DATABASE DESIGN BOOK

5205

Africa, Vol.

Japan, Vols

FOOTAGE

RACK

41 197ESPr 3.12/13

18

6.43

8/7.47

2/-2--

7/18

1/42

MINIMAL MODELING 5D5.16/3-

LEARN HOW TO GET FROM BUSINESS REQUIREMENTS TO A DATABASE SCHEMA

Alexey Makhotkin

449381

DATABASE DESIGN Book

learn how to get from business requirements to a database schema

ALEXEY MAKHOTKIN

2025

© 2025 Alexey Makhotkin All rights reserved.

No part of this eBook may be reproduced, stored in a retrieval system, or transmitted in any form or by any means—electronic, mechanical, photocopying, recording, or otherwise—without the prior written permission of the copyright owner, except for the use of brief quotations in book reviews or scholarly articles.

First Edition (revision: 2025-04-14) This eBook is licensed for your personal use only. Unauthorized distribution is prohibited.

Cover design by Anna Golde (https://www.linkedin.com/in/annagolde/).

Database Design Book: https://databasedesignbook.com/

Minimal Modeling Substack: https://minimalmodeling.substack.com/

Twitter: https://twitter.com/alexeymakhotkin

CONTENTS

1	Intro	oduction 7
	1.1	Who is this book for? 7
	1.2	What level of understanding is this book aimed
		at? 7
	1.3	What you would get from reading this book 8
	1.4	Who I am 9
2	Buil	ding a logical model 10
	2.1	System requirements as input 10
	2.2	Who do we write the logical model for? 11
		2.2.1 For yourself 11
		2.2.2 For software developers 12
		2.2.3 For the project group 12
	2.3	Can I skip the logical model? 13
	2.4	Can I use an ERD (entity-relationship diagram)
		instead? 13
	2.5	Elements of a logical model 14
	2.6	The process 14
	2.7	Anchors: introduction 15
		2.7.1 Example: posts 15
		2.7.2 A list of anchors 16
		2.7.3 Example: invoices 17
		2.7.4 Anchor IDs 18
	2.8	Attributes: introduction 19
		2.8.1 Attributes: definition 19
		2.8.2 Attributes and anchors 21
		2.8.3 Human-readable questions 21
		2.8.4 Example values 22
		2.8.5 Data types and types of data 22
		2.8.6 What is not an attribute? 23
		2.8.7 How to confirm that all the attributes have
		been listed? 23
	2.9	Links: introduction 24
		2.9.1 Links: definition 25
		2.9.2 Links: pair of anchors 26

Links: cardinality 2.9.3 26 Cardinality is a business concern 2.9.4 27 Links: sentences 28 2.9.5 2.9.6 False links and unique pairs of IDs 29 2.10 More on anchors 31 2.10.1 Unique attributes 31 2.10.2 Optional unique attributes 32 2.10.3 External IDs 33 2.10.4 External ID enumeration 34 2.10.5 Several unique IDs 34 2.10.6 Implementing physical IDs 35 2.11 Handling time in the logical model 36 "Hello world" use case: podcast catalog 3 41 **Business requirements** 3.1 41 Anchors 3.2 43 Attributes 3.3 43 3.4 Links 47 Cross-checking the requirements 3.5 50 A diagram 3.6 51 Are we there yet? 3.7 51 3.8 Evolving the system 52 Building a physical schema 53 4 Many table design strategies are possible 4.1 53 Table-per-anchor table design strategy 5 55 Action plan 5.1 55 Anchors: choosing table names 5.2 56 Attributes: choosing column names 5.3 57 Attributes: choose column data types 5.4 59 5.4.1 Recommended physical types for SQL databases: summary 61 Strings 62 5.4.2 Integer numbers 62 5.4.3Monetary amounts 63 5.4.4Numeric values 5.4.563 Yes/no values 64 5.4.6 Either/or/or values 5.4.764 5.4.8 Dates 65 Date with time in UTC timezone 65 5.4.95.4.10 Date with time in a specific timezone 65

