

Degree Name
Project Module Code
ID Number

**Skeem - a database system to deal with
modern needs**

by

Henry Morgan

Supervisor: Dr. C. Hatter

Department of Computer Science
Loughborough University

April/June 2019

Abstract

Skeem is a database

Contents

1	Introduction	6
1.1	Motivation	6
2	Background Information	6
3	Problem Definition	7
3.1	Api Repetition	8
3.2	Storage vs Display	8
3.3	Tight Coupling	9
3.4	Authentication	10
3.5	Lots Of Boilerplate	10
3.6	Bespoke Knowledge	10
3.7	Data Consistency	10
4	Literature Review	11
5	Requirements	11
5.1	High Level Requirements	11
5.2	Specification Gathering	12
5.3	Technical Requirements	13
6	Solution	15
6.1	Requests	16
6.1.1	Fetching Data	17
6.1.2	Mutating Data	19
6.2	Permissions and Authentication	19
6.2.1	Authentication	19
6.2.2	Roles	20
6.2.3	Permissions	20
6.3	Management	21
6.3.1	Cli	21
6.3.2	Gui	21
6.4	The Client	21
6.4.1	Caching	21
6.4.2	Validations	21
6.4.3	Oauth Flows	21
6.5	Live Updates	21

6.6	Plugins	21
6.6.1	Custom Attributes	22
6.6.2	Custom Operation Functions	22
6.7	Documentation	27
7	Methodology	27
7.1	Development Strategy	27
8	Usage	28
9	Implementation	28
9.1	Technologies	28
9.1.1	Typescript	28
9.1.2	Postgres	29
9.1.3	Other Libraries and Services	29
9.2	Project Structure	30
9.2.1	Skeem-Server	30
9.2.2	Skeem-Cil	31
9.2.3	Skeem-Gui	31
9.2.4	Skeem-Client	31
9.2.5	Skeem-Common	31
9.2.6	Es-Qu-El	31
9.2.7	Typer	32
9.2.8	Overseer	32
9.3	Database Usage	32
9.4	The Schema	32
9.5	Migrations	33
9.5.1	Running Migrations	33
9.5.2	Database Diffing	33
9.5.3	Rolling Back	34
9.5.4	Storing and Syncing	35
9.6	Requests	35
9.7	Fetches	35
9.7.1	Operators	36
9.8	Mutations	36
9.9	Sessions and Authentication	36
9.10	File Management	36
9.11	Plugins	36

9.12 Attributes	39
9.13 Operator	39
9.14 Session Provider	39
9.15 File Handlers	39
9.16 Cli	39
9.16.1 Overseer	39
9.17 File Management	39
9.18 Configuration	41
9.19 Documentation	41
10 Testing	41
11 Deployments	42
12 Conclusion	43
12.1 Future Work	43
References	44

1 Introduction

This project focuses on aiding developers building web applications.

1.1 Motivation

I have been a web developer for a number of years have built many applications across a wide range of systems. Over this time I have been involved with the development of many systems spanning a wide range of companies and have refined the techniques needed to create blah blah blah...

2 Background Information

creating a website

- some stats about websites

modern websites must be dynamic and interactive - this is achieved by SPAs

SPAs or single-page applications are a growing trend on the web unlike traditional architectures where each page is a separate resource with its own end-point, its own template and its own request, SPAs combine all the pages into a single end point. When a user navigates to a page they download a large javascript bundle which has the ability to construct any page of the website. The javascript is reads the url the user requested and renders the creates the appropriate views. This means that when a user navigates to another page, the javascript can intercept this and simply render new content without a network request creating a much more responsive interaction.

Note: This is a simplified example of how SPAs tend to work. In reality having a user download the entire code bundle would not be ideal, especially on slower networks. Therefore optimizations are performed such as code splitting where the user only downloads the necessary code to build the requested page and any assets needed to display loading screens, then during idle time or upon navigation download any missing code required to render new pages.

Websites can be decomposed into many separate yet highly dependant parts:

Database a service which allows for the efficient storage and retrieval of large amounts of data.

View a system responsible for which templates can be defined and then populated depending on the specifics of a given request.

API provides an interface between view layer logic and the database

Authentication systems provides methods for confirming a users identity and tracking who is making asking for data in between stateless http requests.

File Hosting services which manages the storage and retrieval of large form data such as images. This

This project will attempt to replace the database, api, and authentication and file hosting subsystems with a single unified system.

3 Problem Definition

This chapter will provide an overview of the problems associated with building web applications. It will offer a high level overview of existing practices, how they currently function, why they exist and tje problems that exist which require solving.

There are many issues with this method of building websites.

3.1 Api Repetition

Consistent and discoverable APIs tend to lead to very repetitive code. If, for example, you need an end point to fetch a list of blogs and you also need one to fetch all the products. What is really different about these routes? The table name in the SQL query and the columns it returns. This duplication of code leads to highly duplicated code, something which is almost always undesirable (???)(<http://www.informit.com/articles/article.aspx?p=457502&seqNum=5>).

There is also a trade off to be made between receiving unnecessary data and have many very similar endpoints. For instance, imagine having an end point requesting a list of blog posts. On the site you wist to display the title of each blog along with a some preview text and the authors name. On the same website you also have a page for viewing an individual blog. This page contains a blogs title, body, its authors name, and a posted date. Here, both pages request very similar data, only differing in one pages need for the posted date. The options here are to either have two very similar end points which return near identical data, or forcing the end user to download the posted at information when they may never actually need to render it. The former leads to having to maintain two distinct APIs where as the latter uses unnecessary data. Neither option is desirable.

