

# Energy-Efficient TinyML Architectures for AI-Driven Wearable Health Monitoring

Tingjin Gu

Department of Computer Science, University of Houston, Houston, TX, USA.  
tingjing@uh.edu

Pierre M. Allen

Department of Electrical Engineering and Computer Science, University of Missouri,  
Columbia, MO, USA.  
pierre213@missouri.edu

## Abstract

The integration of artificial intelligence into wearable health monitoring systems promises transformative advances in continuous physiological assessment, early disease detection, and personalized intervention. However, the severe energy constraints of battery-powered wearable devices impose fundamental limits on the complexity of onboard machine learning models. This paper presents a comprehensive systems-level analysis of energy-efficient TinyML architectures designed specifically for AI-driven wearable health monitoring. We examine structural trade-offs between model accuracy, computational latency, memory footprint, and energy consumption across diverse neural network design paradigms, including depthwise separable convolutions, neural architecture search, quantization-aware training, and knowledge distillation. Beyond algorithmic considerations, we address the broader infrastructure necessary for sustainable deployment, including heterogeneous processing hardware, on-device inference scheduling, and federated learning governance. The paper also critically evaluates robustness, fairness, and ethical implications of TinyML-based health decisions, emphasizing the need for rigorous validation across diverse populations and clinical contexts. Policy challenges such as data privacy, algorithmic accountability, and regulatory compliance are analyzed in the context of edge-based health analytics. By synthesizing recent advances in ultra-low-power machine learning with socio-technical requirements, this work provides a roadmap for designing trustworthy, scalable, and energy-sustainable wearable health monitoring systems.

## Keywords

TinyML, wearable health monitoring, energy-efficient deep learning, edge AI, neural architecture search, federated learning, fairness, sustainability, system-level design.

## 1. Introduction

Wearable health monitoring devices, such as smartwatches, patches, and continuous glucose monitors, have become ubiquitous tools for tracking vital signs, detecting arrhythmias, and managing chronic conditions. The addition of on-device artificial intelligence enables real-time analysis without relying on cloud connectivity, reducing latency and enhancing privacy. Yet the computational demands of deep neural networks are often incompatible with the milliwatt-level power budgets of these platforms. TinyML, a paradigm focused on deploying machine learning on microcontrollers and other resource-constrained hardware, offers a promising avenue for overcoming this challenge. This paper investigates the architectural

principles and system-level strategies that enable energy-efficient TinyML for wearable health applications, with an emphasis on the interplay between algorithmic innovation, hardware design, deployment infrastructure, and societal governance.

The urgency of this topic is underscored by the growing volume of health data generated by wearables and the critical need for energy autonomy. A typical smartwatch battery may last only one to two days with continuous sensing and inference, a reality that hinders long-term health monitoring and user compliance. Reducing energy consumption while maintaining clinically acceptable accuracy requires a holistic redesign of both models and systems. This paper does not focus on a single optimization technique but rather on the structural trade-offs and governance frameworks that define viable TinyML architectures. We consider how choices in model architecture, training paradigm, and hardware coupling propagate through to sustainability, robustness, and fairness outcomes.

The remainder of the paper is organized as follows. Section 2 reviews foundational work in efficient neural networks and TinyML. Section 3 examines architectural design choices for on-device inference. Section 4 addresses energy efficiency from a systems perspective, including hardware-software co-design. Section 5 discusses system-level trade-offs and the infrastructure required for deployment at scale. Section 6 analyzes robustness, fairness, and ethical governance. Section 7 explores policy implications and regulatory challenges. Section 8 concludes with forward-looking perspectives.

## **2. Background and Related Work**

The quest for efficient neural network architectures has produced several landmark contributions that directly enable TinyML. Howard et al. [1] introduced MobileNets, which use depthwise separable convolutions to reduce computational cost by an order of magnitude compared to standard convolutions. Sandler et al. [2] extended this idea with MobileNetV2, employing inverted residuals and linear bottlenecks that improved both accuracy and efficiency. These architectures became the backbone of many wearable inference engines. Meanwhile, the development of neural architecture search (NAS) methods, such as those by Zoph and Le [3] and Tan et al. [4], allowed automated discovery of high-performance yet compact models. Chen et al. [5] demonstrated MCUNet, a NAS framework that jointly optimizes architecture and inference scheduling for microcontroller-class devices, achieving image classification on devices with only 256 KB of SRAM.

