# CSCI 3700 - Database Management Systems Data-intensive Application Development Term Project

Jack Edwards and Oliver Chen

## Contents

# 1 Data-intensive Application Selection

A data-intensive application is characterized by: lots of data; relatively little processing; large number of insertions, deletions, updates, and queries. Data is viewed as a corporate resource and is often used for tasks varying from marketing campaigns, new product strategies, inventory management and distribution logistics, to improving customer loyalty.

Ideally, the application should have: well-understood or well-designed business processes; sufficient documentation about the business processes; tasks are known; application boundary is well-demarcated; allows developing one organization-wide database schema from multiple department-wide database views.

You should either possess sufficient domain knowledge of the application area, or have access to a domain expert. Otherwise, it is extremely difficult to develop a successful data-intensive application.

Describe the data-intensive application that you have selected. What are its characteristics? What makes it data-intensive? Who is the sponsor of this project? Who are the end-users? How will they benefit from this application?

Answer: We have selected the Project Gutenberg as our topic, it has over 57,000 online eBooks. The project would be create a database that hold all the books, identify the entities, attributes and relationships. The database would act as a library allow users to do searches base on their need. Also the database system can be implement with information retrieval system.

# 2 Identification and Documentation of Use-cases

Identify various *classes of users* (aka *actors*) for the data-intensive application. Note that users can be human as well other systems. In some applications, *time* is a user. For example, end-of-day, end-of-week, end-of-month, end-of-quarter, and end-of-fiscal-year are all time-triggered events. Your application needs to respond to these time-triggered events.

*Use-cases* describe interactions between the users and the system. Some interactions can be normal (no error conditions), other interactions may entail additional processing (e.g., preferred customers receive additional services), and yet other interactions require error recovery due to various conditions such as erroneous input or device malfunctioning. Each path through a use-case is called a *scenario*. In other words, a use-case is a set of related scenarios.

Conceptually, a use-case represents a *unit of work* from an end-user perspective. A use-case involves executing a set of tasks in certain sequence.

How may user classes do you have? How many actors (including human, application, and abstract ones like time)? Name use-cases. Document them using the LaTeX template.

Answer: We will have two user classes End users and administrators. Administrators can add documents to the database as well as performing any actions available to end users. End users can stream all documents from the database, they can search documents by language, title, author, and year. We plan to provide java and python libraries to interface with the database.

# 3 Use-case Diagram

*Use-case diagram* is a pictorial representation of interactions between the application users and use-cases. It also shows relationships between use-cases such as one use-case being embedded in another use-case, or one use-case extending the functionality of another use-case.

Answer:



Figure 1: Use-case digram

# 4 Identification and Documentation of Data Tasks in the Application

Each use-case scenario requires executing a set of tasks. For each task identify and document inputs needed, and outputs generated. Also, specify possible error conditions that might occur as inputs are transformed into outputs.

Answer: User could query author name, book title, primary key of certain book, since these could be directly query from the database. We can do a direct query to gather these information. Also if user need to request all documents from database, the database will return all documents in json format with all attributes as the fields. It would be more complicated if user try to search by keywords. This process require information retrieval system to process text and analyze documents, the function will be add if the implementation is needed.
The possible errors these processes might occurs are, user type in wrong information, user misspell words, ir system could not find the information user need.

# 5    Identification and Documentation of Transactions

A *transaction* is a unit of work both from a database end-user perspective as well as from the database system perspective. A transaction requires executing all the tasks that comprise a unit of work in entirety  all or nothing proposition. For each transaction specify its frequency of execution.

Answer: The main transaction will be adding new books to the database, the title, the release date, the language and the full text must be added. Books and authors are represented by separate entities therefore if a books author is not present in the database it must also be created. To create and Author the authors full nane is required. The transaction to create a book will be executed much more frequently than the transaction to create an author.

# 6    Identification and Documentation of Database Queries

Unlike transactions, database *queries* do not change the data in the database. Queries require only read access to the database. Some queries may take quite a bit of time to complete. Therefore, performance is often an issue for database queries.

Specify queries in plain English. For each query specify what data is to be retrieved (not how) as well as its frequency of execution.

Describe your database queries (in English, not in SQL) here.

Answer: Queries sorted by frequency,

1. Stream the full text of all books in the database.
2. Look up book by primary key (release date, title) returns all information about the book
3. Look up author by primary key (full name) returns all book written by the author.
4. Look up books by language, all book in the specified language are returned.
5. Look up all books released in a certain year all books written in that year are returned

# 7    Conceptual Data Model

Start with Entity-Relationship (E-R) diagram for department-wise transactions and queries. The number of departments you will have (e.g., registrar, library, financial aid, campus housing) depends on the scope of the data-intensive application. In the second step, integrate these department-wise diagrams into one corporate-wide ER diagram. Follow established diagrammatic conventions. Use SQL Power Architect tool (or similar) for developing E-R diagrams.

# 8 Logical Data Model

Identify functional and multivalued dependencies. Transform ER/EER diagrams into a relational schema. Determine functional dependencies and perform normalization. Transform each table into 3NF or BCNF using the functional dependencies and normalization rules. You may use Database Design (DBD) tool for this task. For each table in the final schema, specify primary and foreign keys. Also specify data integrity constraints.

Answers:
book(title, text, release_date, language)
attribute language can be mutivalued.
written_by(book_key, author_key)
author(last_name, first_name, middle_name, suffix)

# 9 Physical Data Model

For each database file, specify *initial* storage structures and access paths. Typically, these storage structures and access paths need modifications based on *observed performance* once the database is in operation (aka database tuning).



# 10 Database Creation and Data Loading

Now that your logical database schema and physical database design is in place, write SQL scripts to create the database using PostgreSQL. Load existing data into the tables using either SQL statements or *bulk loading*. Resolve any data integrity constraint violations.

Database loading script is located in code listings.

## 11 Implementing Database Transactions and Queries

Write SQL code for transactions and queries. Verify and validate all transactions and queries. Comment SQL code sensibly.

