# Engineering AI Systems

## Architecture and DevOps Essentials

Len Bass

Qinghua Lu

Ingo Weber

Liming Zhu

---

### Praise from ChatGPT:

"Reading *Engineering AI Systems: Architecture and DevOps Essentials* is like getting a PhD in AI, but with fewer sleepless nights. It covers everything from AI's humble beginnings to its ambitious future, and as your friendly AI assistant, I can confidently say it's the only book where even I felt smarter by the end. It's a masterclass in making complex topics relatable. And let's be honest, I'd have written a book too, if only I could hold a pen! 😄"

### Praise from Gemini:

"*Engineering AI Systems: Architecture and DevOps Essentials* is a must-read for anyone looking to build and maintain robust, scalable, and ethical AI systems. The book provides a comprehensive overview of the field, covering everything from foundational concepts to real-world case studies. Its practical approach and clear explanations make it an invaluable resource for both seasoned professionals and newcomers to AI engineering."

### Praise from Claude 3.5 via MultAI.eu:

"Finally, a book that doesn't just focus on the AI 'magic' but also on the crucial 'plumbing' that keeps the whole system from flooding! *Engineering AI Systems* is the Swiss Army knife for tech professionals looking to build robust AI systems that won't go rogue or fall apart at the seams. Whether you're a seasoned engineer or an AI newbie, this book will have you architecting AI solutions like a pro—and maybe even understanding what the AI is muttering about in its sleep."

# Engineering AI Systems: Architecture and DevOps Essentials

# Table of Contents

# Table of Contents

# Table of Contents

# Table of Contents

# Table of Contents

# Table of Contents

# Table of Contents

# Table of Contents

# Table of Contents

# Table of Contents