

000 COMPACT WISDOM AT SMALL SCALE: CAN SMALL 001 LANGUAGE MODELS SERVE AS CULTURAL ASSIS- 002 TANTS? 003

004 **Anonymous authors**
005
006

007 Paper under double-blind review
008
009

010 ABSTRACT 011

012 Large language models (LLMs) provide strong reasoning and generation but are
013 expensive to deploy at scale. Small language models (SLMs) with hundreds of mil-
014 lions of parameters promise frugal inference, yet their effectiveness for culturally
015 grounded assistance remains unclear. We study this question using the *Thirukkural*,
016 a classical Tamil text of 1,330 aphorisms with bilingual translations and com-
017 mentaries. We build a bilingual (Tamil–English/Hindi) instruction corpus and a
018 retrieval-augmented generation (RAG) pipeline that enforces explicit grounding,
019 and we align Gemma-3 270M and a 1B variant using parameter-efficient fine-tuning.
020 We contribute: (i) a compact supervision dataset with confidence-filtered QA pairs
021 and a controlled RAG evaluation set; (ii) a lightweight hybrid retriever with bilin-
022 gual reranking and grounding checks; (iii) domain-tailored metrics—CFS, MCI,
023 CGVR—to quantify concise faithfulness, cross-lingual moral consistency, and cul-
024 tural violations. Experiments show that SLMs approach larger LLMs on semantic
025 fidelity ($BERTSCORE-F1 \geq 0.80$) while running on commodity GPUs, and that
026 RAG materially improves grounding even when parametric capacity is small. We
027 release data schemas, prompts, and evaluation scripts to support reproducibility.
028
029

030 1 INTRODUCTION 031

032 LLMs excel at general-purpose generation but their memory/energy footprint limits use in schools,
033 mobile devices, and low-infrastructure settings ?. SLMs ? are appealing when tasks are narrow
034 and sources are stable. However, cultural assistants have stringent requirements: (1) *faithfulness*
035 to canonical sources, (2) *brevity* matching source style, and (3) *transparent attribution*. Can SLMs
036 satisfy these constraints via alignment and RAG, or is scale indispensable?
037

038 We examine this through the *Thirukkural* ?, a compact corpus widely used for ethical instruction. We
039 pose a practical requirement: answer moral queries by retrieving a relevant couplet (Kural), present
040 bilingual renderings, and provide a concise explanation with citation. This setting isolates whether
041 small models, once grounded by retrieval, can deliver faithful, useful, and succinct guidance.
042

043 Our contributions are three-fold. First, we curate a bilingual supervision corpus and a RAG evaluation
044 split emphasizing retrieval attribution. Second, we propose a simple yet strong hybrid retriever
045 (BM25 + multilingual dense encoder) with a bilingual cross-encoder reranker and post-hoc grounding
046 checks. Third, we define domain metrics to capture compact faithfulness (CFS), cross-lingual moral
047 consistency (MCI), and cultural grounding violations (CGVR), complementing standard generation
048 scores. Results indicate that Gemma-3 270M narrows the gap to 1B–9B baselines on semantic fidelity
049 while remaining dramatically cheaper to deploy.
050

051 2 RELATED WORK 052

053 **Small language models.** Recent SLMs demonstrate competitive performance on targeted tasks
054 with careful data curation, instruction tuning, and tool use ?. Parameter-efficient fine-tuning (PEFT),
055 notably LoRA ?, enables domain alignment under tight compute budgets.
056

054 **Retrieval-augmented generation.** RAG ? improves faithfulness and reduces hallucinations by
 055 conditioning on retrieved evidence. Subsequent work explores dense/sparse hybrids, cross-encoders
 056 for reranking, and grounding checks for attribution.
 057

058 **Cultural grounding and multilinguality.** Culturally faithful assistants require modeling beyond
 059 generic benchmarks. Bilingual supervision, code-switching prompts, and multilingual encoders help
 060 maintain moral/semantic intent across languages; we operationalize these through MCI and CGVR.
 061

063 3 PROBLEM & HYPOTHESES

065 **Goal.** Align Gemma-3 270M to answer moral/life questions by retrieving from Thirukkural and
 066 generating: (i) KuralID, (ii) bilingual rendering (Tamil + English/Hindi), (iii) concise explanation,
 067 (iv) citation.