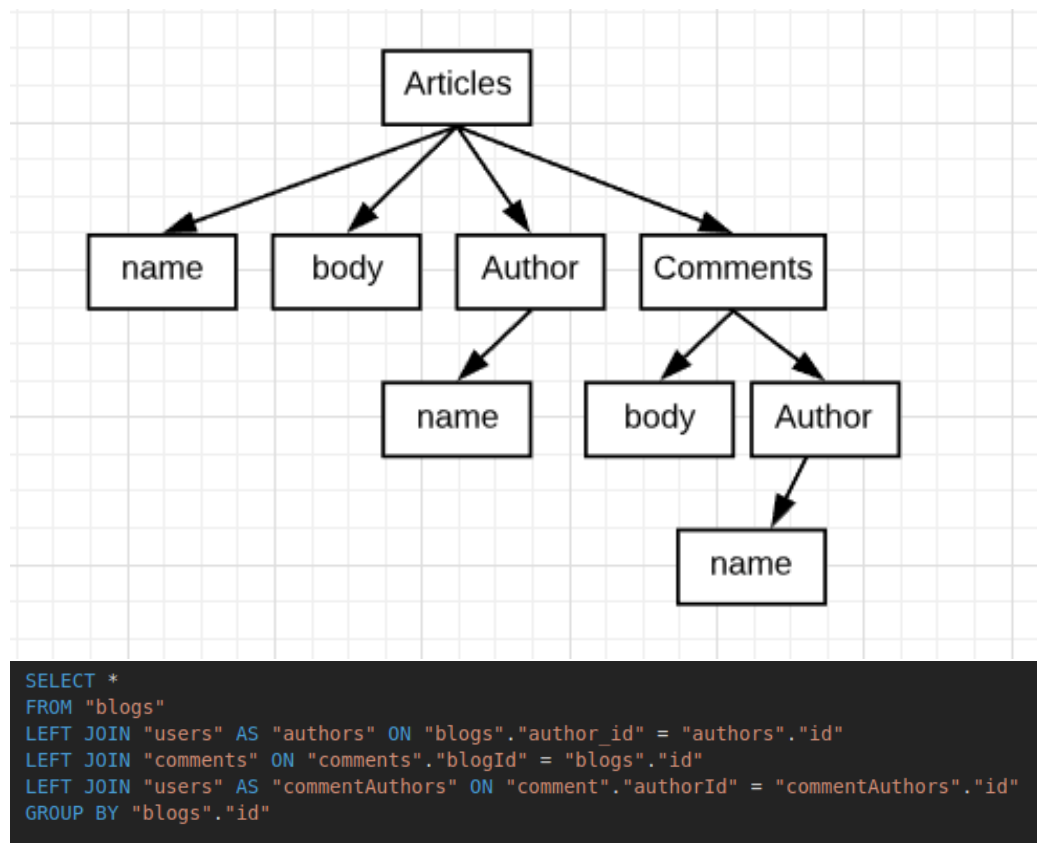
3.2 Storage vs Display

Databases should store normalizes data which, simply speaking, means storing data in flat tables i.e one for articles, one for authors, one for comments and then you storing relational information between different records. Storing data in the fassion is desirable as it greatly removes data duplication which makes updating records easier and removes the risk of desyncing data by ony updating it in a single location.

The issue with this, however, is that data is not displayed like this to the end user. The end user is not presented with a page containing an article and is required to navigate to a separate page displaying the authors name then have to navigate to a third place to read a list of comments. Rather the end user will be presented with a single page containing all the data amalgamated in an easily digestible and pleasant format. The data the user sees can be envisioned as a tree of data: the root being the article itself and

then containing a connected nodes for each comment each having further nodes containing their authors.

The need to request a tree of data from a database is an extremely common and useful thing, however, despite being conceptually simple it can get incredibly complex even when having to traverse only a few levels deep.



3.3 Tight Coupling

APIs tend to be closely related to the underlying database storing the data. If you have a properties table, then you will want a properties API to access the data. If you were to change the name of a column within the database then you would have to remember to update the API to match and this problem grows if you have tables which span multiple APIs.

There is a similar relationship between the client and the API. When an attribute is changed on the API then everywhere using that attribute is required to update simultaneously else risk displaying incorrect responses or worse, completely crash if vital data is changed.

Tight coupling and disparate implementations leads allows for the opportunity for a de-sync which will inevitably lead to bugs.

3.4 Authentication

Authentication is a very simple concept in the abstract but is very easy to get wrong.

3.5 Lots Of Boilerplate

Another issue is that there is a lot of boiler plate

- Slow to set all these systems up

3.6 Bespoke Knowledge

- Lots of bespoke knowledge from many domains means it is hard to train people

3.7 Data Consistency

once you have the data it is important to keep it from being stale.

- Once you retrieve the data from a data source (server) it is important to keep that up to date
- if it gets updated on the server it should be reflected on the client

Although this is not a problem exclusive to SPAs and instead more broadly to websites in general, it is emphasised by SPAs as since they do not need to reload between upon navigation they can potentially cache data more

aggressively further improving performance and responsiveness. The risk of stale data, however, greatly limits this potential.

4 Literature Review

In this chapter I shall discuss existing technologies which solve the defined problem. I shall briefly describe how each solution functions as well as their advantages and limitations.

