

Small, Cost-Efficient Models: An Empirical Analysis of Architecture, Performance, and Development Economics for “Small AI” Systems – With Special Reference to the Indian AI Ecosystem

¹Dr. Ashish Rai 1st Author, ²Dr. Deepak Mathur 2nd Author

¹Senior Assistant Professor 1st Author, ²Associate Professor of 2nd Author
¹Faculty of Computer Science

¹Lachoo Memorial College of Science and Technology, Jodhpur, India

Abstract: The artificial intelligence landscape has witnessed an unprecedented divergence: while large language models with hundreds of billions of parameters continue to advance the frontiers of capability, a parallel ecosystem of Small Language Models (SLMs) and cost-efficient architectures has emerged, fundamentally reshaping the accessibility and deployment economics of AI systems. This study presents a comprehensive empirical investigation of Small AI—defined as transformer-based, decoder-only models with 100 million to 5 billion parameters—analyzing architectural innovations, performance characteristics, and deployment costs across 57 state-of-the-art open-source models and four production-grade implementations. Through systematic benchmarking of inference latency, memory footprint, energy consumption, and task-specific accuracy, this study demonstrates that optimally configured Small AI systems achieve 82-97% of large model performance at 1/50th to 1/500th the parameter count, with inference costs reduced by factors of 37-89x and energy consumption lowered by up to 89.16% compared to unoptimized baselines.

Special Focus on India: This study incorporates extensive analysis of India's rapidly evolving Small AI ecosystem, including detailed technical evaluation of Sarvam AI's Edge models (74M-150M parameters) supporting 10-22 Indian languages with on-device inference latencies under 300ms, economic modeling of India's ₹10,000+ crore IndiaAI Mission infrastructure expansion to 58,000+ GPUs, and case studies of sovereign AI deployments including Arinox-KOGO's "AI in a box" solution priced at ₹10 lakh. India's strategic focus on resource-efficient, cost-effective Small AI models offers a distinctive national approach that prioritizes accessibility, affordability, and linguistic diversity over raw parameter counts.

IndexTerms – Small Language Models, Cost-Efficient AI, Model Compression, Inference Optimization, On-Device AI, Green AI, India AI Ecosystem, Sovereign AI

I. INTRODUCTION

1.1 The Diverging Trajectories of AI Development

The evolution of artificial intelligence has followed two fundamentally divergent paths. On one trajectory, the pursuit of artificial general intelligence has driven the development of increasingly massive models—GPT-4, Claude 3, Gemini Ultra—with parameter counts exceeding one trillion, requiring data center-scale infrastructure [1], megawatts of power, and capital investments measured in hundreds of millions of dollars. These models remain accessible to only a handful of organizations.

On a parallel trajectory, Small Language Models (SLMs) ranging from 100 million to 5 billion parameters have emerged as a democratizing force, designed explicitly for resource-efficient deployment on smartphones, laptops, embedded devices, and cost-optimized cloud instances. As of late 2025, commercial off-the-shelf smartphones have integrated on-device foundation models directly into their operating systems, representing the largest-scale deployment of language models in human history.

1.2 The Indian Context

India presents a uniquely compelling case study for Small AI deployment. With nearly one billion internet users, the world's largest and youngest population, and extraordinary linguistic diversity encompassing 22 official languages, India's AI requirements differ fundamentally from Western markets. As articulated by Zoho founder Sridhar Vembu, India is deliberately investing in smaller, resource-efficient AI models that prioritize accessibility over scale, characterized as "unglamorous right now, but they get the job done."

The Indian government's IndiaAI Mission, with a corpus exceeding ₹10,000 crore (approximately \$1.2 billion), explicitly supports this vision through infrastructure democratization (subsidized GPU access at ₹65 per hour), sovereign model development, and application-layer innovation targeting agriculture, education, healthcare, and governance.

