



UNIVERSITÀ DEGLI STUDI ROMA TRE

Facoltà di Ingegneria
Corso di Laurea Magistrale in Ingegneria Informatica

SNP Web-portal

Autori

U. Buonadonna, D. Sicignani, D. Tosoni

Anno Accademico

2013/2014

Capitolo 1

Database Description

In this chapter we will explain the creation of the database, from code point of view as well as from design point of view.

The **database** is the core part of the back end of the application, because it is the structure that keeps stored all the data necessary for the operation of the portal. Because of its importance, we spent much of our time to decide its structure (tables, relationships, multiplicity, ...) and only after a thorough analysis of results we moved to implementation.

As explained in the introductory chapter, the software implementation of database was performed using **MongoDB** through **Mongoose** (a MongoDB object modeling tool). MongoDB is an open-source document database, and the leading NoSQL database. It is based on *Document-Oriented Storage* and unlike SQL databases, where you must determine and declare a table's schema before inserting data, MongoDB's collections do not enforce document structure. This flexibility facilitates the mapping of documents to an entity or an object.

Mongoose helps to have a more efficient management of the objects; in fact, using mongoose's schemas to define database's objects, permits to treat database's models as commons JavaScript's objects.

Throughout the chapter, we will see the objects that compose the model of our da-

tabase, but, first of all, let us focus on the syntax used to communicate with MongoDB and Mongoose.

1.1 Mongoose's syntax

To realize an object using mongoose, we need to follow a specific syntax.

First of all, we should place all files related to the models (one file for each model) in a specific folder of the project; this folder is as follows:

```
/server/models
```

At the beginning of each file, we have to inform the database that all data passed to builders, but that are not present in the schema that we will define below, must not be stored.

To do this we use the following line of code:

```
1      'use strict';  
2
```

Next, we must specify needed *modules dependencies*. In our project, we have to insert the dependency from mongoose and crypto (a module for encryption of information):

```
1      var mongoose = require('mongoose'),  
2          Schema = mongoose.Schema,  
3          crypto = require('crypto');  
4
```

Now we move to define the *schema*, specifying all necessary components (attributes and relationships). Note that MongoDB will automatically add an ID field to each schema and assigns it using an internal policy.

For example, if we want to model a *User* with an attribute *name* of type *string*:

$\langle \text{Example schema?} \rangle \equiv$

```

1      UserSchema = new Schema ({
2          // attribute
3          name: {
4              type: String,
5              required: true,
6              validate: [validatePresenceOf
7          , 'Name can not be blank']
8          },
9          // relationship
10         patients: [{
11             type: Schema.Type.ObjectId,
12             ref: 'Patient'
13         }]
14     });
15

```

◇

Fragment never referenced.

Next, we can define a number of *additional features* to our model, such as:

- **Validations:** functions called to validate a set of attributes or data specified by the programmer
- **Virtuals:** virtual attributes, not persisted
- **Pre-save hooks:** functions called at the time immediately before the model is saved to the database
- **Methods:** all methods needed for the operations on our model

Finally, we must specify the actual name of the model described by the schema specified as previously shown:

```

1  mongoose.model ( 'User' , UserSchema );

```

1.2 Model Schemas

These are the model schemas that compose our database.

1.2.1 Variants

The model *Variant* represents the key concept that we want to represent in our work: **genomic variants**.

This table will always be accessed during any kind of consultation, in fact are variants what biologists are looking for or insert through the portal.

The centrality of the concept and, therefore, of the table is also demonstrated by the fact that in the database schema (see chapter *Introduction*) Variant table is linked to almost all the other tables.

Let us look at the code components:

Attributes

$\langle \text{Variant's attributes ?} \rangle \equiv$

```

1  /* Variant schema */
2  var VariantSchema = new Schema({
3    chr: {
4      type: String,
5      required: true,
6      validate: [validatePresenceOf, 'CHR cannot be blank']
7    },
8    start: {
9      type: String,
10     required: true,
11     validate: [validatePresenceOf, 'start cannot be blank']
12   },
13   end: {
14     type: String,
15     required: true,
16     validate: [validatePresenceOf, 'End cannot be blank']
17   },
18   ref: {
19     type: String,
20     required: true,
21     validate: [validatePresenceOf, 'Ref cannot be blank']
22   },
23   alt: {
24     type: String,
25     required: true,
26     validate: [validatePresenceOf, 'Alt cannot be blank']
27   },

```

◇

Fragment never referenced.