- GraphQL + Relay
 - Sits on top of an existing
 - Provides a way to query a tree structure of data
 - complicated - non feasible for non technical people
- NoSQL databases
 - Stores arbitrarily shaped json allowing data to be stored in a fashion similar to its usage
 - eliminates the tree problem
- Web framework
 - Rails

5 Requirements

5.1 High Level Requirements

- what would need to be achieved to solve all these issues
- Tree Structures
 - Fetch arbitrary tree structures
 - Query interface which is easily sanitizable such that it could be executed from even dangerous clients without risk of returning non-permitted data.
- Simple:

- Explainable through limited, reasonably sized, help docs
- Simple GUI interface
- Requires 0 knowledge of database structures to use associations
- Includes File management
- User authentication
- Consistency
 - alert clients to changes in data

5.2 Specification Gathering

In order to create a solution which will alleviate these issues I had to ensure that the system achieved everything needed to replace existing systems rather than just add a further system which must be configured, maintained and learnt.

accessed an in-production data base and pulled a list of all interactions with the database

- used by 4'000 unique visitors a day
- 4% are new visitors
- 20'000 registered users

I went through all interactions with the database and records how it was being used:

- attributes
 - has many through
 - has many through where condition
 - has many with condition
 - has many dependent nullify
- Validations
 - presence
 - uniqueness
 - inclusion
 - number greater than
 - uniqueness in scope of attr: value
 - validate uniqueness in scope with condition unless attribute: value
 - Validates on: :create

- association.attribute must = value
 - validates [if/unless] attribute: value
- Callbacks
 - before_validation
 - * default attribute to another attribute if not present
 - * default attributes only on create
 - * default attributes to parameterized other attribute
 - * default attribute to association attribute
 - before_create
 - * self.slug = name.parameterize
 - after_create
 - * update association
 - * send emails
 - * update self
- Scopes
 - where(attr: value)
 - where.not(attr: value)
 - order(attr: :desc)
 - where association count >= 1
 - where association count === 0
 - where association attribute
 - where in associations scope e.g where(tag_id: Tag.published)
 - composing scopes (adding limits)
- Permissions
 - through user association
 - through user association | where(attr: value)

Using this information I obtained the minimum viable feature set needed

5.3 Technical Requirements

combining these two sets of requirements produces the following set:

Must be able to cope with any future requirements and not pigeonhole functionality.

- Create Models
 - store basic types strings, number

- store associations between two models
 - store files
- Fetching
 - attributes
 - * Request primitive attributes such as strings, numbers
 - * request associations
 - provide a filter to a query
 - * request a record given the records id
 - * request a record based on its attributes i.e requesting published records
 - sort queries
 - * by attributes
 - * by associations attributes
 - pagination queries
- Mutations
 - create records
 - Update records
 - delete records
 - add/remove association records
 - upload files
 - validate data
- Sessions
 - Authenticate users
 - Specify users permission to access data
- Consistency
 - Use web sockets to be alerted to updates
- Permissions
 - Specify access (read + write + remove) of users on:
 - * records
 - * attributes
- Provide a way to change production databases safely
- GUI

- Provide a way to create a database
- Provide a way to create a model
- add/update/remove attributes from models
- seed data
- view records for a model

6 Solution

In order to solve the problem I have created a system named Skeem. This chapter acts as a high level overview to Skeem's functionality, and how it is to be used.

Comment:

What is skeem

Send queries from the front end

wraps and maintains a database.

developed a custom query format

Skeem wraps a standard object relational database but augments its functionality by providing a new query interface which is easily sanitizable and so can be executed directly from the client side. This greatly helps to decouple the API and View layers and the view can directly query for its relevant information. This new query interface is designed around tree data structures and so meshes nicely with the needs of the view layer. Skeem fully manages this underlying database and the user is not required to understand how it is laid out or utilized nor do they require any knowledge of SQL to utilize Skeem.

Comment: *auto generate an API based upon the database thus removing some coupling The ability to query from the browser fully eliminates the need for a dedicated API.*

Comment: *built in authentication*

Skeem also provides a full authentication system capable of identifying and tracking users as well as being able to limit their access to resources as necessary.

Comment: *Skeem requires no code. no training required.*

Skeem can be fully setup, configured and maintained without the need to write any code. Instead it is managed through either the command line or through a graphical interface. These interfaces provide instant feedback of any errors that occur when changing things and also provide help information to aid Skeems usage. This helps to solve the issue of training. There are of course still intricacies with using various sub-systems which will require additional help, however, to solve this issue skeem contains a fully set of documentation detailing many aspects.

Comment: *runs a server constantly listening to http requests*

Skeem runs as a service and communication is done via http using a pure json API. This allows browsers to easily communicate with it as all ship with the ability to serialize and parse json data as well as send http requests.

Comment: *websockets for live updates*

Skeem exposes update events through web sockets to enable clients to automatically fetch data when it becomes stale. This allows data to be more aggressively cached and the interactivity of sites increased.

6.1 Requests

Skeems uses JSON to send requests and receive responses. This allows requests to be much more easily sanitized against malicious attempts to manipulate the database in an undesirable way. In other words it helps prevent SQL injection like attacks.

SQL injection is a code injection technique in which malicious SQL statements are inserted into normally safe SQL by having unsanitized input inserted into the query. SQL injection is a major security vulnerability (???)

Due to this sanitization ability queries are safe to execute on the front end and so eliminates the tight coupling between This solves the tight coupling with disparate code issues between the standard client-server model.

All requests are fired to the same end point with a post request and take the form of a JSON object with **type** and **payload**. The type is used to determine

what type of request it is: `fetch`, `mutate`, etc. and the payload contains specific information depending on the type.

Comment: TREE OF DATA

6.1.1 Fetching Data

Fetching data involves pulling data from the models in a structured fashion. A fetch query specifies a single root model name as the key to the query object. The value then specifies exactly what data you want to retrieve, how to filter and sort it and whether you want to split it into pages.