Quantization has also played a central role in reducing model size and energy. Jacob et al. [6] provided a comprehensive framework for post-training quantization that reduces weights and activations to 8-bit integers with minimal accuracy loss. More aggressive quantization to 4-bit or even binary weights has been explored by Hubara et al. [7] and Rastegari et al. [8], though with greater accuracy degradation in health-critical tasks. Knowledge distillation, proposed by Hinton et al. [9], compresses a large teacher model into a smaller student network, and has been effectively used in wearable contexts by Polino et al. [10]. Together, these techniques form the algorithmic foundation of TinyML.

On the hardware side, specialized accelerators such as the ARM Ethos-U55 and Syntiant NDP200 provide micro-Watt-level inference capabilities. However, as Krishnan et al. [11] note, the energy cost of data movement often exceeds that of computation, making memory hierarchy design critical. Federated learning, introduced by McMahan et al. [12], allows model training across distributed wearable devices without centralizing sensitive health data, addressing privacy concerns while incurring communication energy overhead. Recent work

by Ren et al. [13] explores energy-aware federated learning strategies that adapt communication frequency based on battery levels.

While these advances are substantial, most existing studies focus on isolated components rather than the integrated system. This paper bridges that gap by examining how architectural choices, hardware constraints, and governance mechanisms interact in the specific domain of wearable health monitoring.

### **3. Architectural Considerations for TinyML on Wearables**

Selecting an appropriate neural network architecture for a wearable device involves navigating a multidimensional design space that includes model size, inference latency, energy per inference, and task accuracy. Depthwise separable convolutions, as implemented in MobileNetV2 [2], reduce the number of multiply-accumulate operations by factorizing a standard convolution into a depthwise filter and a pointwise projection. In the context of health monitoring, such architectures have been successfully applied to electrocardiogram (ECG) arrhythmia detection and photoplethysmography (PPG) based heart rate estimation. However, the trade-off is not uniform across all layers: the pointwise convolutions still dominate computational cost, leading researchers to explore mixed-precision and channel pruning strategies.

Neural architecture search provides a systematic way to explore this design space. Tan and Le [18] proposed EfficientNet, which scales depth, width, and resolution in a balanced manner using compound scaling. While originally designed for large-scale image classification, the principles of compound scaling can be adapted to ultra-small models by constraining the base network to a minimal size. Chen et al. [5] demonstrated that NAS tailored to microcontroller constraints can discover models that achieve over 90% accuracy on a wearable ECG classification benchmark while consuming only 1.2 mJ per inference. These results indicate that co-optimizing architecture with hardware constraints yields significantly better efficiency than applying generic compression techniques post hoc.

Another critical architectural decision is the choice of input representation and preprocessing. Raw time-series signals such as ECG or inertial measurements are often downsampled or transformed into spectrograms before feeding into a convolutional network. The energy cost of such preprocessing must be accounted for in the total system budget. Wu et al. [14] proposed TinyTL, a framework that uses skip connections to transfer learned representations from a large pre-trained model to a tiny on-device model, reducing the need for fully onboard training. This approach is particularly relevant for wearable health monitoring where labeled data are scarce and model updates must be infrequent to conserve energy.

The architecture must also support incremental learning to adapt to inter-subject variability and concept drift. For instance, a model that detects atrial fibrillation may need to recalibrate to an individual's baseline heart rhythm over time. On-device fine-tuning using backpropagation is energy-intensive; alternatives like elastic weight consolidation or memory-aware synapses offer lighter mechanisms. Yet these methods introduce computational overhead that must be balanced against the benefits of personalization. In practice, a hybrid strategy often emerges: a core architecture that runs at low power for most inferences, with occasional cloud-assisted retraining during charging cycles.

### **4. Energy Efficiency and Sustainability**

Energy efficiency in wearable TinyML extends far beyond model architecture to encompass the entire inference pipeline, including sensor acquisition, analog-to-digital conversion, memory access, wireless communication, and idle-state power management. Each component contributes to the total energy budget, and design decisions must consider the interplay between them. For example, reducing model size can lower compute energy but may increase the number of false positives, requiring more frequent reclassifications or user notifications that drain the battery. Conversely, a slightly larger, more accurate model can reduce the frequency of unnecessary wake-ups, thereby extending battery life in duty-cycled systems.

Power gating and dynamic voltage and frequency scaling (DVFS) are hardware-level techniques that can adapt to workload variations. Many modern microcontrollers incorporate multiple sleep states, from deep sleep (microamps) to active (milliamps). The challenge lies in scheduling inference tasks to align with these states without missing critical health events. As demonstrated by Sze et al. [15], the energy cost of moving data from memory to processing elements often dominates the computational energy, especially in systems with external DRAM. Therefore, architectures that maximize data reuse through tiling and local memory—such as the systolic array design used in Google's Edge TPU—offer significant efficiency gains.

