# Designing Data-Intensive Applications

*The Big Ideas Behind Reliable, Scalable, and Maintainable Systems*

*Martin Kleppmann*

# Table of Contents

# Part III.    Derived Data