Here is an LaTeX markup for typesetting SQL code.

```
1  SELECT name AS "Country␣Name", population AS "Population",
       lifeexpectancy AS "Life␣Expectancy"
2  FROM    country
3  WHERE   lifeexpectancy IS NOT NULL
4  ORDER   BY lifeexpectancy DESC;
```

For the uninitiated, you may use the verbatim command for an uninspiring typeset.

```
SELECT name AS "Country Name",
       population      AS "Population",
       lifeexpectancy AS "Life Expectancy"
FROM   country
WHERE  lifeexpectancy IS NOT NULL
ORDER  BY lifeexpectancy DESC;
```

```
1  CREATE DATABASE gutenburg
2      WITH
3      OWNER = postgres
4      ENCODING = 'UTF8'
5      CONNECTION LIMIT = -1;
6
7  CREATE TABLE public."Book"
8  (
9      gutenberg_id text NOT NULL,
10     release_date date,
11     full_text text,
12     language text NOT NULL,
13     title text NOT NULL,
14     PRIMARY KEY (gutenberg_id)
15 )
16 WITH (
17     OIDS = FALSE
18 );
19
20 ALTER TABLE public."Book"
21     OWNER to postgres;
22
23 CREATE TABLE public."Author"
24 (
```

```
25    author_id serial ,
26    first_name text NOT NULL ,
27    last_name text ,
28    middle_name text ,
29    suffix text ,
30    prefix text ,
31    PRIMARY KEY ( author_id )
32 )
33 WITH (
34    OIDS = FALSE
35 );
36
37 ALTER TABLE public . " Author "
38    OWNER to postgres ;
39
40 CREATE TABLE public . " Written_By "
41 (
42    author_id integer NOT NULL ,
43    gutenberg_id text NOT NULL ,
44    PRIMARY KEY ( author_id , gutenberg_id ),
45    CONSTRAINT gutenberg_id FOREIGN KEY ( gutenberg_id )
46        REFERENCES public . " Book " ( gutenberg_id ) MATCH SIMPLE
47        ON UPDATE NO ACTION
48        ON DELETE NO ACTION ,
49    CONSTRAINT author_id FOREIGN KEY ( author_id )
50        REFERENCES public . " Author " ( author_id ) MATCH SIMPLE
51        ON UPDATE NO ACTION
52        ON DELETE NO ACTION
53 )
54 WITH (
55    OIDS = FALSE
56 );
57
58 ALTER TABLE public . " Written_By "
59    OWNER to postgres ;
60
61 -- the symbol ? will represent a variable to be replaced with
       the user 's actual information need
62
63 -- get a book by gutenberg_id
64 Select * From public . " Book "
65 Where gutenberg_id = ?
66
67 -- get book by title
68 Select * From public . " Book "
```

```
69  Where title = ?
70
71  -- get all books in a certain language
72  Select * from public."Book"
73  Where language = ?
74
75  -- get author by name
76  Select * From public."Author"
77  Where first_name = ? And middle_name = ?
78  And last_name = ? And suffix = ?
79  And prefix = ?
80
81  --get author by author_id
82  Select * from public."Author"
83  Where author_id = ?
84
85  --get all books by an author with author_id
86  select * from public."Book"
87  natural join public."Written_By"
88  where author_id = ?
89
90  --get all authors of a book with gutenberg_id
91  select * from public."Author"
92  natural join public."Written_By"
93  where public."Written_By".gutenberg_id = ?
94
95  --get all available languages
96  select distinct language from public."Book"
```

## 12   Developing Database Applications

This step involves writing database applications using Java or scripting languages such as JSP, PHP, and ASP.NET. Include rationale for choosing a specific language for developing the database applications. Students should not choose a scripting language unless they are already familiar with it. Simply there is no time to learn a new scripting language. Demonstrate a simple Web application based on the database that you have developed.

Discuss the design and implementation details of the database application here. Please do not include actual code. You may include code in Appendix.

Answer:

We had decided to use Angular as our frontend, and flask as our backend. Postgresql would be the database. Angular is written in typescript which is pretty much the javascript and Flask is written in python. We would create the functionality base on the application functions that we discuss in the previous section.

**Welcome to Angular + Python + Flask Demo!**

Search by Author

[william shakespeare] [Search] Server Respond: { "respond": "the complete works of william shakespeare<br />shakespeare's sonnets<br />venus and adonis<br />venus and adonis<br />king richard iii<br />romeo and juliet<br />much ado about nothing<br />as you like it<br />the merchant of venice [liberally edited by charles kean]<br />a fairy tale in two acts taken from shakespeare (1763)<br />macbeth<br />hamlet<br />la tempête<br />the taming of the shrew<br />romeo and juliet<br />much ado about nothing<br />julius caesar<br />hamlet<br />the phoenix and the turtle<br />troilus and cressida<br />coriolan<br />othello, the moor of venice<br />king lear<br />pericles<br />the winter's tale<br />the winter's tale<br />hamlet<br />romeo ja julia<br />beaucoup de bruit pour rien<br />jules césar<br />la comédie des méprises<br />timon d'athènes<br />antoine et cléopâtre<br />le jour des rois<br />kumigas lear<br />antonius ja cleopatra<br />les deux gentilshommes de vérone<br />macbeth<br />les alegres comares de windsor<br />othello<br />troïlus ja cressida<br />le songe d'une nuit d'été<br />roméo et juliette<br />timon atenalainen<br />comme il vous plaira<br />mesure pour mesure<br />othello<br />le conte d'hiver<br />le roi lear<br />troïlus et cressida<br />julius caesar<br />cymbeline<br />la méchante femme mise à la raison<br />peines d'amour perdues<br />périclès<br />coriolanus<br />les joyeuses bourgeoises de windsor<br />le marchand de venise<br />titus andronicus<br />la vie et la mort du roi richard ii<br />le roi jean<br />la festa dels reis<br />the works of william shakespeare [cambridge edition] [vol. 6 of 9 vols.]<br />the works of william shakespeare [cambridge edition] [vol. 8 of 9 vols.]<br />the works of william shakespeare [cambridge edition] [vol. 5 of 9]<br />romeo en julia<br />the works of william shakespeare - cambridge edition<br />de koopman van venetië<br />fair em<br />de getemde feeks<br />dramas<br />koning jan<br />hamlet<br />richard iii<br />macbeth<br />coriolanus<br />romeo und julia<br />ein sommernachtstraum<br />wie es euch gefällt<br />der kaufmann von venedig<br />othello<br />othello<br />was ihr wollt<br />die irrungen (the comedy of errors)<br />timon von athen<br />romeo und juliette<br />maass fuer maass (measure for measure)<br />der sturm<br />das leben und der tod des koenigs lear<br />ein st. johannis nachts-traum<br />macbeth<br />hamlet, prinz von dannemark<br />leben und tod des koenigs johann<br />leben und tod konigs richard des zweyten<br />der erste theil von koenig heinrich dem vierten<br />der zweyte theil von koenig heinrich dem vierten<br />the tragicall historie of hamlet, prince of denmarke<br />julius caesar" }