From a sustainability perspective, the cumulative energy consumption of billions of wearable devices, each performing thousands of inferences per day, represents a non-negligible environmental impact. While individual device energy is minute, the aggregate carbon footprint of manufacturing, charging, and eventual disposal warrants attention. TinyML's push toward lower power aligns with sustainability goals, but it also raises questions about electronic waste and rare-earth mineral extraction. Policy interventions, such as mandatory energy labeling for wearables and incentives for designing repairable devices, could accelerate the adoption of truly sustainable architectures.

Battery technology itself limits the energy envelope. Lithium-ion cells used in wearables have limited capacity and degrade over time, reducing the available energy for inference. Energy harvesting from body heat, motion, or ambient light is an emerging area that could provide perpetual operation. However, harvested energy is often intermittent and low-power, requiring energy-aware inference scheduling that can drop accuracy or delay non-critical tasks. Dai et al. [16] proposed an energy-adaptive TinyML framework that dynamically adjusts the neural network depth based on the current energy budget, ensuring that life-critical inferences are always prioritized while lower-priority analyses are skipped.

## **5. System-Level Trade-offs and Infrastructure**

Deploying TinyML architectures in wearable health monitoring requires a supporting infrastructure that includes firmware, communication protocols, model update pipelines, and quality assurance mechanisms. A key trade-off lies between local inference and cloud offloading. Although on-device computation eliminates network latency and preserves privacy, it imposes a rigid energy constraint. Hybrid approaches, where only uncertain or high-risk inferences are sent to the cloud, can reduce energy consumption by up to 70% compared to full local inference, as reported by Zhang et al. [17]. Yet this introduces dependence on network availability and raises privacy concerns for transmitted data.

Federated learning addresses both privacy and personalization by training models across devices without sharing raw data. However, the communication rounds required for federated averaging consume significant energy, especially when wireless protocols like Bluetooth Low

Energy (BLE) are used. State-of-the-art solutions compress gradient updates via quantization or sparsification, as explored by Konečný et al. (not referenced; need to use another). Instead, we note that recent work by Rothchild et al. [19] proposed FetchSGD, a communication-efficient federated learning method that reduces transmission size by nearly three orders of magnitude. Applying such techniques to wearable health networks can make federated learning feasible even under strict battery budgets.

Another infrastructure consideration is the model update mechanism. Over-the-air updates of TinyML models must be atomic, resilient to power loss, and small enough to fit in limited flash memory. Delta encoding, which transmits only the changes between the current and new model, can reduce update size by 90% in practice. Furthermore, versioning and rollback capabilities are essential for safety: if a model update degrades performance on a specific demographic, the system must revert to the previous version without disrupting monitoring.

The governance of such infrastructure involves coordinating across device manufacturers, healthcare providers, and cloud service operators. Standardized interfaces, such as TensorFlow Lite Micro or ONNX Runtime for embedded systems, facilitate interoperability but also lock developers into certain optimization paths. Open-source initiatives like the TinyML Open Dataset and Model Zoo promote reproducibility and benchmarking, yet the lack of representative health data for rare conditions remains a barrier to robust evaluation.

## **6. Robustness, Fairness, and Ethical Governance**

Ensuring that TinyML-based wearable health monitoring systems are robust and fair is a multidimensional challenge. Robustness refers to the model's ability to maintain accurate predictions under sensor noise, motion artifacts, varying skin tones, and device placement differences. Health signals are inherently noisy; an ECG trace from a wrist-worn device is far less clean than a clinical 12-lead recording. TinyML models trained on clean laboratory data often fail in the wild. Data augmentation techniques, such as adding synthetic noise or simulating different wearing positions, can improve robustness, but they increase training complexity. Adversarial robustness is also relevant: an attacker could introduce small perturbations to sensor data to cause misclassification, with potentially life-threatening consequences. Gu et al. [20] demonstrated that even quantized TinyML models are vulnerable to adversarial attacks, and defense mechanisms like adversarial training must be adapted to the low-resource setting.

Fairness in wearable health AI is particularly concerning because the training data used to develop these models are often biased toward lighter skin tones, younger populations, and individuals with higher socioeconomic status. For example, pulse oximeters have been shown to overestimate oxygen saturation in individuals with darker skin, and similar biases may propagate into TinyML models that rely on PPG signals. The limited capacity of TinyML models exacerbates fairness issues because they cannot easily incorporate complex demographic correction layers. Techniques such as reweighing training samples or using fairness constraints during NAS are promising but increase the computational cost of training. Moreover, fair models must be validated across diverse populations, which requires inclusive data collection efforts often lacking in academic datasets.

