# Mapping Data Flows in Azure Data Factory

Building Scalable ETL Projects in the Microsoft Cloud

Mark Kromer

**apress®**

# Mapping Data Flows in Azure Data Factory

## Building Scalable ETL Projects in the Microsoft Cloud

**Mark Kromer**

Apress®

*Mapping Data Flows in Azure Data Factory: Building Scalable ETL Projects in the Microsoft Cloud*

Mark Kromer
SNOHOMISH, WA, USA

*This book is dedicated to my loving wife Stacy and
our boys Ethan and Jude. Thank you for putting up with
my late hours working on data analytics and writing this book!*

# Table of Contents