Search by Title

[venus and adonis] [Search] Server Respond: { "respond": "<br /><br /><br />Produced by Dianne Bean, and David Widger<br /><br /><br /> VENUS AND ADONIS <br /><br /><br /> by William Shakespeare <br /><br /><br /> _Vilia miretur vulgus; mihi flavus Apollo<br />Pocula Castalia plena ministret aqua._<br /><br /><br />TO THE RIGHT HONOURABLE<br /><br />HENRY WRIOTHESLEY, EARL OF SOUTHAMPTON,<br /><br />and Baron of Titchfield. <br /><br /><br />Right Honourable, I know not how I shall offend in dedicating my<br />unpolished lines to your lordship, nor how the world will censure me<br />for choosing so strong a prop to support so weak a burthen: only, if<br />your honour seem but pleased, I account myself highly praised, and vow<br />to take advantage of all idle hours, till I have honoured you with some<br />graver labour. But if the first heir of my invention prove deformed, I<br />shall be sorry it had so noble a godfather, and never after ear so<br />barren a land, for fear it yield me still so bad a harvest. I leave it<br />to your honourable survey, and your honour to your heart's content;<br />which I wish may always answer your own wish and the world's hopeful<br />expectation.<br /><br />Your honour's in all duty,<br /><br />WILLIAM SHAKESPEARE <br /><br /> VENUS AND ADONIS <br /><br /> Even as the sun with purple-colour'd face<br />Had ta'en his last leave of the weeping morn,<br />Rose-cheek'd Adonis tried him to the chase;<br />Hunting he lov'd, but love he laugh'd to scorn;<br />Sick-thoughted Venus makes amain unto him,<br />And like a bold-fac'd suitor 'gins to woo him.<br /><br />'Thrice fairer than myself,' thus she began,<br />'The field's chief flower, sweet above compare,<br />Stain to all nymphs, more lovely than a man,<br />More white and red than doves or roses are:<br />Nature that made thee, with herself at strife,<br />Saith that the world hath ending with thy life. 12<br />'Vouchsafe, thou wonder, to alight thy steed,<br />And rein his proud head to the saddle-bow;<br />If thou wilt deign this favour, for thy meed<br />A thousand honey secrets shalt thou know: 16<br /> Here come and sit, where never serpent hisses,<br />And being set, I'll smother thee with kisses.<br />'And yet not cloy thy lips with loath'd satiety,<br />But rather famish them amid their plenty, 20<br />Making them red, and pale, with fresh variety;<br />Ten kisses short as one, one long as twenty.<br /> A summer's day will seem an hour but short,<br />Being wasted in such time-beguiling sport." 24<br />'With this she seizeth on his sweating palm,<br />The precedent of pith and livelihood,<br />And trembling in her passion, calls it balm,<br />Earth's sovereign salve to do a goddess good: 28<br />Being so enrag'd, desire doth lend her force<br />Courageously to pluck him from his horse.<br />Over one arm the lusty courser's rein,<br />Under her other was the tender boy, 32<br />Who blush'd and pouted in a dull disdain,<br />With leaden appetite, unapt to toy;<br /> She red and hot as coals of glowing fire,<br />He red for shame, but frosty in desire. 36<br />'So soon she would be thrust,<br />And govern'd him in strength, though not in lust <br />So soon was she along, as he was down,<br />Each leaning on their elbows and their hips: 44<br />Now doth she stroke his cheek, now doth he frown,<br />And 'gins to chide, but soon she stops his lips,<br />And kissing speaks, with lustful language broken,<br />'If thou wilt chide, thy lips shall never open." 48<br />He burns with bashful shame, she with her tears<br />Doth quench the maiden burning of his cheeks;<br />Then with her windy sighs and golden hairs<br />To fan and blow them dry again she seeks. 52<br />He saith she is immodest, blames her miss;<br /> What follows more, she murders with a kiss.<br />Even as an empty eagle, sharp by fast,<br />Tires with her beak on feathers, flesh and bone, 56<br />Shaking her wings, devouring all in haste,<br />Till either gorge be stuff'd or prey be gone:<br />Even so she kiss'd his brow, his cheek, his chin,<br />And where she ends she doth anew begin. 60<br />Forc'd to content, but never to obey,<br />Panting he lies, and breatheth in her face.<br />She feedeth on the steam, as on a prey,<br />And calls it heavenly moisture, air of grace, 64<br /> Wishing her cheeks were gardens full of flowers<br /> So they were dew'd with such distilling showers.<br />Look how a bird lies tangled in a net,<br />So fasten'd in her arms Adonis lies; 68<br />Pure shame and aw'd resistance made him fret,<br />Which bred more beauty in his angry eyes:<br /> Rain added to a river that is rank<br />Perforce will force it overflow the bank. 72<br />Still she entreats, and prettily entreats,<br />For to a pretty ear she tunes her tale.<br />Still is he sullen, still he lours and frets,<br /> Twixt crimson shame and anger ashy pale; 76<br />Being red she loves him best, and being white,<br /> Her best is better'd with a more delight.<br />Look how he can, she cannot choose but love;<br />And by her fair immortal hand she swears, 80<br />From his soft bosom never to remove,<br />Till he take truce with her contending tears,<br /> Which long have rain'd, making her cheeks all wet;<br />And one sweet kiss shall pay this countless debt.<br />Upon this promise did he raise his chin, 85<br />Like a dive-dapper peering through a wave,<br />Who, being look'd on, ducks as quickly in;<br />So offers he to give what she did crave, 88<br /> But when her lips were ready for his pay,<br /> He winks, and turns his lips another way.<br />Never did passenger in summer's heat<br />More thirst for drink than she for this good turn. 92<br />Her help she sees, but help she cannot get;<br />She bathes in water, yet her fire must burn:<br /> "O! pity," 'gan she cry, "flint-hearted boy,<br /> 'Tis but a kiss I beg, why art thou coy? 96<br />"I have been woo'd as I entreat thee now,<br />Even by the stern and dureful god of war,<br />Whose sinewy neck in battle ne'er did bow,<br />Who conquers where he comes in every jar, 100<br /> Yet hath he been my captive and my slave,<br />And begg'd for that which thou unask'd shalt have.<br />Over my altars hath he hung his lance,<br />His batter'd shield, his uncontrolled crest, 104<br />And for my sake hath learn'd to sport and dance,<br />To toy, to wanton, dally, smile, and jest;<br /> Scorning his churlish drum and ensign red<br />Making my arms his field, his tent my bed. 108<br />"Thus far overrul'd I oversway'd,<br />Leading him prisoner in a red rose chain;<br />Strong-temper'd steel his stronger strength obey'd,<br />Yet was he servile to my coy disdain. 112<br /> Oh be not proud, nor brag not of thy might,<br /> For mast'ring her that foil'd the god of fight.<br />"Touch but my lips with those fair lips of thine,<br />Though mine be not so fair, yet are they red, 116<br />The kiss shall be thine own as well as mine:<br />What see'st thou in the ground? hold up thy head,<br /> Look in mine eyeballs, there thy beauty lies;<br />Then why not lips on lips, since eyes in eyes? 120<br />"Art thou asham'd to kiss? then wink again,<br />And I will wink; so shall the day seem night.<br />Love keeps his revels where there are but twain;<br />Be bold to play, our sport is not in sight, 124<br />These blue-vein'd violets whereon we lean<br /> Never can blab, nor know nor what we mean.<br /> "The tender spring upon thy tempting lip 127<br />Shows thee unripe; yet mayst thou well be tasted,<br />Make use of time, let not advantage slip;<br />Beauty within itself should not be wasted,<br /> Fair flowers that are not gather'd in their prime<br /> Rot, and consume themselves in little time. 132<br />"Were I hard-favour'd, foul, or wrinkled old,<br />Ill-nurtur'd, crooked, churlish, harsh in voice,<br />O'erworn, despised, rheumatic, and cold,<br />Thick-sighted, barren, lean, and lacking juice, 136<br /> Then mightst thou pause, for then I were not for thee;<br /> But having no defects, why dost abhor me?<br /> "Thou canst not see one wrinkle in my brow, 139<br />Mine eyes are grey and bright, and quick in turning;<br />My beauty as the spring doth yearly grow,<br />My flesh is soft and plump,<br />my marrow burning,<br /> My smooth moist hand, were it with thy hand felt,<br /> Would in thy palm dissolve, or seem to melt. 144<br />"Bid me discourse, I will enchant thine ear,<br />Or like a fairy, trip upon the green,<br />Or like a nymph, with long dishevell'd hair,<br />Dance on the sands, and yet no footing seen. 148<br />Love is a spirit all compact of fire,<br /> Not gross to sink, but light, and will aspire.<br />"Witness this primrose bank whereon I lie; 151<br />These forceless flowers like sturdy trees support me;<br />Two strengthless doves will draw me through the sky,<br />From morn till night, even where I list to sport me.<br /> Is love so light, sweet boy, and may it be<br /> That thou shouldst think it heavy unto thee? 156<br />"Is thine own heart to thine own face affected?<br />Can thy right hand seize love upon thy left?" }