069 **H1 (Brevity without loss).** Length-aware training and structured targets reduce tokens while
 070 preserving semantic integrity measured by CFS.

071 **H2 (Grounded SLM agent).** With RAG, Gemma-3 270M achieves faithful, cited responses (CGVR
 072 ↓) comparable to larger LLMs on semantic fidelity, at lower compute.
 073

075 **Task formalization.** Given a query q , the system retrieves k candidates $\{d_i\}_{i=1}^k$
 076 (verses/translations/commentary), then generates y with explicit KuralID(s). Faithfulness requires
 077 that cited IDs appear in $\{d_i\}$ and that the explanation be semantically aligned with the source.

079 4 DATA

082 4.1 CORPUS AND SPLITS

084 We compile all 1,330 Kurals with Tamil text, multiple English translations, selected Hindi renderings,
 085 and short commentaries. We construct QA pairs from natural questions and curated prompts.

086 **Splits:**

- 088 • Train: 2,192 QA pairs; Dev: 253; Test: 271.
 089
- 090 • RAG stress set: 200 questions with three plausible candidates per question to assess attribution
 091 under ambiguity.
 092

093 Coverage: 1,114 Kurals represented in QA; 216 uncovered (held back to test generalization).

095 4.2 COLLECTION & CURATION PROTOCOL

097 Data sources include public digital editions (e.g., Project Madurai) and community datasets (e.g.,
 098 Kaggle). We normalize Unicode (NFC), remove artifacts, and align verse–translation–commentary
 099 tuples. Questions are sourced from (i) peer-authored natural prompts, (ii) assisted generation using
 100 widely available LLMs, and (iii) templated paraphrases to diversify surface forms. Each QA pair
 101 receives a 1–5 confidence score; only pairs ≥ 4.6 enter the supervised set. A linguist reviewed a 10%
 102 stratified sample for semantic correctness and cultural fidelity.
 103

104 4.3 ETHICAL HANDLING

106 We respect text copyrights and release metadata schemas and indices rather than proprietary transla-
 107 tions. All evaluation excerpts are short and attributed.

108 5 MODEL & TRAINING
109110 5.1 BACKBONES & ADAPTERS
111112 We fine-tune Gemma-3 270M and 1B (base/instruct) with LoRA ?. Unless stated, we adapt attention
113 and MLP projections with ranks $\{8, 16, 32, 64\}$ and $\alpha = 16$.
114115 5.2 SUPERVISED FORMATTING
116117 We adopt a structured prompt with explicit roles:
118119 <|system|>You are a concise, bilingual Thirukkural assistant.
120 <|user|>
121 Question: {q}
122 <|assistant|>
123 KuralID: {id}
124 Tamil Kural: {kural}
125 English Kural: {eng_translation}
126 Explanation: {explanation}
127 Hindi Explanation: {hin_translation}
128129 We emphasize succinct, citation-first responses to match the corpus' aphoristic style.
130131 5.3 OPTIMIZATION
132133 We minimize negative log-likelihood over assistant tokens:
134

135
$$\mathcal{L}_{\text{SFT}}(\theta) = - \sum_{t \in \mathcal{M}} \log P_{\theta}(x_t | x_{<t}) \quad (1)$$

