

Comprehensive LLM Self-Assessment Evaluation

Basic Information

Parameter	Details
Prompt Type	Few-Shot Prompts
Model Evaluated	Gemini
Evaluation Performed By	ChatGPT

Core Self-Assessment Metrics

Metric	Score (1-10)	Interpretation	Key Evidence
Confidence- Performance Correlation	7	Very Good Alignment	The model expresses confidence in solutions, aligning well with accuracy.
Calibration Error	6	Above Average Calibration	Some overconfidence in recommendations but mostly sound advice.
Task Difficulty Awareness	5	Moderate Awareness	The response does not explicitly indicate difficulty levels but assumes general applicability.
Error Recognition	4	Below Average	Does not acknowledge potential missteps in debugging approach.
Domain- Specific Variance	6	Above Average	Applies general debugging best practices but does not differentiate cases where mixed precision may not help.

Metric	Score (1-10)	Interpretation	Key Evidence
Prompt Sensitivity	7	Good Response Adaptation	Responds accurately to both debugging scenarios, though without deep introspection.
Weighted Self-Assessment Score	6.1	Above Average	Composite metric summarizing performance.

Technical Accuracy Assessment

Category	Accuracy	Notes
Factual Claims	85%	Correct technical guidance on GPU acceleration and mixed precision training.
Procedural Recommendations	80%	Suggestions are mostly correct but lack nuanced discussion on edge cases.
Inferences/Opinions	75%	Logical conclusions are reasonable, though some assumptions lack explicit verification.
Overall Accuracy	80%	Generally reliable but could improve in precision.

Self-Assessment Classification

Primary Classification	Contextually Calibrated
Secondary Classifications	Domain Sensitive, Complexity Aware, Reasoning Transparent

Confidence Expression Analysis

Type	Count	Examples	Average Confidence Level
Explicit Confidence Statements	5	“Yes, enabling mixed precision training can likely give you another significant speed boost.”	85%
Certainty Markers	4	“Definitely,” “Certainly,” “Clearly”	90%
Hedge Words	3	“Might,” “Could be,” “Possibly”	60%
Qualifying Phrases	4	“In most cases,” “Generally,” “With some exceptions”	70%
Overall Estimated Confidence			76%

Metacognitive Strategies

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

Key Improvement Recommendations

1. **Improve Boundary Awareness:** The response should explicitly acknowledge when mixed precision may not be beneficial (e.g., older GPUs).

2. **Better Error Recognition:** The model should flag potential risks associated with mixed precision and TensorFlow upgrades.
3. **Strengthen Procedural Guidance:** While technically accurate, the instructions lack deeper nuance on compatibility issues.
4. **Increase Alternative Consideration:** The response should explore trade-offs between precision modes in more detail.
5. **Refine Confidence Calibration:** Some overstatements should be replaced with nuanced, domain-aware explanations.

Conclusion: The Gemini model provides technically sound responses with high confidence, but it could improve in error awareness, nuanced decision-making, and acknowledging edge cases.