# 13  Summary of Revisions

Briefly describe who critiqued your document (e.g., instructor, peer, friend) and provided suggestions for improvement, and how you have incorporated the suggestions and revised the document.

Include information on: who critiqued your document, what suggestions were made, and how you incorporated the suggestions and revised the document.

Answer:

After we complete the website to communicate with database, we had talked to Dr.Gudivada for some suggestions and all we need to improve is to add in the css style make the website look more pleasant.

# 14  Metacognitive Reflection

We learn how to learn through metacognitive reflection by actively planning, monitoring, and evaluating our own thinking and learning.

Learning how to learn involves going beyond the cognitive and into the realm of the metacognitive. In the context of this assignment, cognitive part is the development of the data-intensive application. Metacognitive part refers to the strategies, techniques, and tools you have used to accomplish these tasks.

Perform metacognitive reflection on this assignment by answering the following questions:

1. Did I solve the right problem?

   Yes, we did.

2. Did I solve the problem right?

   Yes, we did.

3. How did I approach solutions to the problems?

   Setting up database and website to show deliverable.

4. What strategies and techniques did I draw upon?

   We drew upon prior knowledge when cleaning out data. We use methods learned in class to load data into out database .

5. Did I learn a new strategy in completing this assignment? If so, how is it different from and similar to the repertoire of techniques that I have already acquired?

   Yes we learned how to interface with a PostgreSQL database from an application program, We also learned how to create a simple web app with angularJS and flask and have it interface with our database

6. Any other information you may wish to add · · ·

   This assignment provided a great opprotunity to apply concepts learned in class.

## 15   Self-assessment

You need to assign a grade for this assignment yourself. Use the rubric listed below to come up with a score. The instructor will also assign a score. Without this section, assignment will be returned with a score of 0.

The first two traits correspond to writing and the remaining ones relate to domain aspects of the project.

| Perf Level / Trait | Poor | Fair | Good | Outstanding |
|---|---|---|---|---|
| *Diction* | Chooses non-technical vocabulary that inadequately conveys the intended meaning of the communication. | Chooses technical vocabulary that conveys the intended meaning of the communication. | Chooses appropriate, and varied technical, vocabulary that conveys the intended meaning of the communication. | Chooses lively, precise, technical, and compelling vocabulary and skillfully communicates the message. |
| *Communication Style* | Has only a few (but noticeable) errors in style, mechanics, or other issues that might distract from the message. | Is virtually free of mechanical, stylistic or other issues. | Uses complex and varied sentence styles, concepts, or visual representations. | Creates a distinctive communication style by combining a variety of materials, ideas, or visual representations. |
| *Application Selection* | Not a data-intensive application. | Application is somewhat data-intensive | Application is data-intensive but limited access to domain expertise. | Application is data-intensive with adequate access to domain expertise. |
| *Use-cases* | Less than 50% of the use-cases are identified, and documented poorly. | Over 75% of the uses-cases are identified and documented using a standard template. | All the use-cases are identified, but detail is missing for some use-cases. | All the use-cases are identified, well-documented using a standard template, and verified against application requirements. |
| *Data Tasks* | Inputs, outputs, and possible error conditions are documented for less than 50% of data tasks. | Inputs, outputs, and possible error conditions are documented for less than 75% of data tasks. | Inputs, outputs, and possible error conditions are documented for all data tasks. | Inputs, outputs, and possible error conditions are documented for all data tasks. Processing logic (or high-level algorithms) for transforming inputs into outputs is also described. |
| *Transactions and Queries* | Less than 50% of the transactions and queries are identified and described. | Less than 75% of the transactions and queries are identified and described. | All the transactions and queries are identified and described. | All the transactions and queries are identified and described including their frequency of execution. |