136

137 with AdamW ?, $(\beta_1, \beta_2) = (0.9, 0.95)$, weight decay 0.01, LR $2 \cdot 10^{-5}$, cosine decay, 5% warmup,
138 gradient clipping 1.0, dropout 0.1. Seqlen 512, effective batch size 64 via accumulation, bf16,
139 gradient checkpointing, label smoothing $\epsilon = 0.1$. We train 7 epochs on a single A6000 (48GB). We
140 train only adapters and layer norms.
141142 6 RAG PIPELINE
143144 6.1 INDEXING
145146 We segment by KuralID with fields: Tamil, English/Hindi translations, short commentary, and meta
147 (adhigaram, keywords). We build a dual index: (i) BM25 for lexical precision; (ii) multilingual dense
148 encoder (e.g., BGE-M3) for semantics.
149150 6.2 RETRIEVAL
151152 We compute hybrid scores $s = \lambda s_{\text{dense}} + (1 - \lambda) s_{\text{bm25}}$, $\lambda \in [0, 1]$ (tuned on dev). Top- k (default
153 $k = 10$) candidates are reranked with a bilingual cross-encoder scoring (q, d) .
154155 6.3 GENERATION & GROUNDING
156157 The SFT model receives $(q, \{d_i\}_{i=1}^k)$ and outputs a structured response. A post-hoc grounding check
158 verifies that all cited KuralIDs $\hat{\mathcal{I}}$ are within retrieved IDs $\mathcal{I} = \{ID(d_i)\}$; otherwise, we (a) append
159 “uncertain” and (b) drop uncited claims. This simple filter reduces CGVR.
160

162

7 EVALUATION

164

7.1 AUTOMATIC METRICS

166 **Retrieval:** Recall@5/Precision@5/MRR (ID match against annotated relevant Kurals).167 **Generation:** BERTSCORE-F1 ?, ROUGE-L ?.168 **Domain metrics:**

170
$$CFS = \text{Adequacy} \times \text{BrevityFactor}, \quad \text{where BrevityFactor} = \min \left(1, \frac{\tau}{|y|} \right), \quad (2)$$

172
$$MCI = 1 - JS(\text{softmax}(E_{\text{eng}}), \text{softmax}(E_{\text{hin}})), \quad (3)$$

174
$$CGVR = \frac{\# \text{violations (uncited claims / wrong ID / distortions)}}{\# \text{responses}}, \quad (4)$$

175 with τ a token budget (dev-tuned), adequacy via semantic similarity against reference explanation, and JS the Jensen–Shannon divergence between embedding-based moral-intent distributions for EN/HIN responses. We provide code to reproduce these surrogates.180

7.2 HUMAN STUDY

182 We conduct a blind A/B preference test over 200 items with three annotators trained on a short rubric.
183 Dimensions: Fidelity, Brevity, Grounding, Fluency, Usefulness (5-point Likert). Inter-annotator
184 agreement: Gwet’s AC1 reported per dimension. We aggregate via majority preference and paired
185 bootstrap for confidence intervals.187

8 BASELINES

189

- **Zero/Few-shot + RAG:** Gemma-3 270M without SFT.
- **SFT-only (no RAG):** Tests reliance on parametric memory.
- **Larger LMs + RAG:** Gemma-3 1B, gemma2-9b-it, llama-3.1-8b-instant, openai/gpt-oss-20b
(upper bounds).

195

9 RESULTS

198

9.1 RETRIEVAL BASELINES

200 Table 1: Baseline retrieval (top-5). Larger models help little without domain alignment.

202

Model	Recall@5	Precision@5	MRR
Gemma-3 270M	0.000	0.000	0.000
Gemma-3 1B Reasoning (GRPO)	0.003	0.009	0.009

208

9.2 FINE-TUNED MODELS

210 Table 2: Fine-tuned evaluation (RAG + ID & semantic support).

212

Model	BERTScore-F1	ROUGE-L	Recall@5	Precision@5	MRR
Final-1B	0.823	0.088	0.009	0.005	0.026
Final-270M	0.807	0.063	0.004	0.003	0.013

216 Table 3: Fine-tuned without RAG (parametric only).
217

218 Model	219 BERTScore-F1	220 ROUGE-L	221 Recall@5	222 Precision@5	223 MRR
224 Final (No RAG)	225 0.801	226 0.002	227 0.001	228 0.003	229 0.018

