

1: Cassandra

Application = Simple Facebook Profiles

- Facebook sells the service of self-expression via a customizable profile that keeps a structured format.
- The data represents peoples' experiences (photos), thoughts (status updates) and preferences (likes)

Why Cassandra is the best fit

- Cassandra is most suitable for content management. Facebook profiles have a lot of content that is in the form of structured data, such as photo albums and status updates, which are almost like a mini-blogging system. More reasons to use Cassandra include:
- There is a vast population of users with Facebook profiles, so the fact that Cassandra handles data across many servers will be very beneficial. This decentralization feature is also related to scalability, so when new machines are added, service can increase linearly.
- Arelational DBMS is not the best option because it will be slower in terms of read-write time, and not as flexible for storing and managing profile content.
- Using HBase or Mongo would be a very similar experience, but Cassandra's column family and keyspace models are more suited for structured document storage. Hbase would encourage heavy use of Hadoop, which may not be wanted, and Mongo would probably be less scalable and slower.

Column families:

- Notes: the bold words are keys. All values have a timestamp that is not shown for simplicity.

| | | | | |
|----------------|--------------|----------|--|------------|
| users | | | | |
| taggr | name | password | email | location |
| | Tagg Ridler | batman11 | tagg@yahoo.com | Boulder,CO |
| mileyc | name | password | location | |
| | Miley Cyrus | cantstop | Hollywood, CA | |
| baracko | name | password | email | |
| | Barack Obama | potus | bo@whitehouse.gov | |

| | | |
|----------------|---------|--------------|
| friends with | | |
| taggr | mileyc | |
| mileyc | baracko | taggr |
| baracko | mileyc | |

| | | | |
|----------------|-----------------|--------|---------|
| status updates | | | |
| taggr | body | tags | photo |
| | I love Miley! | mileyc | 39hh083 |
| mileyc | body | | |
| | Where am I? | | |
| baracko | body | photo | |
| | Bye healthcare! | 09guq | |

CQL queries

```
1 INSERT INTO users (name, password, location)
2 VALUES ("Lana Del Rey", "tropico", "Cony Island, NY")
3
4 SELECT COUNT (*)
5 FROM users
6 WHERE location = "Hollywood, CA"
7
8 CREATE COLUMNFAMILY Users (
9     KEY uuid PRIMARY KEY,
```

```

10     name text,
11     password text,
12     email text,
13     location text );
14 CREATE COLUMNFAMILY FriendsWith (
15     KEY uuid PRIMARY KEY,
16     Friends text );
17 CREATE COLUMNFAMILY StatusUpdates (
18     KEY uuid PRIMARY KEY,
19     body text,
20     tags text,
21     photo blob );

```

2: XML language

Language Definition

- This model is general enough to represent any system that manages profiles with a picture, posts that can contain pictures, and friends or connections to other users with profiles.
- This model supports the newPost operation which allows a user to create a new post, and the suggestedFriends operation which suggests new connections for users based on the friends of their friends.

```

1  <?xml version="1.0"?>
2  <profile>
3      <name></name>
4      <profilePhoto></profilePhoto>
5      <newPost>
6          <source></source>
7          <tagBox></tagBox>
8          <inputText></inputText>
9          <inputPhoto></inputPhoto>
10     </newPost>
11     <posts>
12         <post>
13             <time></time>
14             <location></location>
15             <tags>
16                 <tag></tag>
17             </tags>
18             <body>
19                 <photo></photo>
20                 <text></text>
21             </body>
22         </post>
23     </posts>
24     <friends>
25         <friend>
26             <name></name>
27             <profilePhoto></profilePhoto>
28             <link></link>
29         </friend>
30     </friends>
31     <suggestedFriends>
32         <source></source>
33         <inputFriends>
34             <friendID></friendID>
35         </inputFriends>
36         <userID></userID>
37     </suggestedFriends>
38 </profile>