- **chr**: it is a five-character string describing the chromosome affected by the variation

- **start, end:** integers that describe the coordinate of the mutation within the chromosome; if they match, it means that we have a punctiform mutation
- **ref:** represents the correct bases, as they would be if there were no mutation
- **alt:** represents what we actually find due to the mutation; a "." or a "-" represent a deletion, while a sequence, for example CTTG..., represents an insertion

Relationships

$\langle \text{Variant's relationships ?} \rangle \equiv$

```
1 // Here start relations
2 gene: {
3     type: Schema.Types.ObjectId ,
4     ref: 'Gene'
5 },
6
7 pathogenicity: {
8     type: Schema.Types.ObjectId ,
9     ref: 'Pathogenicity'
10 },
11
12 patients: [{
13     type: Schema.Types.ObjectId ,
14     ref: 'Patient'
15 }],
16
17 dbSNPs: [{
18     type: Schema.Types.ObjectId ,
19     ref: 'DbSNP'
20 }],
21
22 esps: [{
23     type: Schema.Types.ObjectId ,
24     ref: 'esp'
25 }],
26
27 sequencings: [{
28     type: Schema.Types.ObjectId ,
29     ref: 'Sequencing'
30 }]
31
32 });
```

◇

Fragment never referenced.

- **gene**: links the variant with the gene affected by the mutation
- **pathogenicity**: describes the variant pathogenicity, that is if the variant cause diseases or not
- **patients**: links to patients affected by this variant
- **dbSNPs**: connects the variant to its dbSNPs
- **esps**: Exome Sequencing Project
- **sequencings**: all sequencings affected by this variation

1.2.2 dbSNP

dbSNP is the key concept of the database, that is the **Single Nucleotide Polymorphism Database** (a free public archive for genetic variation within and across different species developed and hosted by the *National Center for Biotechnology Information (NCBI)* in collaboration with the *National Human Genome Research Institute (NHGRI)*).

Let us look at the code components:

Attributes

$\langle \text{dbSNP's attributes ?} \rangle \equiv$

```
1 /* DbSNP schema */
2 var DbSNPSchema = new Schema({
3   dbSNP: {
4     //key
5     type: String,
6     required: true,
7     unique: true,
8     validate: [validatePresenceOf, 'DbSNP cannot be blank']
9   },
10  freqAlt: {
11    type: String,
12    required: true,
13    validate: [validatePresenceOf, 'freqAlt cannot be blank']
14  },
15  freqRef: {
16    type: String,
17    required: true,
18    validate: [validatePresenceOf, 'freqRef cannot be blank']
19  },
```

◇

Fragment never referenced.

- **dbSNP**: identifier of the variation in the dbSNP database
- **freqAlt**: frequency of the alternative allele (variant) in the population 1000 Genomes
- **freqRef**: frequency of the reference allele

Relationships

$\langle dbSNP's\ relationships\ ? \rangle \equiv$

```
1 //Relationship
2 variants: [{
3     type: Schema.Type.ObjectId ,
4     ref: 'Variant'
5 }]
6
7 });
```

◇

Fragment never referenced.

- **variants:** links dbSNP to its variants

1.2.3 Gene

A *Gene* is the molecular unit of heredity of a living organism. Genes hold the information to build and maintain an organism's cells and pass genetic traits to offspring. All organisms have genes corresponding to various biological traits, some of which are instantly visible, such as eye color or number of limbs, and some of which are not, such as blood type, increased risk for specific diseases, or the thousands of basic biochemical processes that comprise life.