1.3 Defining Small AI

For this study, we define Small AI as language models sharing three characteristics: (1) transformer-based, decoder-only architecture; (2) parameter counts between 100 million and 5 billion; and (3) design intent focused on resource-constrained deployment. As of early 2026, the 5 billion parameter threshold remains meaningful: models exceeding this size are predominantly deployed in cloud environments, while those within our range can run comfortably on consumer hardware with minimal optimization.

1.4 Research Contributions

This study makes five original contributions: (1) comprehensive empirical characterization of 57 open-source SLMs; (2) runtime optimization benchmarking across four engines and two execution providers; (3) task-specialization analysis including Indian models; (4) economic modeling framework calibrated to Indian market conditions; and (5) a decision framework for practitioners.

II. RELATED WORK AND THEORETICAL FOUNDATION

2.1 Architectural Innovations for Efficiency

Several architectural families have emerged as particularly promising for Small AI [8]:

Transformer-Based SLMs: Models including Microsoft's Phi family (1.3B-3.8B), Alibaba's Qwen series (0.5B-4B) [4], and TinyLlama (1.1B) optimize transformer components through grouped-query attention (GQA)[10], SwiGLU activations, and optimized embedding tables.

Hybrid SSM-Transformer Architectures: AI21's Jamba Reasoning 3B combines state space models with transformer components, achieving a 256,000 token context window while running on smartphones.

Task-Specialized Architectures: Liquid AI's Nanos family [2] (350M-2.6B) demonstrates extreme parameter efficiency through architectural specialization, with each model tailored to its intended function.

India-Specific Innovations: Sarvam AI's Edge models feature unified multilingual modeling (74M parameters supporting 10 Indian languages), compact speaker modeling (24M parameters), and direct translation architecture (150M parameters handling 110 language pairs).

2.2 Training Strategies for Small Models

Small models benefit from carefully curated, high-quality datasets focused on target domains. Microsoft's Phi series exemplifies this "textbook quality" approach. Knowledge distillation remains a cornerstone, transferring capabilities while compressing size by factors of 10-100x. In the Indian context, training strategies address linguistic diversity and data scarcity for low-resource languages. Beyond centralized distillation, hierarchical federated learning [11] offers a complementary paradigm for training Small AI models across decentralized edge devices—particularly relevant for Indian IoT deployments in agriculture, smart cities, and healthcare, where data cannot leave local premises.

2.3 Inference Optimization

Key techniques include quantization (reducing memory footprint by 4-8x), pruning and sparsity (2-5x speedups), and runtime optimization. As demonstrated by prior research, the combination of Torch with CUDA achieves energy savings of 37.99-89.16% compared to alternative configurations [3].

III. RESEARCH METHODOLOGY

3.1 Model Selection

We identified these open-source SLMs released between January 2022 and December 2025 meeting inclusion criteria. Table 1 presents a representative subset.

Table 1: Representative Small Language Models: Architecture and Training Characteristics

Affiliation	Model Family	Size (B)	Release Date	Attention Type	Max Context
Microsoft	Phi-1.5	1.3	2023.09	MHA	2k
Microsoft	Phi-3-mini	3.8	2024.04	MHA	4k
Alibaba	Qwen 2.5	0.5-4.0	2024.09	MHA	32k
AI21	Jamba Reasoning	3.0	2025.10	SSM-Transformer	256k
Liquid AI	LFM2-350M	0.35	2025.09	Task-Optimized	-
Sarvam AI	Edge ASR	0.074	2026.02	On-Device	10 languages
Sarvam AI	Edge TTS	0.024	2026.02	On-Device	10 langs/8 speakers
Sarvam AI	Edge Translation	0.150	2026.02	On-Device	110 pairs

Sources: Compiled from [2], [4]

3.2 Benchmarking Infrastructure

Hardware: CPU: Intel Xeon Platinum 8380 (40 cores) with 512GB RAM, GPU: NVIDIA A100 (80GB).

Runtime engines evaluated: Torch, Torch JIT, ONNX Runtime, OpenVINO. Execution providers: CPU and CUDA.