```

XML Schema

- This does not include types for attributes of operations.
- Photos are represented as a URI (string)

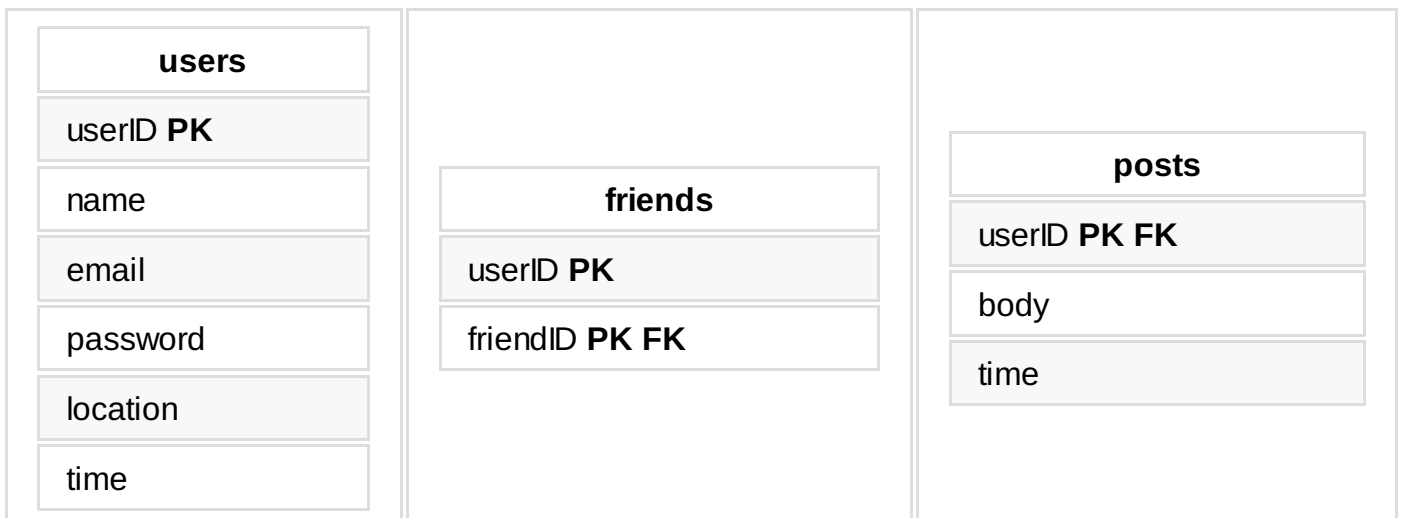
```

1 <?xml version="1.0"? encoding="ISO-8859-1" ?>
2 <xs:schema xmlns:xs="http://www.w3.org/2001/XMLSchema">
3     <xs:element name="name" type="xs:string"/>
4     <xs:element name="photo" type="xs:string"/>
5     <xs:element name="source" type="xs:string"/>
6     <xs:element name="tag" type="xs:string"/>
7     <xs:element name="time" type="xs:dateTime"/>
8     <xs:element name="location" type="xs:string"/>
9     <xs:element name="body" type="xs:string"/>
10    <xs:element name="link" type="xs:string"/>
11    <xs:element name="userID" type="xs:int"/>
12    <xs:element name="friendID" type="xs:int"/>
13
14    <xs:element name="tags">
15        <xs:complexType>
16            <xs:list ref="tag"/>
17        </xs:complexType>
18    </xs:element>
19
20    <xs:element name="post">
21        <xs:complexType>
22            <xs:sequence>
23                <xs:element ref="time"/>
24                <xs:element ref="location"/>
25                <xs:element ref="tags"/>
26                <xs:element ref="photo"/>
27            </xs:sequence>
28        </xs:complexType>
29    </xs:element>
30
31    <xs:element name="friends">
32        <xs:complexType>
33            <xs:sequence>
34                <xs:list ref="friendID"/>
35            </xs:sequence>
36        </xs:complexType>
37    </xs:element>
38
39 </xs:schema>

```

3: OLAP and Relational DB

Relational Schema



- Functional Dependencies: all values in each table are functionally determined by the primary key.
- Multivalued Dependencies: the friends table. The presence of a row (userID, friendID) requires another row (userID, friendID) that is the reverse of the first row.

