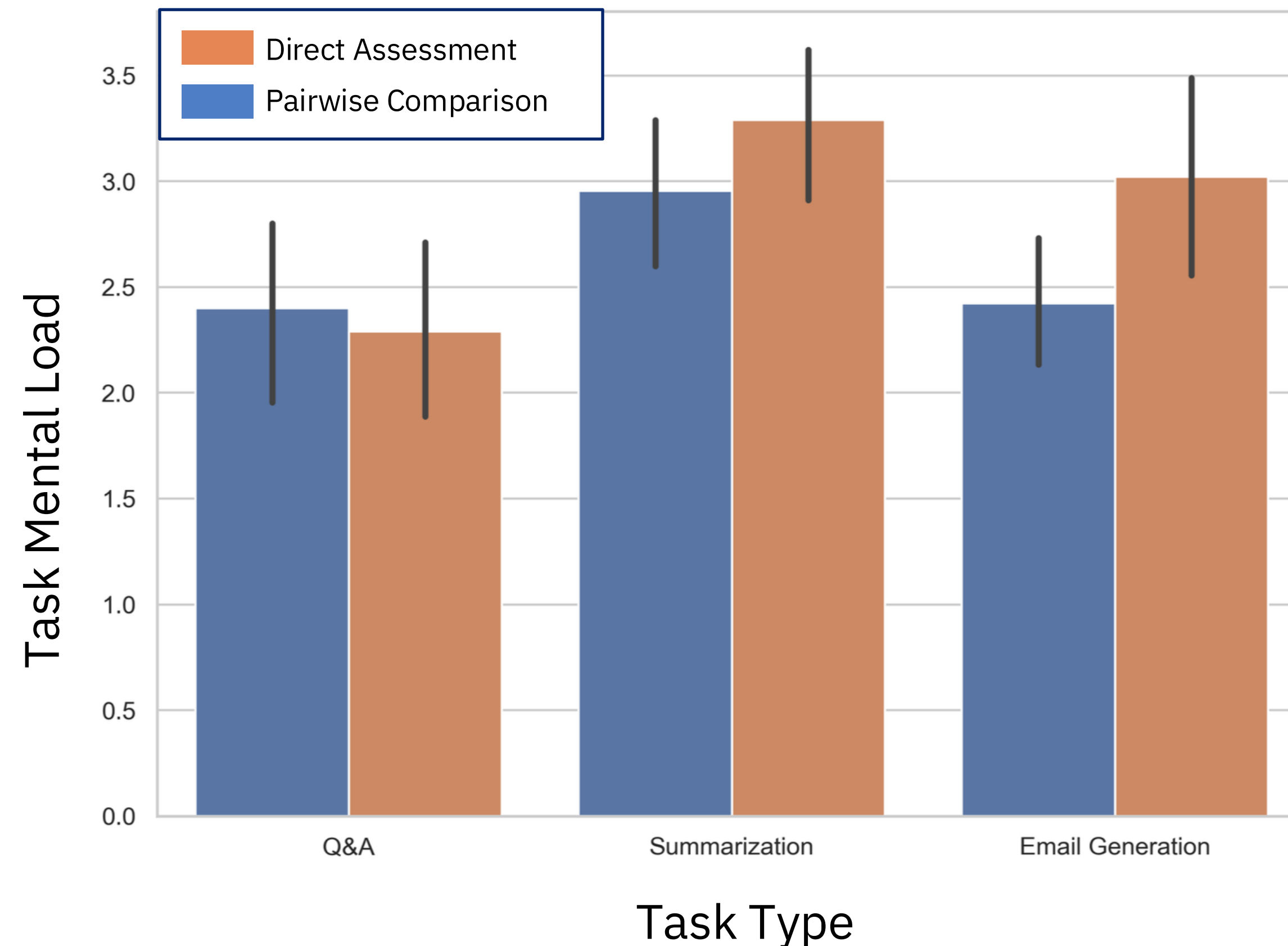Conducted more evaluations in the direct assessment compared to the pairwise comparison

Preferred direct assessment more as it offered more clarity and control

# Task Type (Q&A, Summarization, Email)

**Mental load**: Participants experienced higher mental load in the Summarization task than the Q&A task ($p<0.05$).

# Positional Bias & Explanation

The **positional bias indicator** and the **explanation** were perceived as helpful, prompting participants to revise the criteria or convinced them to adjust their expectations and leading to better human-AI alignment.

*"Positional bias inspired me to simplify the criteria and rephrase them." - P2*

*"After looking at the explanations, I realized that criteria can be improved and be made more specific. For the criteria I chose, I think I agree with the explanations provided, and it helps me understand the ranking better."- P10*

# Discussion

Involve all stakeholders in defining evaluation criteria

Mitigate over-specificity by exposing users to diverse task contexts and data samples

Help users understand how their criteria are processed and provide guidance on refinement

# Takeaway

We introduce EvalAssist (ibm.github.io/eval-assist), an open-source LLM-as-a-judge tool that assist practitioners interactively define evaluation criteria. It supports **two evaluation strategies**, flags **positional bias**, and enables **scalable evaluation**.



EvalAssist

Check out our paper!