Let us look at the code components:

Attributes

$\langle \text{Gene's attributes?} \rangle \equiv$

```
1  /* Gene schema */
2  var GeneSchema = new Schema({
3      gene: {
4          type: String,
5          required: true,
6          validate: [validatePresenceOf, 'Gene cannot be blank']
7      },
8      region: {
9          type: String,
10         required: true,
11         validate: [validatePresenceOf, 'Region cannot be blank']
12     },
13     mutation: {
14         type: String,
15         required: true,
16         validate: [validatePresenceOf, 'Mutation cannot be blank']
17     },
18     annotation: {
19         type: String,
20         required: true,
21         validate: [validatePresenceOf, 'Annotation cannot be blank']
22     },
```

◇

Fragment never referenced.

- **gene**: identifier of the gene
- **region**: the region in which the gene is located
- **mutation**: the mutation that affects the gene
- **annotation**: any annotations

Relationships

$\langle \text{Gene's relationships ?} \rangle \equiv$

```
1      // Relationship
2      variants: [{
3          type: Schema.Type.ObjectId ,
4          ref: 'Variant'
5      }]
6
7  });
```

◇

Fragment never referenced.

- **variants:** links gene to its variants

1.2.4 Pathogenicity

Pathogenicity indicates whether a mutation is pathogenic (able to create damages to the organism) or not.

Let us look at the code components:

Attributes

$\langle \text{Pathogenicity's attributes ?} \rangle \equiv$

```

1  /* Pathogenicity schema */
2  var PathogenicitySchema = new Schema({
3    SIFT: {
4      type: String ,
5      required: true ,
6      validate: [validatePresenceOf , 'SIFT cannot be blank']
7    },
8
9    polyPhen: {
10     type: String ,
11     required: true ,
12     validate: [validatePresenceOf , 'polyPhen cannot be blank']
13   },
14
15   mutationTaster: {
16     type: String ,
17     required: true ,
18     validate: [validatePresenceOf , 'mutationTaster cannot be
19 blank']
20   },
21
22   mutationAssessor: {
23     type: String ,
24     required: true ,
25     validate: [validatePresenceOf , 'mutationAssessor cannot be
26 blank']
27   },
28
29   GERPpp: {
30     type: String ,
31     required: true ,
32     validate: [validatePresenceOf , 'GERPpp cannot be blank']
33   },
34
35   pyoloP: {
36     type: String ,
37     required: true ,
38     validate: [validatePresenceOf , 'pyoloP cannot be blank']
39   },
40
41   SiPhy: {
42     type: String ,
43     required: true ,
44     validate: [validatePresenceOf , 'SiPhy cannot be blank']
45   },
46
47   ...
48 }

```

- **SIFT, polyPhen, mutationTaster, mutationAssessor**: all four are predictors of pathogenicity of the mutation. The numbers are the scores of the pathogenicity.
- **GERPpp, pyoloP, SiPhy**: all three are predictors of conservation of the mutation. The numbers are the scores of conservation.

Relationships

$\langle \text{Pathogenicity's relationships?} \rangle \equiv$

```

1      // Relationship
2      variant: {
3          type: Schema.Type.ObjectId ,
4          ref: 'Variant'
5      }
6
7  });
```

◇

Fragment never referenced.

- **variant**: links pathogenicity to its variant

1.2.5 Sequencing

DNA Sequencing is the process of determining the nucleotide order of a given DNA fragment. The sequence of DNA encodes the necessary information for living things to survive and reproduce, so determining the sequence is therefore useful in fundamental research into why and how organisms live, as well as in applied subjects.