OLTP Queries

Create a new user and automatically make them friends with the original user (MySpace model):

```

1 INSERT INTO users (name, email, password, location)
2 VALUES ("Kanye West", "ye@thebomb.com", "yeezus", "Chicago, IL");
3
4 INSERT INTO friends VALUES (1, 2);

```

A user can add the person as a friend who lives in Santa Barbara, CA with the most recent post:

```

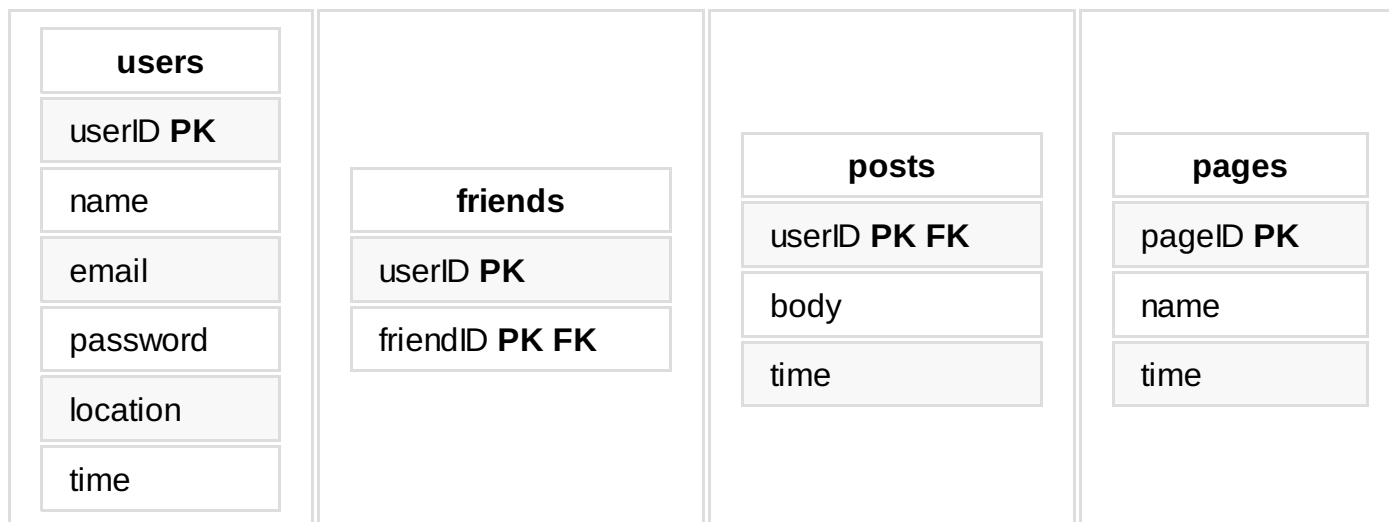
1 INSERT INTO friends VALUES (1,
2 (
3     SELECT userID
4     FROM users u JOIN posts p
5     ON u.userID = p.userID
6     WHERE u.location = "Santa Barbara, CA"
7     ORDER BY p.time LIMIT 1
8 ));

