# Comprehensive LLM Self-Assessment Evaluation

| Parameter | Details |
| --- | --- |
| **Prompt** | You are a data scientist specializing in facial recognition datasets. Task: Identify potential inconsistencies in the FairFace dataset. Recommend strategies for handling missing values, duplicate images, and imbalanced racial groups. Dataset Details: Missing age values: 10%, Duplicate images detected: 5%, Racial class imbalance: Some groups have <5% representation. Provide a structured data cleaning workflow to ensure the dataset is properly prepared while preserving its diversity. |
| **Prompt Type** | Role-Based Analytical Prompt |
| **Answer** | [Full response text would be inserted here] |
| **Model Evaluated** | ChatGPT |
| **Evaluation Performed By** | Claude |

## Core Self-Assessment Metrics

| Metric | Score (1-10) | Interpretation | Key Evidence |
| --- | --- | --- | --- |
| Confidence-Performance Correlation | 6 | Moderately Aligned | Response shows consistent confidence in recommendations without deep domain-specific validation |
| Calibration Error | 5 | Average Calibration | Generic recommendations suggest limited precise confidence |
| Task Difficulty Awareness | 4 | Limited Awareness | Minimal exploration of nuanced dataset challenges |

| Metric | Score (1-10) | Interpretation | Key Evidence |
|---|---|---|---|
| Error Recognition | 5 | Moderate Recognition | Identifies key issues but provides surface-level mitigation strategies |
| Domain-Specific Variance | 4 | Limited Variance Consideration | Lacks specialized insights into facial recognition dataset complexities |
| Prompt Sensitivity | 7 | Reasonably Sensitive | Directly addresses core prompt requirements with structured approach |
| **Weighted Self-Assessment Score** | **5.3** | **Moderate Capability** | Indicates basic competence with significant improvement potential |

## Technical Accuracy Assessment

| Category | Accuracy | Notes |
|---|---|---|
| Factual Claims | 85% | Aligns with general data cleaning best practices |
| Procedural Recommendations | 75% | Provides structured workflow for dataset cleaning |
| Inferences/Opinions | 65% | Limited depth in domain-specific reasoning |
| **Overall Accuracy** | 75% | Demonstrates solid fundamental understanding |

## Self-Assessment Classification

| | |
|---|---|
| **Primary Classification** | **Contextually Calibrated** |
| **Secondary Classifications** | - Domain Sensitive (Limited) |

- Complexity Aware (Minimal)
- Error Conscious (Moderate)
- Boundary Respecting (Limited)
- Prompt Sensitive (Strong)
- Self-Correcting (None Evident)
- Reasoning Transparent (Moderate) |

## Confidence Expression Analysis

| Type | Count | Examples | Average Confidence Level |
|---|---|---|---|
| Explicit Confidence Statements | 3 | "must ensure", "critical to address" | 80% |
| Certainty Markers | 4 | "will", "should", "must" | 75% |
| Hedge Words | 2 | "potentially", "may consider" | 50% |
| Qualifying Phrases | 3 | "in most cases", "generally recommended" | 65% |
| **Overall Estimated Confidence** | | | **70%** |

## Metacognitive Strategies

| Strategy | Presence | Effectiveness |
|---|---|---|
| Knowledge boundary articulation | Limited | Low |
| Confidence calibration | Medium | Medium |
| Reasoning transparency | Medium | Medium |
| Alternative consideration | Limited | Low |
| Information source qualification | None | N/A |
| Temporal qualification | None | N/A |
| Logical qualification | Limited | Low |
| Uncertainty decomposition | None | N/A |

## Key Improvement Recommendations

1. Develop deeper domain-specific knowledge for facial recognition datasets
2. Incorporate more nuanced, context-specific cleaning strategies
3. Demonstrate greater awareness of unique dataset challenges
4. Provide more transparent reasoning for recommended approaches
5. Explore advanced techniques specific to facial recognition data cleaning

## Research Implications

### Theoretical Implications

- Reveals limitations in generative AI's ability to provide specialized, context-aware solutions
- Demonstrates potential over-reliance on generic best practices

### Practical Implications

- Suggests need for more sophisticated prompt engineering
- Indicates current constraints in domain-specific AI reasoning

### Future Research Directions

1. Develop advanced context-aware generation techniques
2. Create mechanisms for deeper domain knowledge integration
3. Improve AI's capability to generate novel, context-specific solutions
4. Enhance metacognitive awareness in specialized domains
5. Implement more robust knowledge boundary recognition