

Category	Metric	Definition / Measurement Approach	Purpose
Automated	Perplexity	Measures log-likelihood on a held-out dataset.	Evaluates model fluency and text quality.
Automated	BLEU (Bilingual Evaluation Understudy)	Compares generated text against multiple reference texts.	Automated NLP metric (n-gram precision).
Automated	ROUGE (Recall-Oriented Understudy)	Evaluates lexical similarity and longest common subsequence (LCS).	Tests accuracy and content overlap.
Automated	Semantic Similarity	Measures cosine similarity using sentence embeddings.	Evaluates contextual relevance and semantic alignment.
Cultural	CSI Score (Cultural Sensitivity Index)	Indirectly weighted by human or model-assisted feedback.	Calculated using human or model-assisted feedback.
Cultural	Stereotype Detection	Identifies potential bias classifiers, keyword detection, and context.	Ensures inclusivity and avoids cultural generalizations.
Cultural	Regional Appropriateness	Evaluates human evaluation or rule-based classifier.	Tests localization and contextual fit.
Biblical Integration	Contextualization Accuracy	Measures expert or rule-based evaluation comparison.	Ensures meaningful and accurate cross-religious analogies.
Biblical Integration	Relevance	Evaluates human scoring on relevance and interpretation.	Tests cross-cultural interpretive value.
Cost	Tokens Used	Number of tokens retrieved from API logs or model output.	Assesses resource usage.
Cost	API Costs	Monetary cost based on token usage and provider pricing.	Enables cost-performance comparison.
Cost	Inference Time	Average time measured per API call or benchmark script.	Evaluates computational efficiency.