A fetch response will always be an array of records, where each record will contain its own id as well as any additional attributes you requested. In some cases you may also retrieve the total record count see the section [Pagination](#).

Attributes

Attributes specify what data you want to receive for each record. Attributes take the form of an array where each element is a string being the name of the attribute you are requesting, or an object with a `name` property and a value of the attribute. This object notation is required for specifying additional configuration such as formatting information like an alternative name for results.

Attributes can be aliased by specifying the `as` property

Listing 1: This query will return all articles each containing the article's id, name, and body. The body will be aliased under the name "text"

```
1 {  
2   articles: {  
3     attributes: ["name", { name: "body", as: "text"  
4   }]  
5 }
```

associations: trees of data....

Filter

Filters all for specifying specific criteria records must meet for them to be returned.

Filters you a tree of “object functions”. This means that each object within a filter operation contains a single key. This key specifies what function you wish to execute e.g equality, less than, empty check. And the value acts as the arguments for the given function.

There are many built in filter functions which cover a broad range of use cases

Sort

Sorting data is an extremely common and essential ability for data retrieval: Most recent tweets, video length, article title. When sorting data you specify what you want to sort **by** and the **direction** you want to sort: either *ascending* or *descending*.

Listing 2: This query will return all articles ordered by the articles "name" attribute.

```
1 {  
2   articles: {  
3     sort: {  
4       by: "name",  
5       direction: "asc"  
6     }  
7   }  
8 }
```

You can also specify an array of sorting criteria. Doing this will sort the data initially by the first item, then resolve conflicts with the next item in the list.

Pagination

Pagination chunks the data into pages. You don't want the end user to download 100'000 records, this would be very slow and wasteful. The returned data will be equivalent to a standard array of records.

Listing 3: This query will return the second page of articles where each page holds 30 records.

```
1 {  
2   articles: {  
3     pagination: {  
4       page: 2,  
5       perPage: 30,  
6     }  
7   }  
8 }
```

Record Count

As well as the actual records, you also get given a count of the number of records you would have gotten if you did not paginate the data. This is useful when wanting to show end users a list of page numbers and allow them to jump to them arbitrarily.

$$totalPages = ceiling(totalRecords/perPage)$$

Retrieving the record count can be disabled by passing the option of `withCount` : `false` to the pagination block.

6.1.2 Mutating Data

6.2 Permissions and Authentication

Controlling what users can access, who can create and edit data is an essential part of all application which make use of a database.

6.2.1 Authentication

Before we can control a users access we must first be able to determine who they are. Skeem provides a couple of ways in which to authenticate someone: stored identifier (email, username, etc) with a password or through an oauth2.0

provider such as Facebook or Google. These methods of authentication are referred to generally as session providers (they provide methods for authentication sessions).

Session providers define a name, model, a type and some configuration dependant on the type selected.

The name is used purely to distinguish between different providers and allows for multiple authentication strategies of the same type. You may have to distinct user sets which are authenticated with different models e.g for a school system you may have one user set named teachers and another for pupils. The model defines where the session provider should look to find the necessary data to check against any credentials provided.

The type defines which session provider to use. Skeem comes with three built in providers: local, facebook, and google.

The local provider authenticates users by storing some identifying attribute and a password in the database itself. Then when an authentication request is made the database is queried for a record with the specified credentials. If a record is found then skeem authenticates the user as that record. The user attribute is most commonly an email address or a username. The password is stored securely using a secure hashing algorithm.

The Facebook and Google providers allow users to authenticate using these services via the familiar “login with XXX” buttons. These providers can specify a list of attributes to extract from the service such as name, email, image.

6.2.2 Roles

With authentication we now have two distinct user states: authenticated and anonymous. These roles can be used to define permissions on fetches and mutates

6.2.3 Permissions

Each model can define a set of permissions based scopes

6.3 Management

There are two predominant methods to managing a skeem application: through the use of a command line interface or via a GUI.

Through the combination of logically named utilities, help text and descriptive error messages, people can use the system with minimal training. This helps to reduce the need for specialized knowledge

6.3.1 Cli

6.3.2 Gui

6.4 The Client

Skeem is not attempting to replace the view layer of websites. Because of this it has an implementation of a client which can be run on the view layer to wrap the requests into simple function calls as well as adding basic checks.

Skeem has a client, written in Javascript, designed to be used with SPAs. The client provides functions which will process queries and send them to the server.

6.4.1 Caching

6.4.2 Validations

6.4.3 Oauth Flows

6.5 Live Updates

6.6 Plugins

APIs are a large and complex systems which cover an incredibly broad range of use cases. It would be almost impossible to foresee every use case of Skeem

and to allow for every possibility. To cope with this, Skeem has the ability to augment functionality by the way of plugins.

Plugins are javascript files located in the project folder. When the server starts these files are loaded, type checked, and inserted into the system.

6.6.1 Custom Attributes

There are many different types of datum which you may want to use which don't fit within the bounds of the built in types.

The attributes built in to skeem use this very system i.e they contain no special functionality which could not be implemented outside of the core code.

6.6.2 Custom Operation Functions

There are a myriad of different and obscure filters you may want to perform within a database. Whilst skeem contains alot of built in operations which can achieve a large variety of results it is implausible that they cover every possible desire.

Therefore, like with attributes, skeem provides the ability to create custom operation function outside of the skeem source and have them loaded in dynamically and used seamlessly with the built-in operations.

Creating An Operation

An operation consists of a single function which must return SQL.

