# Database Internals

## *A Deep Dive into How Distributed Data Systems Work*

*Alex Petrov*

**Database Internals**

by Alex Petrov

# Table of Contents