```

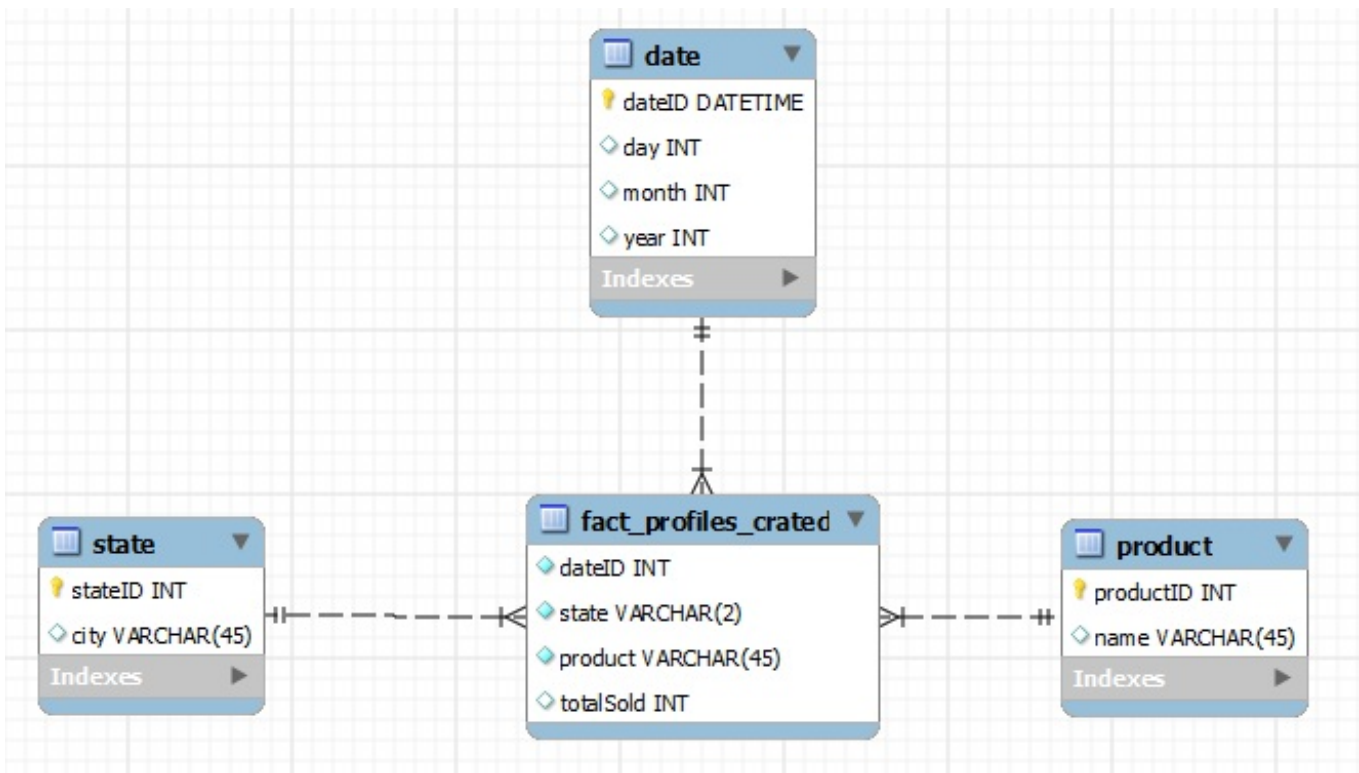
Data Warehousing

Schemas

New relational schema to make this work (I am adding a page table, in reference to facebook pages vs profiles)



Star Schema: for explanation of how it relates to the above schema, see below.



Queries for Trends over Time

- Summarizes "sales" over time by state

```

1 SELECT dateID, state, SUM(totalSold)
2 FROM fact_profiles_created
3 GROUP BY state, dateID
    
```

- Summarizes "sales" of certain product type over time by state

```

1 SELECT product, dateID, state, SUM(totalSold)
2 FROM fact_profiles_created
3 GROUP BY product, dateID, state
    
```

- Summarizes "sales" by product type over a particular month in California

```

1 SELECT product, SUM(totalSold)
2 FROM fact_profiles_created f NATURAL JOIN date d
3 WHERE d.month = 8 AND f.state = "CA"
4 GROUP BY product
    
```