Listing 4: The simplest custom operation - it would always return false and so is utterly pointless.

```
1 function myPointlessOperation() {  
2   return { value: `this will` = 'always be false'`  
3   }  
}
```

In order to create a useful operation the function is passed some variables concerning the request. The most useful of which is the `value` argument. The

`value` contains the data passed to the operation. Using this we can produce a much more useful operation.

The SQL returned is inserted into the query as is and so is essential that it is sanitized prior to being returned.

Listing 5: Returns all the records for which the name matches the value supplied. However there are major issues with this and should not be used.

```
1 function myBadOperation({ value }) {  
2   return { value: `"name" = '${value}'` }  
3 }
```

The above operation shown in figure (FIGURE XXX) would technically do something which could be deemed as useful. You supply a value and it will return all the results for which their name attribute is equal to that value. There are a couple of issues with this operation though which means it should not be used.

The biggest and most critical issue is that it does not sanitize the value. This means it is an entry point of an SQL injection. This is relatively easy to solve though through the use of a sister package of skeem named `es-qu-el`. That is outside the scope of this solution chapter and is discussed in the implementation.

The second issue with the operation is that it makes the assumption that the model you are currently fetching has a column in its table called “name”. This is not always true for obvious reasons. To solve this issue we are passed another useful piece of information - the current model. With this we can search through the models attributes and check for the existence of a name attribute and if it does exist then throw an error.

This is an extremely common need for skeem and as such there exists a helper function to achieve this for you. It is exported from the `skeem-common` package and is called `getAttribute`. This function takes a model and the name of the attribute you wish to retrieve. By using this function you guarantee that the attribute exists and if it does not exist then you can be assured that you will get an error message consistent with that of a built in error.

Listing 6: Checks to see if name actually exists on the model being queried.

```

1  const { getAttribute } = require("skeem-common")
2
3  function mySlightlyBetterOperation({ model, value })
4      {
5      const attribute = getAttribute(model, "name")
6      return { value: `${attribute.name} = '${value}'`
7      }
8  }

```

There is one final issue with this operation and that is it misses the opportunity for optimization. As well as returning the SQL value operations can also return the type of result expected back from the SQL - in this case it would be a “boolean”. Supplying this information from an operation allows Skeem to optimize the SQL query and possibly not even execute anything. For instance consider the following query:

```

1  {
2    filter: {
3      eq: [{ value: "a string" }, { value: 123 }]
4    }
5  }

```

The `value` operation will return the types of string and number (as well as the sanitized SQL value). The `eq` operation then checks these types to see if they are save it is then it will place the values around an equals sign and return it as expected. If, however, they are different then the `eq` operation will return with a type of “false”. If the full filter resolves with the type of “false” skeem will skip executing the query as it knows nothing would be returned. Therefore by returning the correct type skeem can potentially optimize and avoid the database altogether.

Possible types include:

- string
- boolean
- number
- record
- collection
- any - the type could not be determined and so could be anything, this is the default when no type is returned

Listing 7: Checks to see if the attribute actually exists and returns the correct type

```
1  const { getAttribute } = require("skeem-common")
2
3  function myPassableOperation({ model, value }) {
4    const attribute = getAttribute(model, "name")
5    return { value: ` "${attribute.name}" = '${value}'
6      `, type: "boolean" }
7  }
```

It should also be noted that this is a poor use of a custom operation. This operation does not achieve anything that you could not do with the built in operations. Although it does save a few characters, it adds on additional knowledge needed by other people working on a project, the need for testing, and another function which would need maintaining over the life time of the project. Operations should only be added to achieve results which are either not possible with the pre-existing functions or highly impracticable to achieve.

Nested Operations

It is very common for an operation to need to accept an operation object as its value. If you could not compile nested operations then you would not be able to create functions like: `eq`, `lt`, `and`, `not`. This would make things a little tricky. Therefore, along with `model` and `query` you also get supplied with a function named `compile`. This function accepts an operation and returns the `{value, type}` object.

Listing 8: A simple implementation of the `eq` operation.

```
1  const { getAttribute } = require("skeem-common")
2
3  function simpleEq({ compile, value }) {
4    const left = compile(value[0])
5    const right = compile(value[1])
6    return { value: ` "${left.value}" = "${right.value}"`,
7      type: "boolean" }
8  }
```

The Request Context

the final argument passed to an operation is `ctx`. This is the current context for the request. With this it is possible to access information such as app configuration, the current session, and the database connection.

TODO

Using An Operation Plugin

To use an operation you first must create a javascript file within the `<skeem root>/operations`. The name of this file does not matter and can be anything. This javascript file must have a default export of an object where the keys are the names of the operations and their values are the functions as described above.

Listing 9: A full operation plugin file

```
1
2 module.exports = {
3   isANumber: function({ value }) {
4     if (typeof value === 'number') {
5       return { value: true, type: 'boolean' }
6     } else {
7       return { value: false, type: 'boolean' }
8     }
9   }
10 }
11
12 // Usage:
13
14 fetch: {
15   articles: {
16     filter: {
17       { isANumber: 123 }
18     }
19   }
20 }
```

6.7 Documentation

Documetnation exists for skeem which contains guides and examples on how to use the system. This is essential to allow other developers to actually use the system. <https://cd2.github.io/Skeem/#/>

7 Methodology

7.1 Development Strategy

When a new feature would be added it would first have high level tests written aimed to test the final functionality of the feature. For instance when first implementing fetching I wrote tests asserting that given a particular query a specific piece of sql was generated. I would then proceed to implement the feature, using the tests as a guide for when the work was complete.