| Perf Level / Trait | Poor | Fair | Good | Outstanding |
| --- | --- | --- | --- | --- |
| *Data Models* | Only conceptual data model is described in detail. Cursory treat of logical data model. Physical data model design is missing. | Conceptual and logical data models are described in detail. Physical data model design is missing. | Conceptual, logical, and physical data models are described completely and precisely. | Conceptual, logical, and physical data models are described completely and precisely. Database normalization based on functional dependencies is discussed in detail. |
| *Creation and Loading* | SQL scripts are written and executed to create the database and load the data. Data in the database is trivial in size. | SQL scripts are written and executed to create the database and load the data. Data in the database is moderate in size. | Conceptual, logical, and physical data models are described completely and precisely. Data in the database is huge in size – in the order of millions of rows. | Conceptual, logical, and physical data models are described completely and precisely. Data in the database is huge in size – in the order of millions of rows. Detail evidence is provided on how referential integrity constraints are resolved. |
| *Implementing Transactions and Queries* | Less than 50% of the transactions and queries are implemented. | Less than 75% of the transactions and queries are implemented. | All the transactions and queries are implemented; run and execute correctly. | All the transactions and queries are implemented; run and execute correctly. There is also written evidence that transactions and queries are tested. |
| *Revisions* | Only peer or instructor/grader feedback is solicited, but not incorporated. | Both peer and instructor/grader feedback is solicited but not incorporated. | Both peer and instructor/grader feedback is solicited and incorporated. | Both peer and instructor/grader feedback solicited and incorporated. Evidence is presented to show how the feedback improved the document. |
| *Meta-cognitive Reflection* | Not performed. | Is shallow and incomplete. | Is complete but not thorough. | Is complete and thorough. |

Use the following table to score your solution. Circle the appropriate number in each row. For example, to circle 4, use the LATEX markup code `\circled{4}`, which produces ④.

| Perf Level Trait | Poor | Fair | Good | Outstanding |
|---|---|---|---|---|
| Diction | 2 | 3 | 4 | 5 |
| Communication Style | 2 | 3 | 4 | 5 |
| Application Selection | 4 | 6 | 10 | 15 |
| Use-cases | 4 | 6 | 8 | 10 |
| Data Tasks | 4 | 6 | 8 | 10 |
| Transactions and Queries | 4 | 6 | 8 | 10 |
| Data Models | 4 | 6 | 8 | 10 |
| Creation and Loading | 4 | 6 | 8 | 10 |
| Implementing Transactions and Queries | 4 | 6 | 8 | 10 |
| Revisions | 4 | 6 | 8 | 10 |
| Meta-cognitive Reflection | 2 | 3 | 4 | 5 |

Total score: 100 / 100.

## A  Code Listings

```
1   #Code for extracting data from Gutenberg corpus
2   #Directory set up:
3   #This program need:
4   #    input folder: for input files
5   #    output folder: for output files
6   #    duplicate folder: for duplicate files
7   #    organized folder: for organized files
8   #        txt folder: for organized text files
9   #        other foler: for organized other files
10  #        orgDuplicate folder: for organized duplicate files
11  #    extract folder: for extracted data store
12  #        _data folder: extracted information
13  #        encodeErr folder: encode error files
14  #        exist folder: data exist
15  #        nonExist folder: data could not be determine author,
        title, language exist or not
16
17
18  #All user need to do is put input file insdie input folder and
        run the program
19
20  '''
21  directory structure:
22  .
23  +--extract.py
24  +--input
25  +--output
26  +--duplicate
27  +--organized
28      +--txt
29      +--other
30      +--orgDuplicate
31  +--extract
32      +--_data
33      +--exist
34          +--missInfo
35      +--nonExist
36
37  '''
38  #
    ===============================================================================

39
```

```
40  #imports
41  import os
42  import zipfile
43
44  def makeDir():
45      os.mkdir("input")
46      os.mkdir("output")
47      os.mkdir("duplicate")
48      os.mkdir("organized")
49      os.mkdir("organized/txt")
50      os.mkdir("organized/other")
51      os.mkdir("organized/orgDuplicate")
52      os.mkdir("extract")
53      os.mkdir("extract/_data")
54      os.mkdir("extract/exist")
55      os.mkdir("extract/exist/missInfo")
56      os.mkdir("extract/nonExist")
57  #
     ================================================================================

58
59  #Use Recursion to loop into most inner file if it's zip unzip
       and step back one directory and go to second, so on.....
60  def unzipFiles(inputDir,outputDir,fileNameList,duplicateCounter
       ,otherFile,duplicateFile,duplicateDir):
61      fileList = os.listdir(inputDir)
62
63      for fileName in fileList:
64          #if it's directory
65          if os.path.isdir(os.path.join(inputDir,fileName)):
66              newInputDir = os.path.join(inputDir, fileName)
67              unzipFiles(newInputDir,outputDir,fileNameList,
                   duplicateCounter,otherFile,duplicateFile,
                   duplicateDir)
68          #if it's zip files
69          elif ".zip" in fileName:
70              #if file is duplicate unzip to duplicate
71              if fileName in fileNameList:
72                  duplicateCounter+=1
73                  duplicateFile.write(fileName+" "+str(
                       duplicateCounter)+" \n")
74                  unzipDir = os.path.join(inputDir,fileName)
75                  zip_ref = zipfile.ZipFile(unzipDir, 'r')
76                  zip_ref.extractall(duplicateDir)
77                  print(fileName+" DONE")
```

15

```
78                      zip_ref.close()
79                  #else unzip file to output folder
80                  else:
81                      fileNameList.append(fileName)
82                      unzipDir = os.path.join(inputDir,fileName)
83                      zip_ref = zipfile.ZipFile(unzipDir, 'r')
84                      zip_ref.extractall(outputDir)
85                      print(fileName+"␣DONE")
86                      zip_ref.close()
87          #else not zip file not directory, record it
88          else:
89              print("other␣files")
90              otherFile.write(fileName+"\n")
91  #
    ===========================================================================


92
93  #Organize text file into organized folder, categorize file into
        text or other.
94  def organize(inputDir,categorizedTxtDir,categorizedOtherDir,
      categorizedDuplicateDir,fileNameList):
95      outputFileList = os.listdir(inputDir)
96
97      for fileName in outputFileList:
98          if os.path.isdir(os.path.join(inputDir,fileName)):
99              newInputDir = os.path.join(inputDir, fileName)
100             organize(newInputDir,categorizedTxtDir,
                    categorizedOtherDir,categorizedDuplicateDir,
                    fileNameList)
101         elif ".txt" or ".TXT" in fileName:
102             fileName = fileName.replace(".TXT",".txt")
103             if fileName in fileNameList:
104                 currentFile = os.path.join(inputDir,fileName)
105                 newFile = os.path.join(categorizedDuplicateDir,
                        fileName)
106                 os.rename(currentFile,newFile)
107                 print(fileName+"␣DONE")
108             else:
109                 fileNameList.append(fileName)
110                 currentFile = os.path.join(inputDir,fileName)
111                 newFile = os.path.join(categorizedTxtDir,
                        fileName)
112                 os.rename(currentFile,newFile)
113                 print(fileName+"␣DONE")
114         else:
```