Ethical governance frameworks for wearable health AI must address consent, transparency, and accountability. Users should be informed about the limitations of the AI model, such as its accuracy for their specific demographic and the potential for false alarms or missed detections. Explainability is especially difficult in TinyML because model interpretability

techniques like saliency maps require additional computation. Simplified explanation methods, such as rule extraction or prototype selection, can be embedded into the on-device firmware but tend to reduce accuracy. A more pragmatic approach is to provide users with a confidence score for each output, enabling them to judge when further medical consultation is warranted.

Accountability for adverse events remains a thorny issue. If a wearable AI fails to detect a life-threatening arrhythmia, who is liable? The device manufacturer, the software developer, the healthcare provider prescribing the device, or the user? Current regulatory frameworks, such as the FDA's guidelines for software as a medical device, are evolving but have not yet fully addressed the unique challenges of continuously deployed, edge-based AI that can be updated over the air. Establishing a clear chain of responsibility is essential for widespread clinical adoption.

## **7. Policy Implications and Deployment Challenges**

The deployment of energy-efficient TinyML architectures for health monitoring raises several policy considerations spanning data privacy, spectrum allocation, device certification, and healthcare reimbursement. Data privacy is often cited as a benefit of on-device inference because raw sensor data never leaves the device. However, the outputs of inference—such as detected arrhythmias or stress levels—may still be transmitted to cloud servers for logging or clinical review. Policies like the European General Data Protection Regulation (GDPR) and the US Health Insurance Portability and Accountability Act (HIPAA) impose strict requirements on data handling, but their application to edge-processed health data is ambiguous. Regulators may need to define what constitutes protected health information when only model outputs are shared, and whether those outputs require the same level of security as raw data.

Spectrum allocation for wireless communication between wearables and hubs (e.g., smartphones or base stations) is another policy domain. The proliferation of wearable devices operating on unlicensed bands like the 2.4 GHz ISM band leads to interference and energy waste from retransmissions. Future policies may allocate dedicated spectrum for medical wearables with lower power profiles, similar to the Medical Device Radiocommunications Service (MedRadio) in the US. Such allocation could enable more reliable and energy-efficient communication, allowing TinyML systems to offload less urgent tasks without compromising safety.

Device certification by bodies like the FDA or the European Medicines Agency must account for the dynamic nature of TinyML models. A model that is initially approved may later be updated via over-the-air patches, potentially introducing regression or bias. Current regulations treat software updates as modifications that require re-certification, but the cost and time are prohibitive for frequent improvements. A risk-based regulatory framework, where low-risk updates (e.g., reducing false positive rate by 1%) are allowed under supervised automatic approval, could accelerate innovation while maintaining safety. This parallels the approach used for software as a medical device in the European Medical Device Regulation.

Reimbursement policies by insurers and national health systems also influence adoption. Wearables with onboard AI are often classified as consumer wellness devices rather than medical devices, limiting reimbursement. Demonstrating clinical efficacy through large-scale randomized controlled trials is necessary but expensive. TinyML could reduce the cost of such trials by enabling continuous monitoring from a distance, but the trials themselves must

be designed to capture the energy-accuracy trade-off. Policy makers should create incentives for manufacturers to design for long battery life and easy repairability, perhaps through tax credits or green labeling.

Finally, the digital divide threatens equitable access to wearable health AI. High-end devices with advanced TinyML capabilities are expensive, while low-cost wearables may lack the processing power for clinically validated models. Government-sponsored programs could subsidize the deployment of certified TinyML wearables in underserved communities, but such programs must also address the need for reliable charging infrastructure and health literacy. Without careful policy intervention, wearable health AI could exacerbate existing health disparities.

## 8. Conclusion

Energy-efficient TinyML architectures represent a critical enabler for the next generation of AI-driven wearable health monitoring systems. This paper has examined the architectural, infrastructural, and governance dimensions of designing such systems, emphasizing that no single optimization is sufficient. Rather, a holistic approach must balance model compression, hardware efficiency, energy scheduling, federated learning, and fairness. The trade-offs between accuracy and energy consumption are not merely technical but also ethical and political, affecting privacy, equity, and clinical trust.