Listing 10: Example of what the high level tests would assert (not actual tests)

```
1 Given: { articles: {} }
2 Expected: `SELECT "id" FROM "articles"`
3
4 Given: { articles: { attributes: ["name"], filter: {
    eq: [ { attr: 'name' }, {value: 'test'}] } } }
5 Expected: `SELECT "id", "name" FROM "articles" WHERE
    "attr" = 'test'`
```

Once all the tests were passing, if there were additional features which either appeared during implementation or that were initially excluded for simplicity, I would add more high level tests asserting the new functionality. I would then proceed to implement these features until the new tests were passing, adding more tests until development was complete.

When the feature was complete, assert by a suite of passing tests I would begin testing the code at a more granular level. I would select functions which were either complicated or hard to test at a high level (maybe code branches for very specific circumstances) and write specific unit tests.

The specifics of how the tests are written are discussed more in the chapter (Testing)[#testing].

8 Usage

9 Implementation

In this chapter I shall be explaining some of the intricacies of skeem, how it works, and why certain key decisions have been made. This chapter will not cover the implementation of everything skeem offers but instead covers the key aspects for which there were challenges and interesting solutions.

9.1 Technologies

9.1.1 Typescript

Some parts of Skeem run on the client side and thus had to be built in javascript.

Typescript is a superset of javascript which adds typing capabilities to the normally dynamic language (???). This typing information allows compile time code checking which greatly assists improving reliability of code as it helps to ensure against common trivial bugs. Typescript is transpiled into standard, es5 compliant, javascript meaning it is able to run on all modern browsers. This compatibility is essential because if skeem only supported the latest version of chrome then it would instantly eliminate the possibility of using the system on any standard website which has compatibility as a requirement.

Typescript also aids other developers using the system as the typing information is used in most modern IDEs to supply intelligence information allow features such as auto completion, inline error messages, and hints for expected arguments and returns.

NodeJs is a javascript runtime designed to build scalable network applications (???). NodeJs was a logical choice as it allowed writing server code also in

javascript which allowed consistent interfaces to be constructed. It is also much easier to develop a system when writing in a single language as there is less mental energy exerted to convert from one environment to another.

9.1.2 Postgres

Postgres is a object-relational database system (???). Postgres has powerful, inbuilt JSON processing capabilities. It allows for storing JSON objects natively as well as writing queries which inspect the contents of JSON. More importantly Postgres allows for the construction of JSON objects with queries themselves. This makes Postgres a very logical choice when the goal is to create trees of data as the database can pre-format the response greatly reducing the need for much post query processing.

Another feature which postgres offers is an optimizer which can automatically transform subqueries into join statements. This optimization is greatly taken advantage of in skeem as throughout the code base there is not a single joining statement. Further reasons for this are discussed in the chapter on fetch query sql generation.

Postgres also provides the ability to write custom database functions relatively easily. These functions help to encapsulate complex pieces of reusable logic and are optimized by postgres in order to maintain performant queries. They are used in a number of places throughout skeem such as to trigger update messages for live syncing and to format responses in certain circumstances.

9.1.3 Other Libraries and Services

Skeem also takes advantage of a number of prebuilt libraries and services.

NPM is the defacto package manager for node (???), it provides easy hosting and distribution of node packages and is the method skeem uses to manage its publications.

Pusher is a web service specialised in providing real-time functionality to applications (???). It provides simple wrappers around web sockets as well as fallbacks to ensure compatibility across, even outdated, browsers. Skeem uses Pusher send messages to clients in order to enable the live updating capabilities.

Amazon Web Services are a cloud computing platform(???) which provides relatively affordable file hosting and is integrated into skeems file storage capabilities.

React React

TODO

9.2 Project Structure

Due to the size of the skeem code base (pushing 12'000 lines) and the range of environments it runs on (server and the client). It was essential to split up the project into logical parts. However I did not want to fully isolate each section due to the tight coupling of the interactions e.g if the format of the fetch response changes then you have to update the server code to handle this new change as well as the client code to keep the format consistent from a developers point of view. Due to this I structured the code as a monorepo.

“A monorepo is a software development strategy where code for many projects are stored in the same repository” (“Monorepo” 2019).

Then using a tool named Lerna I was able to manage the projects simultaneously. Lerna automatically resolves the dependency order so, for instance, when you attempt to build the project it knows that A depends on B which depends on C and therefore wait for A to build before moving on to B then finishing with C.

Skeem is broken up into 7 packages, 5 directly related packages, each prefixed with “skeem-”, specific to skeem and 2 auxilary packages which were extracted and can provide useful functionallity independently:

9.2.1 Skeem-Server

This package contains the majority of the logic, it contains the implementation for processing requests, creating a database, session authentication, migration creating and validation, etc...

The `skeem-server` package is split into two, confusingly named, parts: a manager and a server. The manager contains the all the functionallity of the

app, it controls loading the schema, producing and executing SQL, updating config, etc. The server simply listens for http requests sent from the client, performs basic format type checking and then calls manager functions in order to create a response.

Server

Manager

9.2.2 Skeem-Cil

This implements a command line interface for interacting with skeem, it implements no fundamental logic and instead acts a wrapper around the server.

9.2.3 Skeem-Gui

Similar to the CLI, the gui acts a wrapper around the server functionality and displays the information in a visual application.

9.2.4 Skeem-Client

Provides functionallity for a client

9.2.5 Skeem-Common

Holds common functionallity needed between packages, such as error messages

9.2.6 Es-Qu-El

provides helper functions to generate and sanitize SQL statements

9.2.7 Typer

Typer standardizes type checking across the app.