Let us look at the code components:

Attributes

$\langle \text{Sequencing's attributes ?} \rangle \equiv$

```

1  /* Sequencing schema */
2  var SequencingSchema = new Schema({
3      //patient id and date are keys
4      patientId: {
5          type: String ,
6          required: true ,
7          validate: [validatePresenceOf , 'patient cannot be blank']
8      },
9      date: {
10         type: String ,
11         required: true ,
12         validate: [validatePresenceOf , 'date cannot be blank']
13     },
14     patientHealthStatus: {
15         type: String ,
16         required: true ,
17         validate: [validatePresenceOf , 'patientHealthStatus cannot be
18         blank']
19     },
20     sequencerName: {
21         type: String ,
22         required: true ,
23         validate: [validatePresenceOf , 'sequencerName cannot be blank
24         ' ]
25     },
26     sequencerModel: {
27         type: String ,
28         required: true ,
29         validate: [validatePresenceOf , 'sequencerModel cannot be
30         blank' ]
31     },
32     referenceGenome: {
33         type: String ,
34         required: true ,
35         validate: [validatePresenceOf , 'referenceGenome cannot be
36         blank' ]
37     },
38     detail :{
39         qual: {
40             type: String ,
41             required: true ,
42             validate: [validatePresenceOf , 'qual cannot be blank']
43         },
44         filter: {

```


- **patientId**: identifier of the patient.
- **date**: the date when sequencing was performed.
- **patientHealthStatus**: if the patient is diseased or not.
- **sequencerName**: name of the sequencer.
- **sequencerModel**: model of the sequencer.
- **referenceGenome**: the reference genome.
- **qual**: describes the quality of the sequencing.
- **filter**: reference allele.
- **genotype**
- **genotypeQuality**: quality of the genotype.
- **readsDepth**: the depth of the reads.
- **ref**: string that indicates whether the quality of the variant is good (PASS).
- **altFilterReads**: number of reads that do not contain the mutation to their inside.
- **genotypeLikelihood**: probability that the determined genotype is correct.
- **haplotypeScore**: indicate the presence of misaligned reads in the neighborhood of the variant.
- **strandBias**: all four are predictors of pathogenicity of the mutation.

Relationships

$\langle \textit{Sequencing's relationships ?} \rangle \equiv$

```

1      },
2      genotype: {
3          type: String,
4          required: true,
5          validate: [validatePresenceOf, 'genotype cannot be blank']
6      },
7      genotypeQuality: {

```

◇

Fragment never referenced.

- **variants**: links sequencing to its variants

1.2.6 ESP

ESP is the part of the database that represents locally the interesting portion of the database of the **Exome Sequencings Project**, the goal of which is to discover novel genes and mechanisms contributing to heart, lung and blood disorders by pioneering the application of next-generation sequencing of the protein coding regions of the human genome across diverse, richly-phenotyped populations and to share these datasets and findings with the scientific community to extend and enrich the diagnosis, management and treatment of heart, lung and blood disorders.

Let us look at the code components:

Attributes

$\langle esp's\ attributes\ ? \rangle \equiv$

```

1  /* Esp schema */
2  var EspSchema = new Schema({
3      ESP6500_ALL: {
4          type: String,
5          required: true,
6          validate: [validatePresenceOf, 'ESP6500_ALL cannot be blank']
7      },
8      ESP6500_AA: {
9          type: String,
10         required: true,
11         validate: [validatePresenceOf, 'ESP6500AA cannot be blank']
12     },
13     ESP6500_EA: {
14         type: String,
15         required: true,
16         validate: [validatePresenceOf, 'ESP6500_EA cannot be blank']
17     },

```

◇

Fragment never referenced.

- **ESP6500_ALL**: frequency of the mutation in the EVS database and in the overall population
- **ESP6500_AA**: frequency of the mutation in the EVS database and in the American/African population
- **ESP6500_EA**: frequency of the mutation in the EVS database and in the European /American population

Relationships

$\langle esp's\ relationships\ ? \rangle \equiv$

```
1      //Relationship
2      variants: [{
3          type: Schema.Type.ObjectId ,
4          ref: 'Variant'
5      }]
6
7  });
```

◇

Fragment never referenced.

- **variants:** links esp to its variants

1.2.7 Patient

A *Patient* describes a patient in the real world. Can only be created by a user with Admin privileges.

Let us look at the code components:

Attributes The Patient object has no attributes to be specified; it only has an identifier (**ID**) that is added from MongoDB by default.