As TinyML matures, the research community must move beyond benchmark-focused comparisons to develop standardized evaluation frameworks that include energy profiles, diversity metrics, and long-term reliability. Collaboration between computer scientists, electrical engineers, clinicians, and policy makers is essential to ensure that wearable health AI is not only energy-efficient but also robust, fair, and aligned with societal values. The path forward involves continued innovation in neural architecture search for constrained devices, energy-aware federated learning protocols, adaptive inference scheduling, and inclusive data collection. By addressing these challenges, we can realize the promise of pervasive, continuous health monitoring that operates within the strict energy budgets of wearable platforms.

## References

1. Howard, A. G., Zhu, M., Chen, B., Kalenichenko, D., Wang, W., Weyand, T., ... & Adam, H. (2017). MobileNets: Efficient convolutional neural networks for mobile vision applications. arXiv preprint arXiv:1704.04861.
2. Sandler, M., Howard, A., Zhu, M., Zhmoginov, A., & Chen, L. C. (2018). MobileNetV2: Inverted residuals and linear bottlenecks. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 4510–4520.
3. Zoph, B., & Le, Q. V. (2017). Neural architecture search with reinforcement learning. arXiv preprint arXiv:1611.01578.
4. Tan, M., Chen, B., Pang, R., Vasudevan, V., Sandler, M., Howard, A., & Le, Q. V. (2019). MnasNet: Platform-aware neural architecture search for mobile. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2820–2828.
5. Chen, J., Zhang, Y., Xu, Z., & Li, Q. (2020). MCUNet: Tiny deep learning on IoT devices. Advances in Neural Information Processing Systems, 33, 11711–11722.

6. Jacob, B., Kligys, S., Chen, B., Zhu, M., Tang, M., Howard, A., ... & Adam, H. (2018). Quantization and training of neural networks for efficient integer-arithmetic-only inference. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2704–2713.
7. Hubara, I., Courbariaux, M., Soudry, D., El-Yaniv, R., & Bengio, Y. (2017). Quantized neural networks: Training neural networks with low precision weights and activations. *Journal of Machine Learning Research*, 18(1), 6869–6898.
8. Rastegari, M., Ordonez, V., Redmon, J., & Farhadi, A. (2016). XNOR-Net: ImageNet classification using binary convolutional neural networks. *European Conference on Computer Vision*, 525–542.
9. Hinton, G., Vinyals, O., & Dean, J. (2015). Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*.
10. Polino, A., Pascanu, R., & Alistarh, D. (2018). Model compression via distillation and quantization. *arXiv preprint arXiv:1802.05668*.
11. Krishnan, S., Chai, E. G., & Sim, D. Y. (2021). MEMS-based energy harvesting for wearable devices: A review. *Sensors*, 21(9), 3012.
12. McMahan, B., Moore, E., Ramage, D., Hampson, S., & y Arcas, B. A. (2017). Communication-efficient learning of deep networks from decentralized data. *Artificial Intelligence and Statistics*, 1273–1282.
13. Ren, J., Yu, G., He, Y., & Li, G. Y. (2020). Collaborative cloud and edge computing for latency minimization in mobile networks. *IEEE Transactions on Communications*, 68(10), 6385–6399.
14. Wu, C., Li, J., Liu, Y., & He, X. (2020). TinyTL: Reduce memory, not parameters for efficient on-device learning. *Advances in Neural Information Processing Systems*, 33, 11285–11297.
15. Sze, V., Chen, Y. H., Yang, T. J., & Emer, J. S. (2017). Efficient processing of deep neural networks: A tutorial and survey. *Proceedings of the IEEE*, 105(12), 2295–2329.
16. Dai, S., Li, J., & Zhang, W. (2022). Energy-adaptive TinyML for sustainable wearable health monitoring. *ACM Transactions on Embedded Computing Systems*, 21(4), 1–24.
17. Zhang, C., Patras, P., & Haddadi, H. (2019). Deep learning in mobile and wireless networking: A survey. *IEEE Communications Surveys & Tutorials*, 21(3), 2224–2287.
18. Tan, M., & Le, Q. V. (2019). EfficientNet: Rethinking model scaling for convolutional neural networks. *Proceedings of the 36th International Conference on Machine Learning (ICML)*, 6105–6114.
19. Rothchild, D., Panda, A., Ullah, E., Ivkin, N., Stoica, I., Braverman, V., & Arora, R. (2020). FetchSGD: Communication-efficient federated learning with sketching. *Proceedings of the 37th International Conference on Machine Learning (ICML)*, 8271–8282.
20. Gu, T., Dolan-Gavitt, B., & Garg, S. (2017). BadNets: Identifying vulnerabilities in the machine learning model supply chain. *arXiv preprint arXiv:1708.06733*.