		5.4.11 Timezone names 65
		5.4.12 Binary blobs 66
	5.5	Links 66
		5.5.1 One-to-many (1:M) links 66
		5.5.2 Many-to-many (M:N) links 68
	5.6	Podcast catalog: a complete physical schema 70
		5.6.1 CREATE TABLE statements 73
		5.6.2 Audio file and cover images 74
	5.7	Physical ID design 75
		5.7.1 Case study: a tiny CMS 75
		5.7.2 The maximum number of items 77
		5.7.3 Reaching the maximum number of items 77
		5.7.4 Space taken by IDs 78
		5.7.5 Disk space is time 79
		5.7.6 Storage density 79
		5.7.7 UUIDs as anchor IDs 80
		5.7.8 Countries, currencies, languages: well-known
		anchors 81
		5.7.9 Countries, currencies, and languages in
		your business 82
	5.8	Handling time in the physical model 83
6	Othe	er table design strategies 86
	6.1	Table design concerns86
	6.2	What is the recommended table strategy? 87
	6.3	Table-per-anchor, revisited88
	6.4	Side tables 89
		6.4.1 Naming and composition of side tables 90
		6.4.2 A stopgap table 91
	6.5	JSON columns 92
7	Seco	ondary data 94
	7.1	Cached column example 95
	7.2	There is no free lunch 96
	7.3	Cached column is not an attribute 97
	7.4	Discovering secondary data 97
8	Evol	lving your database 99
	8.1	Elementary database migrations 100
	8.2	Table rewrite 101
	8.3	Adding an attribute 104
		8.3.1 Step 1: Update logical model 105

Step 2. Run database migration 8.3.2 106 8.3.3 Step 3. Update code 106 8.4 Dealing with table rewrite 107 Movie tickets: repeated sales pattern 109 9 **Business requirements** 109 9.1 9.2 Per-department modeling 110 Movies department 9.3 112 Maintenance department 9.4 112 Movie schedule department 9.5 114 9.6 Tickets department 116 9.7 A diagram 119 Conclusion 9.8 119 10 Books and washing machines: polymorphic data pattern 120 10.1 Business requirements 120 10.2 Per-department modeling 121 10.3 Key insight: generic anchor vs specific anchors 121 10.4 Multiplexing on item type 123 10.5 Tangled links 124 10.6 A diagram 125 10.7 Table-per-anchor approach 126 10.8 Polymorphic table design strategy 127 10.8.1 JSON-based columns 127 10.8.2 Physical schema 127 10.8.3 Storing links in JSON 128 10.8.4 Table design concerns, revisited 129 10.8.5 Documenting physical storage 129 11 Practicalities 133 11.0.1 Document-based catalog 133 11.0.2 Spreadsheet-based catalog 133 11.0.3 How much to write 135 11.0.4 Lightweight designs 137

INTRODUCTION

This is a book on database design. You can find book updates and extra material at https://databasedesignbook.com/.

1.1 WHO IS THIS BOOK FOR?

The goal of this book is to help you get from a vague idea of what you need to implement (e.g., "I need to build a website to manage schedule and instructor appointments for our gym") to the comprehensive definition of database tables.

If you're a beginner, this book is for you. If you know that you need a database for your project, but are not sure how to begin, this is for you, too. If you've learned something about database design theory, but not sure if your design would work: you'll find help here. And if you struggle with some parts of the business domain, while the rest is clear, the approach laid out in this book will give clarity.

1.2 WHAT LEVEL OF UNDERSTANDING IS THIS BOOK AIMED AT?

The first part of the book is completely database-agnostic. All you need to understand is how your business works. This does not have to be a business, either: just some problem that necessitates a database.

The second part of the book requires some level of familiarity with common databases: how the tables are created, what physical data types exist, what are primary key and index, how the tables would be queried, and how to insert and update data.

We begin by assuming that you will use one of the traditional relational database servers, such as MySQL or PostgreSQL.

You do not need to understand database theory, such as normal forms.

