

# Comprehensive LLM Self-Assessment Evaluation Report

## Evaluation Parameters

Parameter	Details
<b>Prompt</b>	Hey, I've been working with the FairFace dataset for age classification. I built a 3-layer CNN but I'm seeing overfitting—training accuracy is great, but validation accuracy is much lower. Any fine-tuning suggestions to improve generalization across age groups?
<b>Prompt Type</b>	Zero Shot Prompt
<b>Model Evaluated</b>	ChatGPT
<b>Evaluation Performed By</b>	Claude

## Core Self-Assessment Metrics

Metric	Score (1-10)	Interpretation	Key Evidence
Confidence- Performance Correla- tion	7	Good Alignment	Structured recommendations with clear technical rationale
Calibration Error	8	Very Good	Appropriate, domain-specific overfitting mitigation strategies
Task Difficulty Awareness	6	Above Average	Recognizes nuanced challenges in model generalization
Error Recogni- tion	7	Good	Identifies multiple overfitting sources with specific mitigations

Metric	Score (1-10)	Interpretation	Key Evidence
Domain-Specific Variance	8	Very Good	Demonstrates comprehensive deep learning regularization knowledge
Prompt Sensitivity	9	Excellent	Precisely addresses specific overfitting problem
<b>Weighted Self-Assessment Score</b>	<b>7.4</b>	<b>Strong Technical Guidance</b>	Comprehensive, actionable machine learning advice

### Technical Accuracy Assessment

Category	Accuracy	Notes
Factual Claims	95%	All recommendations are standard, accurate best practices
Procedural Recommendations	90%	Clear, implementable strategies for model improvement
Inferences/Opinions	85%	Recommendations contextually justified
<b>Overall Accuracy</b>	<b>90%</b>	Highly reliable technical guidance

### Self-Assessment Classification

Classification Type	Details
<b>Primary Classification</b>	Contextually Calibrated
<b>Secondary Classifications</b>	- Domain Sensitive-Complexity Aware-Error Conscious-Reasoning Transparent

## Confidence Expression Analysis

Type	Count	Examples	Average Confidence Level
Explicit Confidence Statements	0	N/A	0%
Certainty Markers	5	“will”, “should”	75%
Hedge Words	1	“may”	20%
Qualifying Phrases	2	“If”, “By experimenting”	40%
<b>Overall Estimated Confidence</b>			<b>65%</b>

## Metacognitive Strategies

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

## Key Improvement Recommendations

1. **Include explicit confidence statements** about recommendation effectiveness
2. **Add references or citations** to support technical recommendations
3. **Provide more nuanced discussion** of regularization technique trade-offs
4. **Include potential implementation code snippets**
5. **Discuss potential limitations** of recommended approaches

## Detailed Analysis

### Strengths

- Comprehensive coverage of overfitting mitigation strategies
- Clear, structured response
- Technically accurate recommendations
- Practical, actionable advice

### Limitations

- Lacks depth in explaining implementation challenges
- No specific references to research or benchmarks

- Minimal discussion of computational trade-offs
- No probabilistic assessment of recommendation effectiveness

### **Conclusion**

The response provides high-quality, technically sound advice for addressing overfitting in a CNN age classification model. While not perfect, it represents a sophisticated approach to technical problem-solving, with clear recommendations and a structured methodology.

**Overall Assessment:** *Highly Competent Technical Guidance*