Relationships

$\langle \text{Patient's relationships ?} \rangle \equiv$

```
1  /* Patient schema */
2  var PatientSchema = new Schema({
3
4      // Relationship
5      variants: [{
6          type: Schema.Type.ObjectId ,
7          ref: 'Variant'
8      }],
9
10     family: {
11         type: Schema.Type.ObjectId ,
12         ref: 'Family'
13     }
14
15 });
```

◇

Fragment never referenced.

- **variants:** links patient to its variants
- **family:** links patient to its family

1.2.8 Family

A *Family* represents a set of one or more *patients* (modeled in patient schema). It can be created only by an user with admin privileges. This schema allow to group patients with similar genomic mutations as well as real families (father, mother, son...).

Let us look at the code components:

Attributes

$\langle \textit{Family's attributes ?} \rangle \equiv$

```

1  /* Family schema */
2  var FamilySchema = new Schema({
3
4      name: {
5          type: String ,
6          required: true ,
7          validate: [validatePresenceOf , 'Name cannot be blank']
8      },

```

◇

Fragment never referenced.

- **name:** describes family's name

Relationships

$\langle \textit{Family's relationships ?} \rangle \equiv$

```

1      //Relationship
2      patients: [{
3          type: Schema.Type.ObjectId ,
4          ref: 'Patient'
5      }]
6
7  });

```

◇

Fragment never referenced.

- **patients:** links family to its members

1.2.9 User

The *User* object describes a user of the system, a *biologist* who accesses the portal to carry out one of the supported operations (described in the introductory chapter).

Let us look at the code components:

Attributes

$\langle \text{User's attributes ?} \rangle \equiv$

```

1  /**
2   * User Schema
3   */
4  var UserSchema = new Schema({
5      name: {
6          type: String,
7          required: true,
8          validate: [validatePresenceOf, 'Name cannot be blank']
9      },
10     email: {
11         type: String,
12         required: true,
13         match: [/.\+\@\.\+\.\+/, 'Please enter a valid email'],
14         validate: [validatePresenceOf, 'Email cannot be blank']
15     },
16     username: {
17         type: String,
18         unique: true,
19         validate: [validatePresenceOf, 'Username cannot be blank']
20     },
21     roles: {
22         type: Array,
23         default: ['authenticated']
24     },
25     hashed_password: {
26         type: String,
27         validate: [validatePresenceOf, 'Password cannot be blank']
28     },
29     provider: {
30         type: String,
31         default: 'local'
32     },
33     salt: String,
34     facebook: {},
35     twitter: {},
36     github: {},
37     google: {},
38     linkedin: {}
39 });

```

◇

Fragment never referenced.

These attributes are all self-explanatory.

Methods

$\langle \text{User's validation?} \rangle \equiv$

```
1  /**
2   * Validations
3   */
4  var validatePresenceOf = function(value) {
5      // If you are authenticating by any of the oauth strategies, don't
6      // validate.
7      return (this.provider && this.provider !== 'local') || value.length;
8  };
```

◇

Fragment never referenced.

$\langle \text{User's methods?} \rangle \equiv$

```

1  /**
2   * Virtuals
3   */
4  UserSchema.virtual('password').set(function(password) {
5      this._password = password;
6      this.salt = this.makeSalt();
7      this.hashPassword(password);
8  }).get(function() {
9      return this._password;
10 });
11
12 /**
13  * Pre-save hook
14  */
15 UserSchema.pre('save', function(next) {
16     if (this.isNew && this.provider === 'local' && this.password && !
17         this.password.length)
18         return next(new Error('Invalid password'));
19     next();
20 });
21
22 /**
23  * Methods
24  */
25 UserSchema.methods = {
26     /**
27      * HasRole - check if the user has required role
28      *
29      * @param {String} plainText
30      * @return {Boolean}
31      * @api public
32      */
33     hasRole: function(role) {
34         var roles = this.roles;
35         return roles.indexOf('admin') !== -1 || roles.indexOf(role)
36             !== -1;
37     },
38     /**
39      * IsAdmin - check if the user is an administrator
40      *
41      * @return {Boolean}
42      * @api public
43      */

```

Functions performed by these methods are explained in the code comments.