1.3 WHAT YOU WOULD GET FROM READING THIS BOOK

The idea of the book is to make every step of table design actionable.

You will learn **how to extract the logical model** from a freeform description of the problem. This will also help to clarify your understanding of the problem, because the process forces you to be more precise and reduces hand waving.

You will learn **how the logical model is translated into a physical table design**. There are a dozen possible ways to represent logical models in physical space. First we discuss one such design strategy (a table per anchor); later we'll see why one might want to occasionally use a different one .

The book will help you answer a common question: **how can I prove, even to myself, that my design is complete**? Also, you will be able to defend this design, be it at a job interview or in a database exam.

A logical model is a very good way to explain your design to other people, for example to your team members or business stakeholders. This way of presenting your design immediately answers many questions that people may have. Because of that they can give better feedback and confirm that you're on the same page.

You will learn to **think separately about logical and physical aspects of your database**. This brings a lot of clarity: your problem lies either on the one or the other side, and you can save a lot of complexity if the other side is fixed.

Our approach is strongly rooted in relational modeling. However, it is **designed to accommodate non-relational NoSQL systems**, and even more specialized technologies such as Amazon S₃. As mentioned earlier, the logical model is not dependent on any specific database server system: you just have to choose a suitable table design strategy.

1.4 WHO I AM

My name is Alexey Makhotkin. I have been working with databases for more than 25 years in various roles: as a software engineer, database administrator, team lead, head of software engineering. I've built and helped to build dozens of schemas over the years.

In 2021 I started the "Minimal Modeling" substack: https://minimalmodeling.substack.com/, trying to summarize what I have learned. This book is another step. You can contact me at squadette@gmail.com.

BUILDING A LOGICAL MODEL

Suppose that you want to **build a software system that stores some data in a database**. It may be one of hundreds possible kinds of software systems, such as:

- ecommerce: placing and fulfilling orders;
- money-tracking software, such as personal budget tracker or a business accounting system;
- social network with users, posts, and comments;
- tracking your progress and achievements in a role-playing game;
- etc., etc.

In each of those systems there is some sort of underlying database, with tables, columns, rows and datatypes. Your task is to design and build this database. Then the system will use it to store and retrieve the data.

How do you design and build this database? This chapter, as well as the entire book, aims to answer exactly this question.

2.1 SYSTEM REQUIREMENTS AS INPUT

What exactly are we building? In this book, it is your responsibility to specify the system that you are building. You need to gather and provide answers about what the system is supposed to do.

The point of this chapter is that your understanding of the requirements will gradually evolve.

In the beginning you have just **an informal text that roughly specifies what the system does**. Imagine that you sit down with somebody and explain what the system is supposed to do. They

"HELLO WORLD" USE CASE: PODCAST CATALOG

Let's design our first toy but realistic use case: a catalog of podcasts.

We're going to design an MVP: a minimum viable product. More complex use cases would be found later in the book.

We discuss how to extend this scope with more interesting features (e.g., categories) later in this section.

We're only going to build a logical model here. In the second part of the book we'll cover the physical design directly based on this model.

3.1 BUSINESS REQUIREMENTS

Let's use a screenshot of an episode from Apple Podcast and see what information we have here (see the next page).

Here is a list of elements that we can see here:

- cover image;
- air date (16 January 2025);
- episode number (26);
- length (1 hr 3 min);
- episode title ("*Debanking explained*");
- name of the show ("Complex Systems with Patrick McKenzie");
- episode description ("*In this episode*, ... ").

If you scroll down, there will be a bit more information, but let's begin with this. Oh wait, the actual audio file should prob-

BUILDING A PHYSICAL SCHEMA

So, now that we have a catalog of anchors, attributes and links, we can build a physical schema: the actual database tables. When you have a catalog, 75% of work is done. Constructing database tables is pretty straightforward at that point.

4.1 MANY TABLE DESIGN STRATEGIES ARE POSSIBLE

Imagine that you want to build a relatively simple, but non-trivial application, for example, a to-do list. You can write down a set of functional requirements for the app and you can define a logical model (exactly as described in the previous chapters). You do all that without any actual software development first.