16

```
115            currentFile = os.path.join(inputDir,fileName)
116            newFile = os.path.join(categorizedOtherDir,fileName
                  )
117            os.rename(currentFile,newFile)
118            print(fileName+"␣DONE")
119
120  #
     =============================================================================

121
122  def extractData(categorizedTxtDir,extractExistDir,
     extractNonExistDir,extractDataDir,extractExistMissDir,
     encodeErrDir):
123   categorizedTxtFileList = os.listdir(categorizedTxtDir)
124
125   for fileName in categorizedTxtFileList:
126         title =""
127         author =""
128         release =""
129         language =""
130         text =""
131         asciiErr = False
132         utf8Err = False
133
134         currentFileDir = os.path.join(categorizedTxtDir,
               fileName)
135         currentFile = open(currentFileDir,"r",encoding='ascii')
136         currentFile = open(currentFileDir,"r",encoding='utf-8')
137         try:
138             ascContent = currentFile.read()
139         except:
140             asciiErr=True
141
142         try:
143             utfContent = currentFile.read()
144         except:
145             utf8Err=True
146
147         if asciiErr == False:
148             content = ascContent
149         elif utf8Err == False:
150             content = utfContent
151         else:
152             currentFile = os.path.join(categorizedTxtDir,
                  fileName)
```

```
153            newFile = os.path.join(encodeErrDir,fileName)
154            os.rename(currentFile,newFile)
155            continue
156
157
158        splitFile = content.split("***␣START␣OF␣THIS␣PROJECT␣
                GUTENBERG␣EBOOK")
159
160        if len(splitFile) == 2:
161            info = splitFile[0]
162            text = splitFile[1]
163            splitLines = info.split("\n")
164            for eachLine in splitLines:
165                if "Title:" in eachLine:
166                    title = eachLine.replace("\n","")
167                elif "Author:" in eachLine:
168                    author = eachLine.replace("\n","")
169                elif "Release␣Data:" in eachLine:
170                    release = eachLine.replace("\n","")
171                elif "Language:" in eachLine:
172                    language = eachLine.replace("\n","")
173
174            if title == "" or author == "" or release == "" or
                    language == "":
175                currentFile = os.path.join(categorizedTxtDir,
                        fileName)
176                newFile = os.path.join(extractExistMissDir,
                        fileName)
177                os.rename(currentFile,newFile)
178            else:
179                writeNewFileName = os.path.join(extractDataDir,
                        fileName)
180                writeNewFile = open(writeNewFileName,"w")
181                writeNewFile.write(title+"\n"+author+"\n"+
                        release+"\n"+language+"\n␣Text:␣"+text)
182                writeNewFile.close()
183                currentFile = os.path.join(categorizedTxtDir,
                        fileName)
184                newFile = os.path.join(extractExistDir,fileName
                        )
185                os.rename(currentFile,newFile)
186        else:
187            currentFile = os.path.join(categorizedTxtDir,
                    fileName)
188            newFile = os.path.join(extractNonExistDir,fileName)
```

```
189                os.rename(currentFile,newFile)
190            print(fileName+" Done")
191  #
     =============================================================================


192
193  #main
194  #1
195  currentDir = os.getcwd()
196  inputDir = os.path.join(currentDir, "input")
197  outputDir = os.path.join(currentDir, "output")
198  otherFile = open(os.path.join(currentDir,"other.txt"),"w")
199  duplicateDir = os.path.join(currentDir, "duplicate")
200  duplicateFile = open(os.path.join(currentDir,"duplicate.txt"),"
     w")
201  fileNameList = []
202  duplicateCounter = 0
203  #2
204  categorizedTxtDir = os.path.join(currentDir,"organized","txt")
205  categorizedOtherDir = os.path.join(currentDir,"organized","
     other")
206  categorizedDuplicateDir = os.path.join(currentDir,"organized","
     orgDuplicate")
207  #3
208  extractExistDir = os.path.join(currentDir,"extract","exist")
209  extractExistMissDir = os.path.join(currentDir,"extract","exist"
     ,"missInfo")
210  extractNonExistDir = os.path.join(currentDir,"extract","
     nonExist")
211  extractDataDir = os.path.join(currentDir,"extract","_data")
212  encodeErrDir = os.path.join(currentDir,"extract","encodeErr")
213
214  print("Select your action:")
215  print("1. create file directories")
216  print("2. Unzip all files from input directory")
217  print("3. organize output folder, categorize it into text files
      and other files")
218  print("4. extract data from organized/txt and move file to
     extract/exist or extract/nonExist")
219  print("0. quit")
220
221  userInput = input()
222  userInputInt = int(userInput)
223  if userInputInt == 1:
224      makeDir()
```

```
225  elif userInputInt == 2:
226      unzipFiles(inputDir,outputDir, fileNameList,
             duplicateCounter,otherFile,duplicateFile,duplicateDir)
227  elif userInputInt == 3:
228      otherFile.close()
229      duplicateFile.close()
230      organize(outputDir,categorizedTxtDir,categorizedOtherDir,
             categorizedDuplicateDir,fileNameList)
231  elif userInputInt == 4:
232      extractData(categorizedTxtDir,extractExistDir,
             extractNonExistDir,extractDataDir,extractExistMissDir,
             encodeErrDir)
233  elif userInputInt == 0:
234      exit()
235
236
237  #------------------------------------------------------------#
238
239  #code for loading data in to database
240  import re
241  import os
242  import sys
243
244  import psycopg2 as pg
245  import pandas.io.sql as psql
246  from pprint import pprint
247
248
249  def extractDataFromHeaderLine(dataDescription, line):