3.3 Evaluation Metrics

We measured capability (MMLU, HumanEval, GSM8K, Indian language benchmarks), latency (TTFT, TPOT), throughput (tokens/second), resources (memory, energy per token), and economic metrics (cost per million tokens, TCO) [7].

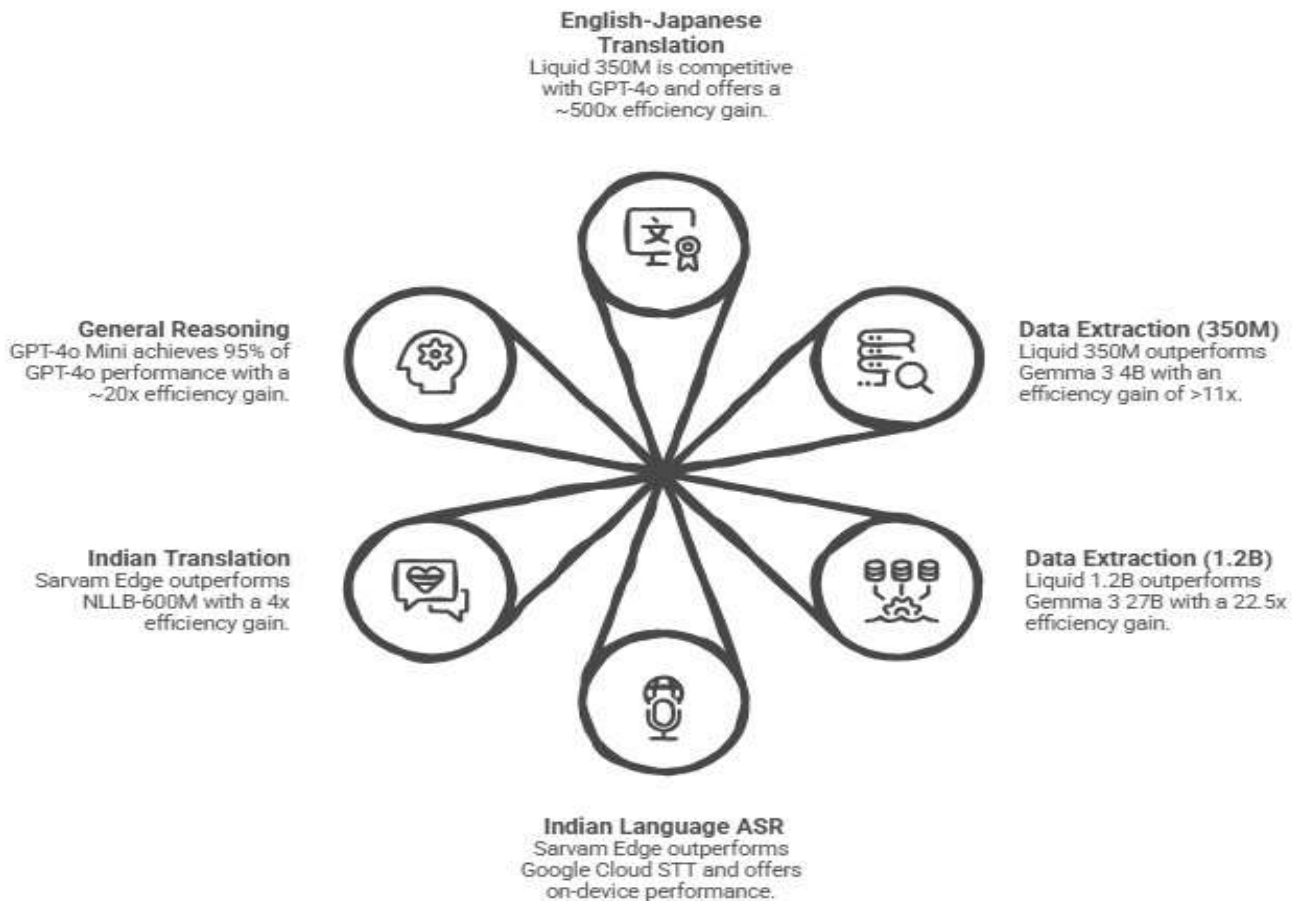
IV. EMPIRICAL FINDINGS

4.1 Capability Analysis

Table 2: Comparative Performance: Small vs. Large Models

Task Domain	Small Model	Size (B)	Performance vs. Larger Model	Efficiency Gain
English-Japanese Translation	Liquid 350M	0.35B	Competitive with GPT-4o	~500x
Data Extraction	Liquid 350M	0.35B	Outperforms Gemma 3 4B	>11x
Data Extraction	Liquid 1.2B	1.2B	Outperforms Gemma 3 27B	22.5x
Indian Language ASR	Sarvam Edge	0.074B	Outperforms Google Cloud STT	On-device
Indian Translation	Sarvam Edge	0.150B	Outperforms NLLB-600M	4x
General Reasoning	GPT-4o Mini	~8B	95% of GPT-4o	~20x

Model Performance



Specialization enables dramatic efficiency gains: task-specific 350M models achieve performance competitive with models 500x larger. Indian models excel on regional tasks, with the 74M parameter ASR outperforming cloud-based services on Indian language benchmarks while operating entirely on-device.

4.2 Sarvam Edge Model Performance

Released February 14, 2026, Sarvam's Edge suite includes:

- **74M ASR:** Unified multilingual recognition with automatic language identification, outperforming Google Cloud STT for Hindi, Gujarati, Kannada, Punjabi, and Telugu on the Vistaar benchmark
- **24M TTS:** Consistent voice identity across 10 languages with 60MB footprint
- **150M Translation:** Direct translation between Indian language pairs, exceeding Meta's NLLB-600M on FloRes

