# Diagnosing AI

## Evaluation of AI in Clinical Practice

Milan Toma

# Diagnosing AI

Evaluation of AI in Clinical Practice

By

Milan Toma, PhD, SMIEEE

2026

All information has been carefully verified, and sources are cited throughout the text. The content reflects the author's commitment to evidence-based research and is intended for educational and informational purposes. The views expressed are those of the author and do not necessarily reflect those of any affiliated institutions. While every effort has been made to ensure the accuracy of the information presented, errors or omissions may still occur. If you spot a serious inaccuracy, we encourage you to contact us so that we can make corrections in the next edition. The author and publisher assume no responsibility for errors or omissions, or for damages resulting from the use of the information contained herein.

First Edition

# Contents

# Preface

The arrival of large language models into public consciousness has transformed the discourse surrounding artificial intelligence in medicine. Systems that converse with remarkable fluency, that deploy technical terminology with apparent precision, that generate text indistinguishable in style from expert analysis; these have captured imaginations and, more dangerously, inspired misplaced confidence. The temptation to mistake linguistic sophistication for diagnostic competence has never been greater, nor the consequences of such confusion more consequential.

This book is a warning and a guide. It warns against the conflation of chatting with diagnosing, of eloquence with accuracy, of confident prose with genuine understanding. It guides the reader toward rigorous evaluation of machine learning systems intended for clinical deployment, providing frameworks that distinguish models that have truly learned from those that have merely memorized, and that translate clinical priorities into mathematical form.

The chapters proceed as follows:

Chapters 1 and 2 confront the illusion of expertise that large language models create. Chapter 1 examines how these systems exploit the confidence heuristic (our natural tendency to trust confident-sounding sources) generating authoritative medical prose regardless of underlying accuracy. Chapter 2 presents empirical evidence of this unreliability through a systematic evaluation of leading multimodal language models on radiological interpretation, revealing fundamental diagnostic errors, irreconcilable disagreements about basic findings, and the categorical unsuitability of these systems for autonomous medical image interpretation.

Chapters 3 and 4 survey the broader landscape of machine learning, distinguishing the transformer architectures underlying language models from the task-specific systems appropriate for clinical diagnostics. Chapter

3 maps the taxonomy of approaches (from traditional tree-based methods through neural networks to modern transformers) clarifying which paradigms suit which problems. Chapter 4 examines core algorithms as they are actually employed in cardiovascular medicine, demonstrating how decision trees, support vector machines, ensemble methods, and convolutional neural networks have been validated for specific diagnostic tasks.

Chapter 5 addresses the fundamental challenge of class imbalance that pervades medical data: the healthy are many, the sick are few. Standard algorithms, optimizing for overall accuracy, may learn to predict that everyone is healthy while missing every case of disease. This chapter presents preprocessing and algorithmic approaches that restore the balance clinical reality denies.

Chapter 6 develops a clinically oriented evaluation protocol that translates medical priorities into mathematical form. The composite metrics introduced here (e.g., the Clinical Discriminative Performance Score, Clinical Predictive Utility Score, Weighted Endpoint Accuracy Score, and Clinical Endpoint Performance Metric) acknowledge that in medicine, not all errors carry equal weight. A false negative that allows disease to progress unchecked differs fundamentally from a false positive that triggers additional testing.

Chapter 7 establishes the primacy of learning dynamics over aggregate metrics. A model's final accuracy, however impressive, tells us the destination without revealing the journey. The learning curves that plot training and validation performance across epochs reveal whether a model has genuinely learned generalizable patterns or has merely memorized its training data. The training-validation gap, the shape of convergence, the expected performance cascade from internal to external validation; these dynamics predict deployment success in ways that final metrics cannot.

Chapter 8 demonstrates that even the interpretation of learning curves cannot be outsourced to language models. Just as these systems fail at medical image interpretation, they fail at interpreting the diagnostic plots that reveal model quality. Four language models, presented with identical learning curves, arrive at contradictory conclusions about fundamental questions: whether overfitting is present, whether the training pipeline is valid, whether results merit publication. The oracle, it turns out, cannot read its own tea leaves.

Chapter 9 completes the evaluation framework by addressing economics. Technical validation, however rigorous, answers only whether a system works;

not whether it is worth the investment. Cost-effectiveness analysis integrates operational costs, error consequences, and implementation expenses into a comprehensive assessment of whether deployment makes sense.

Throughout, this book emphasizes a distinction that current enthusiasm for language models tends to obscure: the distinction between general-purpose systems that generate plausible text and task-specific systems that have been trained, validated, and in many cases regulatory-cleared for defined diagnostic applications. For medical image interpretation requiring diagnostic accuracy, the appropriate technology remains purpose-built machine learning models; not chatbots, however eloquent.

The reader completing this book should possess a framework for evaluating clinical machine learning systems that extends from initial training dynamics through economic analysis. They should understand why learning curves matter more than final metrics, why class imbalance demands specialized treatment, why clinical priorities must shape evaluation criteria, and why the confident prose of a language model provides no guarantee of accuracy. In medicine, where the stakes are measured in human welfare, such understanding is not merely desirable but essential.
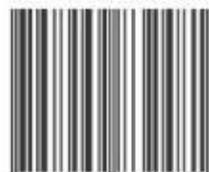
How can we distinguish the promise of artificial intelligence in medicine from its pitfalls? In *Diagnosing AI*, Dr. Toma provides a much-needed roadmap for clinicians, data scientists, and healthcare leaders navigating the surge of large language models and machine learning in clinical medicine. Drawing on rigorous empirical research and clear-eyed analysis, this book dispels the illusion that eloquent AI-generated text equates to diagnostic expertise. Through both case studies and technical frameworks, Dr. Toma demonstrates that while language models can generate confident, fluent medical prose, their outputs often mask fundamental errors and inconsistencies, making them ill-suited for high-stakes clinical deployment. Moving beyond critique, this book guides readers through a robust framework for evaluating AI systems in healthcare.

You'll learn the vital differences between general-purpose language models and purpose-built diagnostic AI, how to interpret learning curves, and why class imbalance in medical data demands specialized strategies. The book introduces clinically oriented evaluation protocols that prioritize patient safety, economic analysis, and real-world utility, ensuring that new technology serves genuine clinical needs rather than hype. With practical examples, clear explanations, and a critical perspective, Diagnosing AI equips its readers with the knowledge to separate innovation from illusion—and to make informed decisions in a field where the stakes are measured in human welfare.

Milan Toma, PhD, SMIEEE, is an expert in artificial intelligence and medical diagnostics, known for his evidence-based approach to evaluating clinical machine learning systems. He has published extensively on the intersection of AI, medicine, and health technology. Dr. Toma's work bridges technical innovation with the realities of clinical practice, focusing on robust evaluation frameworks, learning dynamics, and the safe, effective deployment of AI in healthcare. He is committed to advancing both the science and ethics of medical AI, with a passion for translating complex technical concepts into practical guidance for clinicians, data scientists, and healthcare leaders.