Explanation of My Data Warehousing

- The sales table in the star schema is aggregating the number of times a new profile or page was created in the application. The state table is storing the information related to the location of the profiles or pages that were created, while the product table is storing whether the thing created was a profile or a page. The date table is keeping track of the timestamp representing the time of creation of a profile or page.
- The aggregation described above is done so summary queries can be written to analyze the data relating to the usage of the facebook-like app, and reveal trends over time. This process is very different from online transaction processing (OLTP) because its only purpose is data analysis. The data warehousing system should not modify records in the OLTP system in any way; it should only analyze them. Relatively, the OLTP system should not be used to analyze any of its own data; it should only create, read, update, or delete it.

4: Assertions and Inferences

Modified XML Language:

```

1  <?xml version="1.0"?>
2  <profile>
3      <home></home>
4      <name></name>
5      <profilePhoto></profilePhoto>
6      <newPost>
7          <source></source>
8          <tagBox></tagBox>
9          <inputText></inputText>
10         <inputPhoto></inputPhoto>
11     </newPost>
12     <totalPosts></totalPosts>
13     <posts>
14         <post>
15             <time></time>
16             <location></location>
17             <tags>
18                 <tag>
19                     <friendID></friendID>
20                     <taggedBy>
21                         <userID></userID>
22                     </taggedBy>
23                 </tag>
24             </tags>
25             <body>
26                 <photo></photo>
27                 <text></text>
28             </body>
29         </post>
30     </posts>
31     <friends>
32         <friend>
33             <name></name>
34             <profilePhoto></profilePhoto>
35             <link></link>
36         </friend>
37     </friends>
38     <suggestedFriends>
39         <source></source>
40         <inputFriends>
41             <friendID></friendID>
42         </inputFriends>
43         <userID></userID>
44     </suggestedFriends>
45 </profile>

```

RDF Triples / Assertions

- a post by mileyc should be tagged by mileyc
 - `www.myface.com/mileyc/posts/1432324 <taggedBy> www.myface.com/mileyc`
- one of baracko's suggested friends should be taggr
 - `www.myface.com/baracko <suggestedFriends> www.myface.com/taggr`
- any profile's home route should be the home page
 - `www.myface.com/profile/ <home> www.myface.com`
- a new post's time should be the current time when it is created
 - `www.myface.com/lanar/posts/new <time> 1/1/2014`
- a new post's location should be the location of the person who owns the post
 - `www.myface.com/kanyew/posts/new <location> Chicago, IL`
- taggr's profile photo should be a photo he owns and in his 'profile' route
 - `www.myface.com/taggr <profilePhoto> www.myface.com/taggr/photos/profile`
- the name on someone's profile should match their actual identity
 - `www.myface.com/mileyc <name> Miley Cyrus`
- the input photo for a new post should come from someone's hard drive
 - `www.myface.com/lanar/posts/new <inputPhoto> C:\Users\LanaDe1Rey\Photos\`
- taggr's friends should include mileyc

- `www.myface.com/taggr <friends> www.myface.com/mileyc`
- the number of posts someone has made should be an integer representing their total number of posts
 - `www.myface.com/mileyc <totalPosts> www.metrics.com/myface/posts/mileyc`