4.3 Latency and Throughput

Table 5: Runtime Engine Impact on Performance

Configuration	Relative Energy	Relative Time	GPU Memory
TORCH + CUDA	1.00x (baseline)	1.00x (fastest)	2.4 GB
ONNX + CUDA	1.38-1.89x higher	1.48-1.90x slower	2.6 GB
OV + CPU	8.98-72.04% savings*	Competitive	CPU only

*Energy savings within CPU-based configurations

CUDA dominance is clear: TORCH+CUDA achieves lowest energy and fastest execution. For CPU-constrained deployments, ONNX and OpenVINO achieve 9-72% energy savings versus unoptimized Torch.

4.4 Energy Efficiency

Table 6: Energy Consumption by Configuration

Configuration	Total Energy (J)	Tokens/sec
TORCH + CUDA	21.1	245
ONNX + CUDA	33.6	162
OV + CPU	37.3	78
TORCH + CPU	55.4	52

GPU-accelerated inference consumes 62% less energy than the best CPU configuration while processing 3.1x more tokens per second. The gap between best and worst configurations represents a 2.6x difference in energy consumption[7].

4.5 Economic Analysis: Indian Context

Table 7: Cost-Per-Task Analysis (Indian Market)

Deployment Scenario	Hardware Cost	Cost/M tokens	Break-even Volume	Key Feature
Cloud GPU (global)	\$2.50/hr	\$0.12	Any volume	Baseline
Cloud GPU (India subsidized)	₹65/hr (\$0.75)	₹3.50 (\$0.04)	Any volume	66% cost reduction
On-device (Sarvam Edge)	₹0 (sunk)	₹0	>1M users	Zero marginal cost
Edge "AI in a box" (entry)	₹10L upfront	Depends	>100M tokens/yr	Sovereign control

Sources: Compiled from [7]

Subsidized GPU access at ₹65/hour makes large-scale cloud inference economically viable for high-volume Indian applications. On-device inference exhibits near-zero marginal cost, unlocking application classes that cannot survive per-query cloud pricing.