The next step is to design the physical table schema. But let's do a thought experiment: suppose that we hire a dozen teams of software developers, give them the requirements document, and ask them to implement it as they see fit.

They are free to choose their favorite database server and to design tables any way they like. The only restriction is that the requirements must be implemented fully, without changes, or additional features, and so on.

There would most probably be twelve different physical schemas. Some of them may look very similar to each other, but a couple would probably look really exotic.

Physical schemas would be different because different teams would use different **table design strategies**. A table design strategy is a set of rules that prescribes how to store anchors, how to store attributes, and how to store links.

TABLE-PER-ANCHOR TABLE DESIGN STRATEGY

Table-per-anchor strategy is probably the **most straightforward way** of building tables:

- one table per anchor;
- one table per many-to-many links;
- each attribute is a column in its anchor's table;
- each only-one-to-many link is a column of a table for the "many" anchor.

This approach corresponds to the third and fourth normal forms as defined in the relational theory. It is roughly what they would teach you in the beginner level university courses on database design.

The resulting tables are clean, understandable, logical and error-proof. From the physical point of view, this design is neutral, it is not optimized for any specific access pattern, except for understandability.

5.1 ACTION PLAN

Creating tables from the logical model is straightforward. Here are the next actions:

- Step 1. Insert table names in the list of anchors;
- Step 2. For each attribute, add the column name and choose the data type;
- Step 3. For each many-to-many link, choose the name of the table;
- Step 4. For each only-one-to-many link, choose the name of the column;

OTHER TABLE DESIGN STRATEGIES

6.1 TABLE DESIGN CONCERNS

There are four main concerns that underlie every table design strategy.

First, the result tables must be able **to store all the data defined in the logical schema**. This requirement may sound trivial, but sometimes people cut corners or get confused. As a result, database schema may just be incomplete. One way to make an incomplete database schema is to skip defining a logical schema.

Second, you need to understand **the effort needed to add a new attribute, link or anchor** to the database. Note that here we only talk about adding a new element, and not about other ways of database evolution (changing, deleting, etc.). Adding a new element is probably the most common operation, and it is expected to be roughly the easiest.

Third, we need to consider how the database would **handle the commonly expected read-only queries, performance-wise**. Queries could be classified in three groups:

- Naturally performant. For example, retrieving some attributes by a set of IDs is virtually always the cheapest possible operation, because it is naturally optimized in almost all databases.
- **Easily optimizable.** For example, if you need to query by a value of some attribute, you can easily add an index on the corresponding table column.
- **Specifically optimized.** Some queries are so important for the system performance that you need to optimize the database specifically for them. You need to be careful with

SECONDARY DATA

Anchors, links and attributes together are considered primary data. In other words, they are the source of truth.

For example, when we store the information about an order made by a customer, it is stored in anchor tables, attribute columns and link tables. If we lose this information, for example a table row is accidentally deleted, the system would think that this is how it is. If we delete one item from the order, we'll send a package without this item. If we delete information about the order, the system won't recognize the order number. This is primary data. If it gets deleted accidentally, you need to restore from backups.

There is another sort of data: secondary data. This is the data that was copied from the primary data, duplicated, reorganized, flattened, post-processed and so on. There are many ways to introduce and use the secondary data. If it gets deleted accidentally, it can be recalculated from the primary data (theoretically, at least).

Here is an incomplete list of things that we consider secondary data:

- cached (pre-computed) columns;
- pre-aggregated tables;
- denormalized columns and tables;
- ML models generated and updated based on data;
- materialized feeds;
- flat tables;
- full-text indexes;
- copies of data in different databases;
- data caches such as Memcached or in-memory Redis;

EVOLVING YOUR DATABASE

Inevitably, sooner or later in your project's lifecycle, you will have to change the structure of your database.