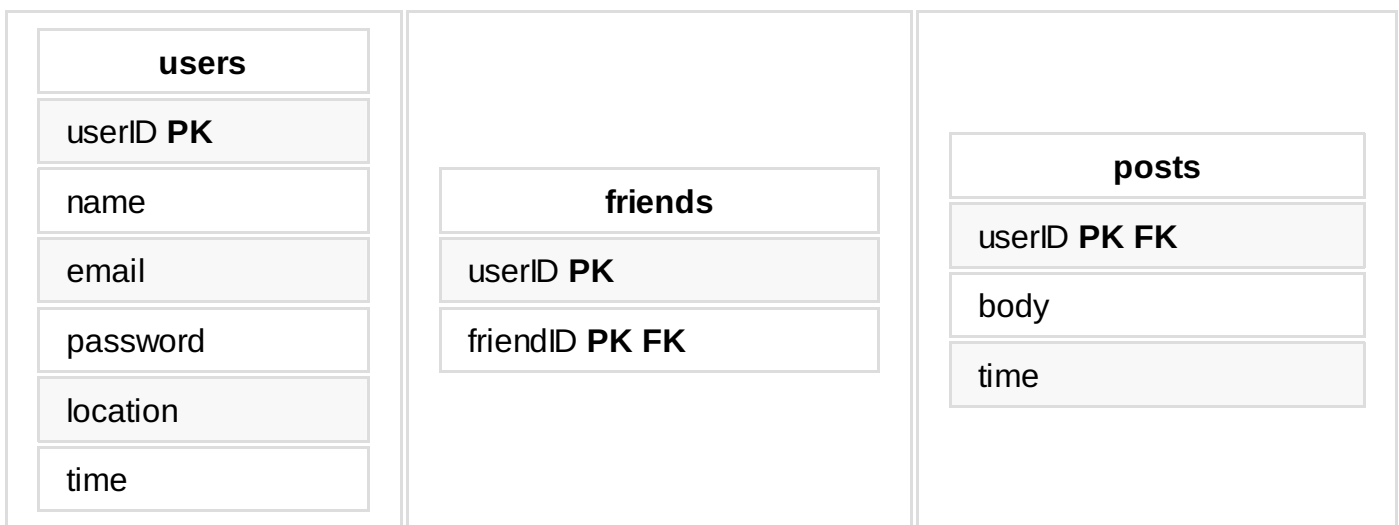
Inferences

- All of the left values in the triples are a route in the app, while all the right values are either a string, local file, or a different route within the app. Thus I can infer that many of the data elements in profiles do not rely on outside web services.
- All of the right values that relate to documents saved to the database are routes within the app. Thus I can infer that many of the documents that make up a profile are contained in the app's document database.
 - A photo that is part of a new post that hasn't been saved yet is a file on someone's hard drive, but a profile photo that has already been saved as a profile document has a route contained within the app's domain.