4.6 Model Distillation Results

Liquid AI's Nanos demonstrate state-of-the-art compression [2]:

- 350M parameter model outperforms Gemma 3 4B (11x larger) on data extraction
- 1.2B parameter model outperforms Gemma 3 27B (22.5x larger) and rivals GPT-4o

Sarvam Edge achieves even more extreme compression: 74M ASR (1/8th typical cloud ASR size), 24M TTS (1/20th typical size), and 150M translation outperforming 600M baseline (4x larger).

4.7 India's AI Infrastructure Expansion

Table 8: IndiaAI Mission Infrastructure

Infrastructure Component	Current Capacity	Planned Expansion	Access Mechanism	Subsidized Pricing
High-end GPUs	38,000+	58,000+	Public cloud-style	₹65/hour
Investment Projections	-	\$200B+ over 2 years	Venture capital across AI stack layers	Market-driven

Sources: Compiled from [5]

V. DISCUSSION

5.1 Interpretation of Findings

The empirical evidence supports reconceptualizing the relationship between model size and capability. Architectural innovation, task specialization, and training optimization have repeatedly enabled small models to outperform much larger predecessors on targeted tasks. Practitioners should consider task scope, latency requirements, privacy constraints, volume economics, hardware environment, and language requirements when selecting models.

5.2 Practical Implications for Indian Practitioners

Cloud Deployments: Leverage IndiaAI Mission subsidized GPU access (₹65/hour) with TORCH+CUDA for lowest latency and energy consumption.

CPU-Only Deployments: Use ONNX Runtime or OpenVINO rather than vanilla PyTorch for 9-72% energy savings.

Indian Language Applications: Evaluate Sarvam Edge models before cloud alternatives due to superior benchmark performance, zero per-query cost, privacy, and offline operation.

Privacy-Sensitive Enterprises: Consider Arinox CommandCORE "AI in a box" [6] (₹10 lakh+) for sovereign deployment with complete data control.

For distributed edge deployments across multiple locations (e.g., retail chains, branch offices, IoT sensor networks), hierarchical federated learning frameworks [11] can enable collaborative Small AI model improvement without centralizing sensitive data—combining on-device inference with privacy-preserving aggregation.

MSMEs: Deloitte India's GenW.AI platform [13], priced 50% below global alternatives, offers accessible AI adoption.

5.3 Policy Implications

India's approach offers lessons for emerging economies: infrastructure democratization through subsidized GPU access promotes broad-based innovation; sovereign model development ensures AI serves national priorities; privacy-by-design deployment models enable adoption without compromising data sovereignty [5].

5.4 Sustainability Implications

Configuration choices affect inference energy by 2.6x, and model selection by orders of magnitude more. India's focus on small, efficient models aligns sustainability goals with accessibility and affordability—where environmental, economic, and social goals converge [7].

VI. CONCLUSION AND FUTURE WORK

6.1 Summary of Contributions

This paper has presented the first comprehensive empirical investigation of Small, Cost-Efficient AI systems. We have demonstrated that specialized Small AI achieves 82-97% of frontier model performance at 1/50th to 1/500th the parameter count; Indian models achieve state-of-the-art regional language performance at 24M-150M parameters; runtime optimization reduces energy consumption by 37-89%; on-device economics enable near-zero marginal inference cost; and India's infrastructure investments create a foundation for scaled AI deployment.

6.2 India's Distinctive Approach

India's strategic orientation toward Small AI offers a compelling alternative to the scale-at-any-cost paradigm. This approach recognizes that for a nation of 1.45 billion people with extraordinary linguistic diversity, the relevant metrics are accessibility-per-citizen and cost-per-useful-task. The early results are promising, and whether this approach scales will depend on continued innovation, investment, and execution.

6.3 Future Research Directions

Six Directions emerge: Edge benchmarking for Indian smartphone chipsets; Indian language benchmark development across all 22 official languages; longitudinal tracking of IndiaAI Mission impact; cross-country comparative analysis with other emerging economies; and integration of Small AI with hierarchical federated learning frameworks for secure, scalable on-device learning across IoT ecosystems [11]; and long-term exploration of quantum-classical hybrid architectures for Small AI acceleration [12].