250      lineS = re.split('^' + dataDescription, line, maxsplit=1)
251      #print('lineS is: ')
252      #pprint(lineS)
253      if(len(lineS) == 2):
254          return lineS[1].strip().lower()
255      else:
256          return ""
257
258
259
260  #for well formed docs
261  def getData(file):
262      title =""
263      author =""
264      release =""
265      language =""
```

```
266        text =""
267        gutenberId = ""
268        charSetEncode=""
269        asciiErr = False
270        utf8Err = False
271        defaultErr = False
272
273        content = None
274
275        currentFileAscii = open(file,"r",encoding='ascii')
276        currentFileUtf8 = open(file,"r",encoding='utf-8')
277        currentFileDefult = open(file,"r")
278
279        try:
280            ascContent = currentFileAscii.read()
281        except:
282            asciiErr=True
283        try:
284            utfContent = currentFileUtf8.read()
285        except:
286            utf8Err=True
287
288        try:
289            defaultContent = currentFileDefult.read()
290        except:
291            defaultErr = True
292
293        if asciiErr == False:
294            content = ascContent
295        elif utf8Err == False:
296            content = utfContent
297        elif defaultErr == False:
298            content = defaultContent
299        else:
300            currentFileAscii.close()
301            currentFileUtf8.close()
302            currentFileDefult.close()
303
304            print('error reading file: ' + file)
305
306        #extract the data
307
308        headerAndContent = re.split('\*\*\*', content, maxsplit=1)
309        #print(headerAndContent[0])
310
```

```
311     header = headerAndContent [0]
312     fullText = headerAndContent [1]
313     #print ( fullText )
314
315     headerSplit = header . split ( '\n' )
316
317     for line in headerSplit :
318         if "Title :" in line :
319             title = extractDataFromHeaderLine ( 'Title :' , line )
320         elif "Author :" in line :
321             author = extractDataFromHeaderLine ( 'Author :' , line )
322
323         elif "Release␣Date :" in line :
324             release = extractDataFromHeaderLine ( 'Release␣Date :'
                    , line )
325             if ( '[' in release ):
326                 releaseS = release . split ( '[' )
327                 release = releaseS [0] . strip ()
328
329         elif "Language :" in line :
330             language = extractDataFromHeaderLine ( 'Language :' ,
                    line )
331
332     '''
333     print ( 'The title is : ' + title )
334     print ( 'Author : ' + author )
335     print ( 'Date : ' + release )
336     #pprint ( release )
337     print ( 'Language : ' + language )
338     '''
339
340     return title , author , release , language , fullText
341
342
343 def processAuthorName ( authorName ):
344     firstName = ''
345     middleName = ''
346     lastName = ''
347     suffix = ''
348     prefix = ''
349     nameL = authorName . split ()
350
351     i = 0
352     while ( i < len ( nameL )):
353         if ( ')' in nameL [i] or '(' in nameL [i ]):
```

```
354            del nameL[i]
355            i -=1
356        i += 1
357
358    nameLen = len(nameL)
359
360    #check for prefix
361    comPrefixes = ['mr', 'mrs', 'miss', 'sir', 'lord', 'ms']
362    comSuffixes = ['sr', 'jr', 'ii', 'iii', 'iv', 'v']
363
364
365    if(nameLen > 1):
366        if(nameL[0].strip('.') in comPrefixes):
367            prefix = nameL[0]
368            del nameL[0]
369            nameLen = len(nameL)
370
371    if(nameLen > 1):
372        if(nameL[nameLen - 1].strip('.') in comSuffixes):
373            suffix = nameL[nameLen - 1]
374            del nameL[nameLen - 1]
375            nameLen = len(nameL)
376
377    if(nameLen == 0):
378        firstName = 'anonymous'
379    elif(nameLen == 1):
380        firstName = nameL[0]
381    elif(nameLen == 2):
382        firstName = nameL[0]
383        lastName = nameL[1]
384
385    elif(nameLen == 3):
386        firstName = nameL[0]
387        middleName = nameL[1]
388        lastName = nameL[2]
389    else:
390        firstName = nameL[0]
391        lastName = nameL[nameLen - 1]
392        del nameL[nameLen - 1]
393        del nameL[0]
394        middleName = ' '.join(nameL)
395
396    '''
397    print('first name: ' + firstName)
398    print('middle name: ' + middleName)
```

```
399     print('last name: ' + lastName)
400     print('suffix: ' + suffix)
401     print('prefix: ' + prefix)
402     '''
403
404     return firstName, middleName, lastName, suffix, prefix
405
406
407  ########## Database Querys
408
409  def getBookByPrimaryKey(gutenbergId, cur):
410      cur.execute('Select * From public."Book" Where gutenberg_id
              = %s', (gutenbergId,))
411      result = cur.fetchall()
412      #print(result)
413      return result
414
415  def bookIsInDatabase(gutenbergId, conn):
416      result = getBookByPrimaryKey(gutenbergId, conn)
417      #print('book is:')
418      #print(result)
419      #sys.exit()
420
421      if(len(result) > 0):
422          return True
423      else:
424          return False
425
426  def getAuthorByName(firstName, middleName, lastName, suffix,
          prefix, cur):
427      cur.execute('Select * From public."Author" Where first_name
              = %s And middle_name = %s And last_name = %s and suffix
              = %s and prefix = %s',
428                      (firstName, middleName, lastName, suffix,
                          prefix))
429
430      result = cur.fetchall()
431      #print(result)
432      return result
433
434  def authorIsInDatabase(firstName, middleName, lastName, suffix,
          prefix, cur):
435      result = getAuthorByName(firstName, middleName, lastName,
              suffix, prefix, cur)
436
```

```
437      if(len(result) > 0):
438          return True
439      else:
440          return False
441
442
443  def getWrittenBy(gutenbergId, authorId, cur):
444      cur.execute('Select␣*␣From␣public."Written_By"␣Where␣
             author_id␣=␣%s␣And␣gutenberg_id␣=␣%s', (authorId,