9.2.8 Overseer

declarative CLI generator used to power Skeems CLI.

9.3 Database Usage

- Skeem wraps a Database
 - why? performance
- Built in tables
- Upgrading

9.4 The Schema

- what is it
 - models
- attributes
 - an attribute in the abstract simply comprises of a set of simple methods get, set, migrate.
 - * these are then
 - simple types such as string, date, number, boolean
 - skeem offers more complex types
 - * These types are where skeem really shines
 - * association
 - * file
 - see the chapter on file manager for more specifics as to how the files get processed and stored

9.5 Migrations

- Migrations
 - what do they do?
 - * the manipulate the Schema
 - they dont manipulate the database directly because then they can be run in order to assert there are no errors

Migrations are the way the schema is mutated -they provide a simple and repeatable method for manipulation.

Migrations do not directly effect the database instead they modify the schema, and then the differences between the old and the new are calculated and a set of transformation steps are generated. The reason for this split is so that migrations can be executed without affecting the database which allows any errors to be detected without having to...

Migrations are comprised of a type and some data.

- Migrations work to mutate the schema
- They allow multiple developers to work on a single application
- A record of all changes made to the database
- Migrations serve incremental, reversible changes to the schema
- Migrations are stored in files within a folder named migrations located in the root of a skeem project
- This allows for migrations to be transferred between computers
- only stored in files to transfer computers, the ones which get executed are actually stored in the database

9.5.1 Running Migrations

9.5.2 Database Diffing

After a new schema has been produced for an application, a list of change steps must be realised in order to mutate the existing database. This happens

are migrations are run or when the application is initialized. In the latter case the empty schema is used as the old schema.

The first step of this process is to compare the new schema and the old one in order to find what is actually different. Because the models and providers only exist in the abstract, as opposed to the db property which is backed by database tables, only the db of the schema is diffed.

The first step in diffing the dbs is to isolate which tables are new, which have been removed and which have been **potentially** updated. The name field of the table is used to link the old and the new schema. If a name exists in the old schema but not in the new, then the table is marked as deleted. Similarly, if a name exists in the new schema but not the old it is marked as created. If the name exists in both then the table is marked as potentially updated and undergoes a further diff.

For each new table the appropriate `CREATE TABLE` SQL query is generated and appended to a list of all pending database queries. Like wise for each removed table a `DROP TABLE` query is produced and appended to the list.

For each updated table each column

The list of sql commands is then executed. To do this, first, a new transaction is created within the database. This means if an error occurs within the mutate steps the database can be fully restored to prior to the mutations. Without this then the schema could become out of sync with the database. The list of sql statements is the run sequentially. After each step has been executed the schema in the database is replaced by the new schema. Finally, a commit message is sent to Postgres informing it to proceed with the mutations. By updating the schema within the same transaction ensures synchronicity between the tables and the schema.

9.5.3 Rolling Back

Migrations all have the ability to be reversed.

When a migration is executed it has the option of returning some data. This data will be passed back to the migration when rolling back. This is used, for example, when creating an association the migration will return: the name of the added attribute and the name of the created joining table. With this

information the migration is then able to fully undo any effects it had thus restoring the previous schema.

Rolling back is not guaranteed to revert the database completely to its state prior to the migration, instead it simply guarantees the structure of the database will be identical. This is because some migrations such as deleting an attribute are lossy in nature and thus are not purely reversible.

9.5.4 Storing and Syncing

- migrations are stored in the database along with additional information such as whether they have been executed and what data they returned after being executed.
- They are stored in files to allow them to be transferred between systems to be run on a different machine
- also gets committed to version control software allowing multiple developers to be working on the same project

9.6 Requests

- The server hosts
- listens by http
 - so the client can communicate
- request and response formats
 - why pure json
- authentication token

9.7 Fetches

- they use the Schema
- build sql
- execute a single query
- what does the sql look like

9.7.1 Operators

9.8 Mutations

- what are they
- structure of the query
- how sql is generated

9.9 Sessions and Authentication

- Providers
 - What are they
- Local Provider
 - what it is
 - how it works
- Oauth Providers
 - built in ones
- once authenticated a session token is generated which is used to identify. this is encrypted with a secret key using the JSON web token standard

9.10 File Management

- What are file providers
- How do files get stored
- How are files set

9.11 Plugins

Plugins allow skeem to cover a much wider use case than would likely be possible if I had to implement all edge cases manually. By forcing system features to use a plugin style system it helps to better define the edges of systems and to creating cleaner more targeted sub-systems. This helps to prevent the core of skeem to become bloated with extremely specific features (how).

Note: There are examples of very specific features within the code base such as the `associationEquals` operation. It was, in fact, the addition of features such as this which prompted the need for the plugin system.

Every pluggable feature (attributes, session providers, operators and file handlers) all work off the same abstracted code and each one simply specifies a set of configuration. The abstracted code handles the file loading and parsing and exposes key functions to the specific implementations.

Each pluggable feature specifies the following:

- name** this is used when logging debug information.
- builtinsFolder** this specify the name of the folder where the built in features can be located.
- folderName** the folder name of where the external plugins will be located relative to the application root.
- validateExports** a function which will get passed a file when it has been imported and parsed. This function should return a boolean indicating whether the exported contents of the file is valid. For instance checking whether a file handler contains a store and retrieve method.
- merge** a function responsible for merging all the individual files together into one final object.

Abstracting the plugin system like this helped to enforce a consistent plugin style (which can aid developers when writing plugins), reduce duplicated complexity, and to isolate a critical system feature to be tested separately.