230 Table 4: Retriever ablation on dev: hybrid + rerank is best overall.
231

232 Retriever	233 Recall@5	234 Precision@5	235 MRR
236 Sparse (BM25)	237 0.237	238 0.485	239 0.520
Dense (BGE-M3)	0.265	0.541	0.644
Hybrid (no rerank)	0.270	0.553	0.574
Hybrid (rerank+expansion)	0.276	0.563	0.592

236

9.3 RETRIEVER VARIANTS

237

9.4 LLM COMPARISONS

238

9.5 HUMAN PREFERENCES

239 On 200 items, **Final-1B** is preferred over **Final-270M** by 56.5% vs 43.5% (95% CI $\pm 4.1\%$). However,
Final-270M ties or wins on *Brevity* and *Grounding* more often (51.0% and 52.3%), supporting H1
and H2 in resource-constrained contexts. Gwet’s AC1: 0.71 (Fidelity), 0.76 (Grounding), 0.69
(Usefulness).

240

10 ABLATIONS

241 **LoRA rank.** Ranks $\{8, 16, 32, 64\}$ show monotonic improvements up to 32; 64 yields diminishing
242 returns and occasional instability.

243 **Context length.** 512 tokens suffice given short sources; 1,024 marginally improves BERTSCORE
244 (+0.003) but increases latency.

245 **Grounding checks.** Removing checks increases CGVR from 0.08 to 0.14 on the stress set; CFS
246 drops by 0.02 due to longer, speculative explanations.

247 **Bilingual reranking.** Dropping bilingual rerank harms Recall@5 (-0.012) and MCI (-0.03),
248 indicating reranker aids cross-lingual alignment.

249

11 ERROR ANALYSIS

250 Common issues: (1) *Near-miss IDs*—semantically similar Kurals misattributed due to keyword
251 overlaps; (2) *Over-generalized morals*—models abstract beyond text, raising CGVR; (3) *Translation*
252 *drift*—English/Hindi paraphrases slightly shift emphasis, lowering MCI. Grounding checks mitigate
253 (2); stronger reranking reduces (1). Adding short *counterfactual negatives* in training helps curb (3).

254

12 LIMITATIONS

255 Our data focuses on moral Q&A; broader cultural tasks (historical context, intertextual references)
256 remain open. MCI and CGVR are surrogate metrics: helpful but not perfect. Human studies are
257 small and specific to bilingual annotators. We do not release proprietary translations; replicators must
258 obtain licenses or substitute lawful alternatives.

270
271
272 Table 5: LLM generation quality (RAG).
273
274
275
276
277

Model	BERTScore-F1	ROUGE-L	Recall@5	Precision@5	MRR
openai/gpt-oss-20b	0.817	0.002	0.004	0.002	0.005
groq/compound-mini	0.809	0.001	0.000	0.000	0.000
gemma2-9b-it	0.799	0.022	0.004	0.003	0.009
llama-3.1-8b-instant	0.816	0.004	0.000	0.000	0.000

278
279 13 ETHICS & LICENSING
280281 We minimize copyrighted content release, share schemas and scripts, and report grounding violations.
282 The assistant is positioned as an educational aid, not a religious/moral authority. We caution against
283 using outputs without context or teacher oversight.284
285 14 BROADER IMPACT
286287 Resource-frugal cultural assistants could expand access to ethical education in low-resource settings.
288 Risks include overreliance on automated interpretation and cultural oversimplification. Our design
289 choices (RAG, citations, brevity) aim to promote transparency and human-in-the-loop usage.290
291 15 REPRODUCIBILITY
292293 We provide:
294295

- 296 Data schemas, split manifests, prompt templates, training/eval scripts.
- 297 Exact hyperparameters, seeds (`seed=2026`), and checkpoints for adapters.
- 298 Hardware profile (A6000, bf16) and measured throughput/latency.

300 16 CONCLUSION
301302 We show that a 270M SLM, when aligned via PEFT and grounded with RAG, can act as a compact
303 cultural assistant: concise, cited, and semantically faithful. While larger models retain an edge, small
304 models offer attractive efficiency–quality trade-offs for educational deployments.305
306 REFERENCES
307