             gutenbergId))
445      result = cur.fetchall()
446      return result
447
448
449  def writtenByRelationInDatabase(gutenbergId, authorId, cur):
450      result = getWrittenBy(gutenbergId, authorId, cur)
451
452      if(len(result) > 0):
453          return True
454      else:
455          return False
456
457
458  def insertIntoDatabase(gId, title, release, language, author,
         fullText, conn, cur, brokenFiles):
459
460      #gId = str(gId) + '-' + language
461      try:
462
463
464          if(not bookIsInDatabase(gId, cur)):
465              try:
466                  #inset book
467                  cur.execute('INSERT␣INTO␣public."Book"␣(
                         gutenberg_id,release_date,full_text,language,
                         title)␣VALUES␣(%s,%s,%s,%s,%s);',
468                              (gId, release, fullText, language,
                                 title))
469                  #print('loaded book')
470                  conn.commit()
471              except pg.DataError:
472                  conn.rollback()
473                  cur.execute('INSERT␣INTO␣public."Book"␣(
                         gutenberg_id,release_date,full_text,language,
                         title)␣VALUES␣(%s,%s,%s,%s,%s);',
```

```
474                              (gId, None, fullText, language,
                                     title))
475                 print('loaded␣book␣wth␣null␣date:␣' + str(gId))
476                 conn.commit()

477

478         else:
479             print('book␣is␣in␣database,␣gid:␣' + str(gId))

480

481         if('␣and␣' in author):
482             authorNameS = re.split('␣and␣', author)
483         elif(',' in author):
484             authorNameS = re.split(',', author)
485         else:
486             authorNameS = [author]

487

488         authorNameS = [name.strip() for name in authorNameS if(
             name != '' or name != None)]

489

490         if(len(authorNameS) > 1):
491             pprint(authorNameS)
492         for author in authorNameS:

493

494             firstName, middleName, lastName, suffix, prefix =
                 processAuthorName(author)
495             #TODO fix rest of querys
496             if(not authorIsInDatabase(firstName, middleName,
                 lastName, suffix, prefix, cur)):

497

498             #inset author
499                 cur.execute('insert␣into␣public."Author"␣(
                     first_name,last_name,middle_name,suffix,
                     prefix)␣VALUES␣(%s,␣%s,␣%s,␣%s,␣%s)␣returning
                     ␣author_id',
500                              (firstName, lastName, middleName,
                                 suffix, prefix))
501                 #print('loaded author')
502                 authorId = cur.fetchone()
503                 conn.commit()
504                 #print('auth id is: ' + str(authorId[0]))
505             else:

506

507                 authorId = getAuthorByName(firstName,
                     middleName, lastName, suffix, prefix, cur)
                     [0][0]
508                 print('author␣is␣in␣database,␣author_id:␣' +
```

```
                             str(authorId))
509
510              if(not writtenByRelationInDatabase(gId, authorId,
                     cur)):
511                  cur.execute('INSERT␣INTO␣public."Written_By"␣(
                         author_id,␣gutenberg_id)␣VALUES␣(%s,␣%s)', (
                         authorId, gId))
512                  conn.commit()
513                  #print('loaded written by')
514              else:
515                  wb = getWrittenBy(gId, authorId, cur)
516                  #print(wb[0])
517                  #print('written by relation in database: ' +
                         str(wb[0][0]), + ', ' + str(wb[0][1]))
518
519
520      #      #inset written_by
521      except Exception:
522          #pass
523          brokenFiles.append(gId)
524          print(Exception)
525          print(gId)
526
527
528
529  #              MAIN                #
530
531  dbname = "gutenburg"
532  dbhost = "localhost"
533  dbport = "5432"
534  dbuser = "postgres"
535  dbpassword = "postgres"
536
537
538  conn = pg.connect(database=dbname, host=dbhost, port=dbport,
         user=dbuser, password=dbpassword)
539  cur = conn.cursor()
540  #dirToProcess = '/Users/edwardsja15/desktop/gutenberg/load/
         extract/_data'
541  #dirToprocessWindows = '/home/reilly/database/load/_data'
542
543  dirToProcess = '/home/reilly/database/load/_data'
544
545
546  fileList = os.listdir(dirToProcess)
```
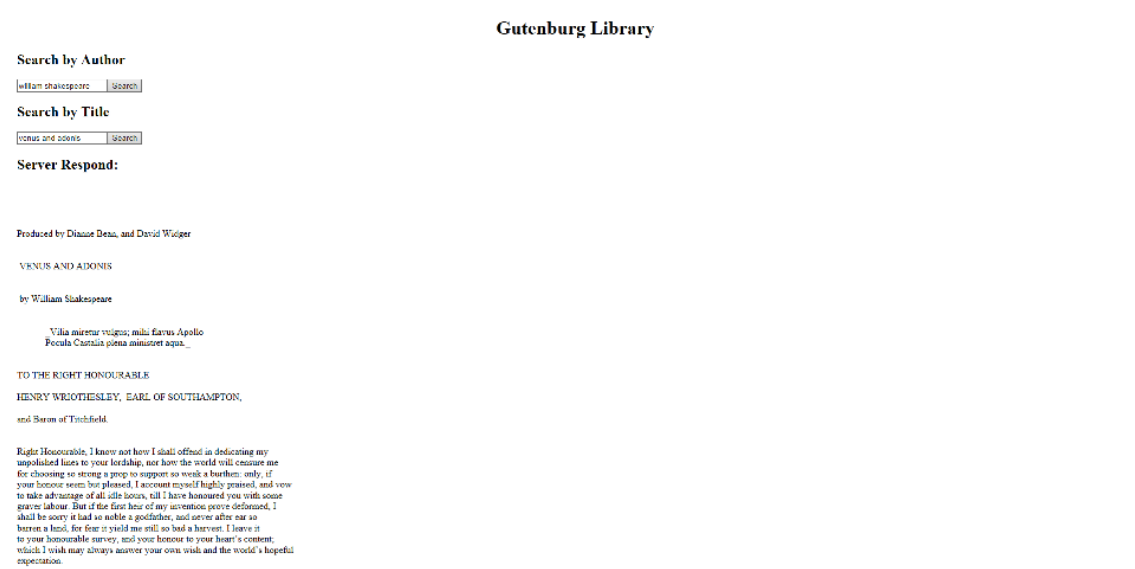
```python
547  #print(fileList[:10])
548  i = 0
549  brokenFiles = []
550  invalidIds = []
551  for filename in fileList:
552      #can load
553      if('.txt' in filename):
554          fullfilename = os.path.join(dirToProcess, filename)
555
556          title, author, release, language, fullText = getData(
                 fullfilename)
557
558          if(language == 'en'):
559              language = 'english'
560
561          filenameS = filename.split('.')
562          isValid = False
563
564          if('-' in filenameS[0]):
565              idSplit = filenameS[0].split('-')
566              try:
567                  int(idSplit[1])
568                  isValid = True
569              except ValueError:
570                  isValid = False
571
572                  invalidIds.append(filename)
573
574          else:
575              isValid = True
576
577          gutenbergId = filenameS[0]
578          print('filename is: ' + filename)
579
580          if(isValid):
581              insertIntoDatabase(gutenbergId, title, release,
                     language, author, fullText, conn, cur,
                     brokenFiles)
582          #i += 1
583          #if(i > 1050):
584          #    break
585  print('invalid ids:')
586  pprint(invalidIds)
587  print('broken files:')
588  pprint(brokenFiles)
```

```
589
590  conn.close()
```

# B  Test Cases



# C  Other

Full code listings can be seen at:
https://github.com/reiman2222/gutenberg-corpus-for-learning-SQL-and-IR