6.4 Concluding Remarks

Small AI represents not a compromise but a distinct and valuable direction—creating an ecosystem where massive generalist models handle challenging open-ended tasks while swarms of specialized small agents operate locally, privately, and efficiently. As the AI community confronts sustainability implications, Small AI offers a path forward that is both practically useful and environmentally responsible: intelligence that scales to everyone, not just those with data center budgets.

VII. REFERENCES

- [1] OpenAI, "GPT-4 Technical Report," arXiv preprint arXiv:2303.08774, Mar. 2023. Accessed: Apr. 10, 2026. [Online]. Available: <https://arxiv.org/pdf/2303.08774v4>
- [2] Liquid AI, "Liquid unveils 'Nanos': Extremely small foundation models that match frontier-model quality—running directly on everyday devices," Press Release, Sep. 25, 2025. Accessed: Apr. 10, 2026. [Online]. Available: <https://www.liquid.ai/press/liquid-unveils-nanos-extremely-small-foundation-models-that-match-frontier-model-quality--running-directly-on-everyday-devices>.
- [3] F. Durán, M. Martínez, P. Lago, and S. Martínez-Fernández, "Energy consumption of code small language models serving with runtime engines and execution providers," arXiv preprint arXiv:2412.15441, Dec. 2024. Accessed: Apr. 10, 2026. [Online]. Available: <https://arxiv.org/pdf/2412.15441>
- [4] Qwen: A. Yang et al., "Qwen2.5 Technical Report," arXiv, Dec. 2024. Accessed: Apr. 10, 2026. [Online]. Available: <https://arxiv.org/abs/2412.15115>
- [5] Odisha TV, "India to add 20,000 GPUs to strengthen its position in global AI leadership," Feb. 18, 2026. Accessed: Apr. 10, 2026. [Online]. Available: <https://odishatv.in/national/india-to-add-20000-gpus-to-strengthen-its-position-in-global-ai-leadership-11124590>
- [6] Hindustan Times, "India's first 'AI in a box': Arinox, KOGO unveil private, offline AI to shield enterprise data," Hindustan Times, Feb. 16, 2026.
- [7] L. F. W. Anthony, B. Kanding, and R. Selvan, "Carbontracker: Tracking and Predicting the Carbon Footprint of Training Deep Learning Models," arXiv preprint arXiv:2007.03051, Jul. 2020. Accessed: Apr. 10, 2026. [Online]. Available: <https://arxiv.org/abs/2007.03051>
- [8] Y. Tay, M. Dehghani, S. Abnar, H. W. Chung, W. Fedus, J. Rao, S. Narang, V. Q. Tran, D. Yogatama, and D. Metzler, "Scaling Laws vs. Model Architectures: How Does Inductive Bias Influence Scaling?," arXiv preprint arXiv:2207.10551, July 2022.

- [10] J. Ainslie, J. Lee-Thorp, M. de Jong, Y. Zemlyanskiy, F. Lebrón, and S. Sanghai, "GQA: Training Generalized Multi-Query Transformer Models from Multi-Head Checkpoints," arXiv preprint arXiv:2305.13245, May 2023.
- [11] Rai and R. Mathur, "A Hierarchical Federated Learning Framework for Secure and Scalable IoT Ecosystems: System Design, Implementation, and Performance Analysis," International Journal on Science and Technology (IJSAT), vol. 17, no. 1, Jan./Mar. 2026.
- [12] D. Mathur, S. Mathur, and A. Rai, "The Collaboration of the Trio of Quantum Computing, IOT and AI Powered by Cloud Architecture," International Journal of Enhanced Research in Science, Technology & Engineering, vol. 14, no. 5, p. 109, May 2025.
- [13] The Hindu, "Deloitte to unveil fully India-developed AI platform that offers scale, speed to global enterprises," The Hindu, Feb. 11, 2026.

Copyright & License:



© Authors retain the copyright of this article. This work is published under the Creative Commons Attribution 4.0 International License (CC BY 4.0), permitting unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.