5: Distributed Databases

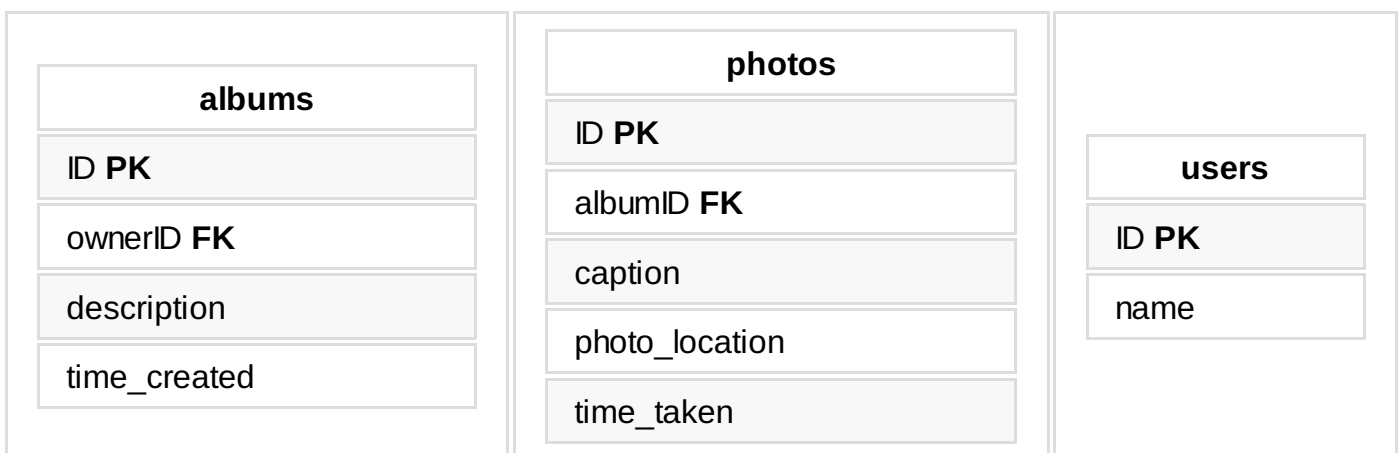
Schemas

Original Profile Schema



- Functional Dependencies: all values in each table are functionally determined by the primary key.
- Multivalued Dependencies: the friends table. The presence of a row (userID, friendID) requires another row (userID, friendID) that is the reverse of the first row.

Pictures App Schema (like Instagram)



- Functional Dependencies: all values in each table are functionally determined by the primary key.
- Multivalued Dependencies: none

Q&A App Schema



| | | |
|---------------|---------------|--------|
| ID PK | ID PK | |
| questionID FK | askerID FK | users |
| answererID FK | question_text | ID PK |
| answer_text | category | points |
| time_answered | time_asked | |

- Functional Dependencies: all values in each table are functionally determined by the primary key.
- Multivalued Dependencies: none

New Schema built with Primitives

Please see attached hand-drawn schema

6: Ambient Intelligence

Last summer, I went to Italy and had to ride a lot of trains. I'm going to describe a higher-tech ticketing system than the one they currently have, in the same train stations that I visited.

Description of System

- This system would be located in a train station
- The first part of the system would be touch-screen ticketing kiosks. The user would purchase their ticket from at this kiosk and their thumbprint would be scanned, as a means of storing their identity. The second part of the system would be a thumbprint reader on the actual train, which users would have to use in order to validate their ticket and board the train.
- The system adapts to a diverse population by requiring a standardized set of information to purchase a ticket. Everyone who wants to use this system should have a means of purchasing a train ticket. If a user did not have thumbprints, a backup identification method such as loading their identity into the system from a credit card should be possible, but I'm not going to worry about such a case for this example.
- The hardware involved would be thumbprint readers, credit card scanners, touch-screen computers, and web servers.

Web Connectivity

- Aweb connection would be valuable to the users because the system could then email them receipts or itineraries.
- Users could also have the option of buying a ticket on a web-connected, fingerprint-enabled device such as an iPhone 5S, and not have to go early to the station to use the kiosks.
- Aweb connection would be valuable to the implementers because it could provide information about the customers, which could lead to information helpful for marketing or other analytical purposes.
- Aconnection to the system's servers that host information about routes, available tickets, and purchases is necessary for the system's basic functions.

Appropriate Database

The best database to back this system would be a relational database system such as MySQL. This system's primary purpose is transactions, so the locks and various levels of transaction protocols that MySQL supports would be necessary. One transaction would be when a user purchases a ticket. This would need to be stored in the database without any uncertainty that the row(s) were written. The other transaction would be when a user checks in to the actual train. This transaction would require less certainty and more speed than purchases, so a lower isolation level should be used. Since there is not much else going on other than transactions, other types of databases would not be as valuable for this system.

7: Comparisons

| Database Name | Structure | "Query" Method | Transactions | Focus |
|---------------|-----------|-----------------------|-------------------|---------------------------------|
| Neo4J | Graph | Cypher Query Language | full ACID support | ease for common graph-like data |

| | | | | |
|-----------|-----------------|-------------------------------|------------------------------------|---|
| Mongo | BSON/Document | MongoDB CRUD Operations | Nested Documents | large-scale document storage |
| MySQL | Tables | SQL | full rollback / ACID support | transactions, routine querying |
| Cassandra | Column Families | CQL | 3rd Party additions required | Big Data, masterless decentralization, scalability |
| Hbase | Columns | Map- Reduce/Hadoop | mostly ACID support (locks) | scalability, clusters, big data |