# The Definitive Guide to Azure Data Engineering

## Modern ELT, DevOps, and Analytics on the Azure Cloud Platform

Ron C. L'Esteve

*The Definitive Guide to Azure Data Engineering: Modern ELT, DevOps, and Analytics on the Azure Cloud Platform*

Ron C. L'Esteve
Chicago, IL, USA

*For Mom and Dad.*

# Table of Contents

xvi

# About the Author

**Ron C. L'Esteve** is a professional author residing in Chicago, IL, USA. His passion for Azure Data Engineering originates from his deep experience with designing, implementing, and delivering modern Azure data projects for numerous clients. Ron is a trusted technology leader and digital innovation strategist, responsible for scaling key data architectures, defining the road map and strategy for the future of data and business intelligence (BI) needs, and challenging customers to grow by thoroughly understanding the fluid business opportunities and enabling change by translating them into high-quality and sustainable technical solutions that solve the most complex challenges and promote digital innovation and transformation. He applies a practical and business-oriented approach of taking transformational ideas from concept to scale. Ron is an advocate for data excellence across industries and consulting practices and empowers self-service data, BI, and AI through his contributions to the Microsoft technical community.

# About the Technical Reviewer

**Greg Low** is one of the better-known database consultants in the world. In addition to deep technical skills, Greg has experience with business and project management and is known for his pragmatic approach to solving issues. His skill levels at dealing with complex situations and his intricate knowledge of the industry have seen him cut through difficult problems.

Microsoft has specifically recognized his capabilities and appointed him to the Regional Director program. They describe it as consisting of "150 of the world's top technology visionaries chosen specifically for their proven cross-platform expertise, community leadership, and commitment to business results."

Greg leads a boutique data consultancy firm called SQL Down Under. His clients range from large tier 1 organizations to start-ups.

Greg is a long-term Data Platform MVP and considered one of the foremost consultants in the world on SQL Server and Microsoft data-related technologies. He has provided architectural guidance for some of the largest SQL Server implementations in the world and helped them to resolve complex issues. Greg was one of the two people first appointed as SQL Server Masters worldwide. Microsoft use him to train their own staff.

For several years, Greg served on the global board for the Professional Association for SQL Server. He is particularly proud of having helped it triple the size of its community and, more importantly to him, taken it from being 90% US based to being a truly global community with 60% of chapters outside the United States.

A talented trainer and presenter, Greg is known for his ability to explain complex concepts with great clarity to people of all skill levels. He is regularly invited to present at top-level tier 1 conferences around the world. Greg's SQL Down Under podcast has a regular audience of over 40,000 listeners.

Outside of work and family, Greg's current main passion is learning Mandarin Chinese, and he is determined to learn to read, write, speak, and understand it clearly.

# Acknowledgments

Writing this book has been both a solitary and accompanied journey with sacrifices and victories along the way. Thank you to all who have supported me on the path to completing this book.

# Introduction

With the numerous cloud computing technologies being at the forefront of the modern-day data architectural and engineering platforms, Microsoft Azure's cloud platform has contributed over 200 products and services that have been specifically designed to solve complex data challenges, empower self-service data engineering, and pave the way for the future of data and AI.

Navigating through these many offerings in the Azure Data Platform can become daunting for aspiring Azure Data Engineers, architects, consultants, and organizations that are seeking to build scalable, performant, and production-ready data solutions. This book is intended to uncover many of the complexities within the Azure data ecosystem with ease through structured end-to-end scenario-based demonstrations, exercises, and reusable architectural patterns for working with data in Azure and building highly performant data ingestion and ELT pipelines.

As Azure continues to introduce numerous data services to their ever-growing and evolving platform, this book will demystify many of the complexities of Azure Data Engineering with ease and introduce you to tried, tested, and production-ready patterns and pipelines for data of all different volumes, varieties, and velocities.

Additionally, you will be introduced to the many capabilities of bringing value and insights to your data through real-time and advanced analytics, continuously integrating and deploying your data ingestion pipelines, and getting started with many Azure data services to help you progress through your journey within the Azure Data Engineering ecosystem.