Listing 11: The code required to define the attribute plugin system

```
1 export const loadAttributes = createPluginLoader<
  IAttribute>({
2   name: "attributes",
3   builtinsFolder: "./builtins",
4   folderName: CUSTOM_ATTRIBUTES_PATH,
5   validateExports(_filename, exps) {
6     // attributes should be an object with three
        keys: "migrate", "get", "set"
7
8     if (!isObject(exps)) {
9       return false
```

```

10     }
11
12     const expectedKeys = ["migrate", "get", "set"]
13     const isValid = Object.keys(exps).every(key =>
        expectedKeys.includes(key))
14     return isValid
15 },
16 merge(acc, filename, exps) {
17     acc[filename] = exps
18     return acc
19 }
20 })

```

Loading a plugin executes the following algorithm:

1. combine the `builtinsFolder` variable and the directory the function was defined in, in order to find the full path for the built ins.
2. if this path doesn't exist, then skip steps 2-6.
3. load a list of all files in this directory.
4. for each javascript file require it
5. pass the contents to the `validateExports` method along with the file name. If this function returns false, throw an error
6. otherwise, store the contents along with the file name in an array.
7. combine the `folderName` variable and skeem's root directory to gain a full path for where external elements are stored.
8. repeat steps 2-6 using this new path.
9. for each element in the loaded array, pass it to the merge function.
10. return the result of the merge.

When skeem is started it calls all the loaded functions created by the plugin system. These functions will then load all the built elements and any external ones, merge them, and then store them with the manager. Then when a component wants to use one of the plugin elements it references it through the manager.

9.12 Attributes

9.13 Operator

9.14 Session Provider

- just talk about them in the context of wanting to implement twitter

9.15 File Handlers

9.16 Cli

- How its structured
- what commands are available
- The cli interacts with skeem by means of a text based command line. The was split into an actual definition of the interface and the parsing and processing of the command line text.
- Command Definition:
 - Commands are defined declaratively as an object containing: a name, some help text, a list of accepted options and arguments (each with a name, a type and an optional default) and a function to call when the command is executed. Skeem currently contains 14 commands.

9.16.1 Overseer

9.17 File Management

- How its built
 - react

```
henry@henry-OptiPlex-3060:~$ skeem -h
Commands:
  init <name>
    initializes a skeem project in the current directory
  start
    starts the server
  seedCreate
    create a seeding task
  seedRun <name>
    run all seed files
  migrationsCreate <type> <data> [--check=true]
    creates a new migration
  migrationsRun [<n>=all] [--reset=false]
    runs migrations
  migrationsList
    Lists all migrations
  migrationsSync
    sync migrations files with the database migrations
  migrationsRollback [<n>=all]
    rolls back migrations
  schema <name>
    show the current schema
  fetch <query>
    perform a fetch query
  mutate <query>
    perform a mutate query
  upgrade
    upgrades the database to the latest version
  gui [--open=false]
    starts the gui in the current folder
henry@henry-OptiPlex-3060:~$
```

Figure 1: CLI Help screen

9.18 Configuration

- this is how skeem identifies the project root, and allows you to execute commands even when within subfolders

9.19 Documentation

- why it is important
 - many parts of skeem seem very complicated
 - lots of stuff can be hard to know exactly what to do
 - contains examples
- docsify
- built into the repo and automatically deploys

10 Testing

Tests are an essential part of any software project especially those providing some critical functionality to users - Skeem is no exception.

Tests were written in Jest.

- Tests were written in jest.
- Tests targeted functionality rather than implementation. However tests were written for smaller parts of the system when functionality became too complex or when the underlying functions were critical - such as loading the schema from the database.
- For a complete list of all tests please see the appendix.
- CI
 - Due to skeem being used in production it was essential that it was not only tested but that testing was constantly carried out. By using CircleCI tests are automatically run when a change is deployed to the git repository.
- Coverage
 - Code coverage reports show how many lines of the code are touched by the tests this is very useful to ensure all the code branches are tested and perform as expected.

- Code coverage ended up at 46% at the end of the project.
- Code quality
 - Codeclimate is a service which analyzes code and detects “code smells”. “A code smell is any characteristic in the source code of a program that possibly indicates a deeper problem.” - Wikipedia. This includes problems such as:
 - * Cognitive complexity: how complicated is the code to understand.
 - * File and function length: does the file or function contain too many lines (only counting actual lines of code, ie. not comments or blank lines)
 - * Duplication: are large parts of the code duplicated in multiple places
 - Codeclimate then predicts the amount of time it would take to fix this technical debt. At the project end, skeem contained 147 code smells with a predicted clean up time of 2 months.

11 Deployments

- Throughout the development i has the opportunity to deploy skeem on real world applications. At this point skeem is in production use in two apps and in employed on 3 further, currently in development, projects.
- Resooma
 - Resooma is a bills consolidation company focussing primarily on university students.
 - 88 distinct models, 700 attributes, 50'000'000 users
- Quote generator and Stock management tool for Enterprise Security Distributions Norwich
 - Tracks more than XXX quotes for customers concerning more than 20'000 products. Used heavily
- II
 - Invoice fraud detection using machine learning
- Voluble
 - Messages API
- Rolecall
 - Job tracking and communication platform aimed at contract work-

ers

There are a wide range of projects using skeem. Very flexible.

12 Conclusion

Skeem has turned out to be a very successfull project already helping out a wide range of projects

12.1 Future Work

References

“Monorepo.” 2019. *Wikipedia*. Wikimedia Foundation. <https://en.wikipedia.org/wiki/Monorepo>.