There are several practical reasons for the database to evolve:

- normal development of **new features**;
- changing the definition of existing database elements: anchors, attributes and links. This happens when the **business requirements change** in such a way that the original implementation is incompatible with them.
- improving database performance or scalability by **reorganizing physical schema**;
- introducing secondary representations, such as **derived columns and pre-aggregated tables**;
- **improving organizational scalability**: for example, making it easier to add new attributes by introducing a JSON column;
- **removing data** that is no longer needed, or no longer wanted;
- **migrating to a different database** for performance reasons, compliance reasons or for organizational scalability;
- **building data warehouses**: derived representations of primary data for better reporting, data analytics and so on.

Any database evolution requires changing three things:

- logical schema;
- application code;
- physical database itself.

MOVIE TICKETS: REPEATED SALES PATTERN

Let's discuss how to model selling movie tickets. This use case illustrates a very common pattern that I call repeated sales pattern. We will only focus on the logical model. We're only going to talk about essential parts of this pattern, omitting some details.

"Repeated sales" pattern is obviously important to model many businesses: not only movie tickets, but also flight and train tickets, hotel reservations, cargo transportation and other examples. Outside of sales, this pattern arises in a well-known textbook case of "course assignments" that is often used in teaching database modeling.

9.1 **BUSINESS REQUIREMENTS**

Suppose that our business is a movie theater, with multiple screens. Inside every auditorium there are a number of seats, physical chairs. We want to sell tickets to people willing to occupy those chairs for a few hours at certain times of a specific day.

Our movie theater also buys rights to show various movies, as they are released.

Every day there are multiple showtimes at each auditorium. Anyone can buy a ticket for a certain time of the day at a certain date, when a certain movie will be shown.

For simplicity, we assume that every movie ticket costs exactly the same, no matter what time slot, movie, or a screen. It's not hard to add this improvement on top of the basic logical model we'll build.

Moreover, we won't even bother with payments. Generally speaking, there are several ways how you can get a ticket:

BOOKS AND WASHING MACHINES: POLYMORPHIC DATA PATTERN

10.1 BUSINESS REQUIREMENTS

Imagine an e-commerce website that sells a lot of different things. For example, books and washing machines, but also clothes, bicycles, LED lamps, vinyl records, groceries, gift cards, etc., etc.

For each type of item we want to help users find the best item that we offer. For example, people shopping for a washing machine may want to search by the vendor name, capacity, physical dimensions, and for some unique features that the washing machine industry has invented.

People buying books want to search by the author name, cover type, publication date, etc. This is an entirely different set of properties. (Note that we do not say "different set of attributes", because author information requires links too. This is something that is often omitted in discussions of this pattern.)

Same for the bicycles, groceries and so on.

But there are some things that are common between books, washing machines and so on: they can all be ordered and they have a price. You can place a single order for two items: a book and a washing machine, and this order will be delivered to you, maybe in two shipments.

Let's apply our approach to this sort of business and see how it works. Here, in addition to the logical model, we will also extensively discuss the physical table design.

PRACTICALITIES

You can experiment with the approach presented in this book using any tool that you're familiar with. You need to maintain three tables: list of anchors, list of attributes and list of links. Also, you need a place to keep a textual description of your project.

You also need to maintain the physical database schema somewhere. It depends on the way you develop your system. It may be an SQL file in your source repository, or an ORM definition, or maybe you directly design the tables in your database using some sort of UI, graphical or text-based. It's possible that you're only interested in the logical schema so you skip this part.

11.0.1 Document-based catalog

It seems that the most natural choice here is some word processor, such as Google Docs (that's what I use). Tables in this book are formatted for readability, due to demands of e-book readers. For the actual table design process you need something simple.

On the following page you can see some example Google Docs screenshots (before typesetting).

11.0.2 Spreadsheet-based catalog

Another type of tool that you could use is a spreadsheet, for example Google Sheets.

Create three tabs: anchors, attributes and links. Input the column names, and make the header "frozen", so that it stays on the screen when you scroll the rest of the document. Enable text wrapping on all the columns, and set text alignment to top.