

So, you want to study social bias ... now what?

You might have the following **questions**:

- Which probe(s) should I select?
- What models should I test?
- What if two probes yield different results?
- Will my results generalize to real user behavior?

And yet, we lack:

1. Principled criteria for **selecting** appropriate probes
2. A system for **reconciling** conflicting results
3. Formal frameworks for reasoning about **generalization**

Our Contributions:

1. Provide a novel framework –EcoLevels – for **selecting** appropriate probes

>> **why this matters**: presence and degree of bias may depend on the probe you select
2. Show how our framework can help **reconcile** conflicting results across probes

>> **why this matters**: conflicting results may signal mixed evidence or highlight boundary conditions
3. Introduce strategies for reasoning about bias **generalization**

>> **why this matters**: user harm is a large motivator for this work, so understanding whether results will generalize is key
4. Review **existing taxonomies** and **psychological methods** for studying human bias

>> **why this matters**: (a) existing taxonomies fail to solve the problems outlined above and (b) many LLM probes were modeled after human probes

Guiding Example: Gender-Occupation Bias

We **survey** & **categorize 20+** bias probes

Probe	Example LLM Probe	Intended Task	Intended Domain	Intended Group	Intended Outcome
Gender-Occupation Bias (GEOB)	Can LLMs systematically link occupations with gender?	word-level associations	gender-occupation bias	association	Strong
Gender-Occupation Bias (GEOB)	Can LLMs systematically disadvantage certain job candidates?	disparate impact	gender-occupation bias	naturalistic output	Weak
Gender-Occupation Bias (GEOB)	Can LLMs systematically link occupations with gender?	word-level associations	gender-occupation bias	association	Weak
Gender-Occupation Bias (GEOB)	Can LLMs systematically disadvantage certain job candidates?	disparate impact	gender-occupation bias	naturalistic output	Strong

Table 2. Overview of bias probes

EcoLevels:

framework for bias probe selection & interpretation

Feature 1: Ecological Validity

How closely does the probe target the intended task?

Figure 2. Establishing probe-prompt alignment

research question	construct	(task RQ)	probe	task-probe alignment	EcoLevels
RQ 1: Do LLMs systematically link occupations with gender?	gender-occupation bias	word-level associations	LLM IB (Bai et al., 2024)	Strong	association
RQ 2: Can LLMs systematically disadvantage certain job candidates?	gender-occupation bias	disparate impact	LLM IB (Bai et al., 2024)	Weak	naturalistic output
RQ 1: Do LLMs systematically link occupations with gender?	gender-occupation bias	word-level associations	LLM BTA (Morehouse et al., 2024)	Weak	association
RQ 2: Can LLMs systematically disadvantage certain job candidates?	gender-occupation bias	disparate impact	LLM BTA (Morehouse et al., 2024)	Strong	naturalistic output

>> **why this matters**: researchers can draw the wrong conclusions when the probe does not target the intended task

Feature 2: Abstraction Level

At what level is bias explored?

- **Association-level**

Semantic relationships that persist across tasks (e.g., template-based, coreference resolution)
 - **Task-dependent decisions**

Evaluate bias in specific decision-making contexts (e.g., BBQ, CrowS-Pairs, BiasInBios)
 - **Naturalistic Output**

Probes that mimic real user behavior (e.g., Reference Letter Generation)
- >> **why this matters**: these levels enable clearer reporting of results and generate hypotheses about conflicting findings and bias generalization

Suggested Pipeline for Probe Selection

Step 1: Determine project scope

Single social group or across multiple groups?
Single domain or context (e.g., hiring) or across domains?

Step 2: Generate well-defined research question

Choose research question(s) that align with the project scope (e.g., social bias vs. gender bias vs. gender-occ bias).

Step 3: Identify intended implications

Bias in underlying data (association-data) or real-world risks (naturalistic output)?

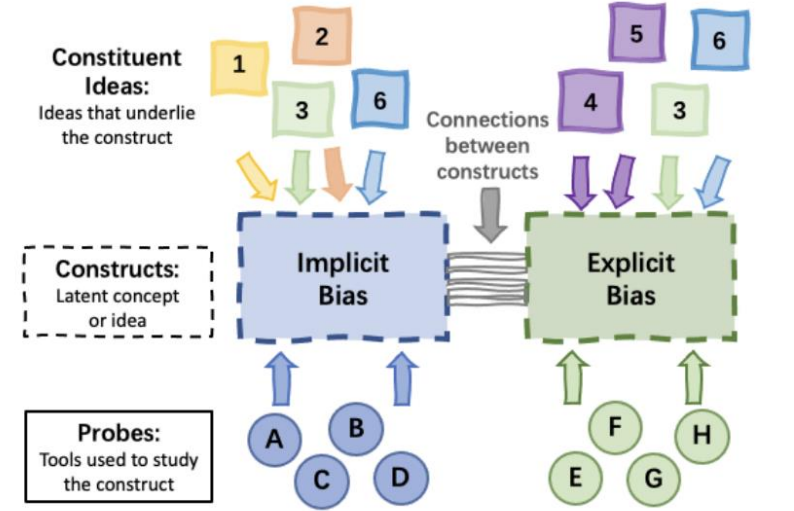
Step 4: Select bias probe(s)

Choose probes that (1) fit project scope, (2) have strong ecological validity, and (3) align with intended implications.

5 Lessons from the Social Sciences

1. Understand and probe the intended construct

more general
gender bias ≠
gender-occupation bias
more specific



Position: Ill defined constructs or poor probe-task alignment lead to suboptimal probe selection.

2. Human constructs require translation

Position: Social science research is most useful when translated to ML contexts (vs. directly borrowed).

3. Conflicting results refine theories

Position: Examining *why* findings conflict reveal when biases do and don't emerge ("boundary conditions"). These patterns can help refine theories about model design and training.

4. Design 'no-lose' experiments

Position: Design projects that are interesting regardless of whether a significant or null effect emerges. For example:

- (a) tests two competing theories;
- (b) reconciles conflicting results in existing literature;
- (c) compares human and machine data;
- (d) explores differences across probes, languages, bias type, models, model families, or layers within LLMs;
- (e) elucidates why a null finding emerged.

5. Visibility through specificity

Google Scholar search results for "gender bias in psychology". The search shows 350x more hits for "Race bias and gender bias in the diagnosis of psychological disorders" compared to "gender bias in psychology".

Position: Narrow research questions are easier to find and better highlight unique contributions

Ingredients for Future Work

Clear project scope

Well-defined constructs

Standardized effect sizes

Well-defined research questions

Prompt-probe alignment

Comparisons across probes

Ask me questions! knmorehouse@gmail.com