

Google Docs System Design

⇒ How are documents stored?

- ↳ Refer Part-I for starting with Google Docs
- ↳ subscribe @msDeepSingh on YouTube.

Storage Requirements

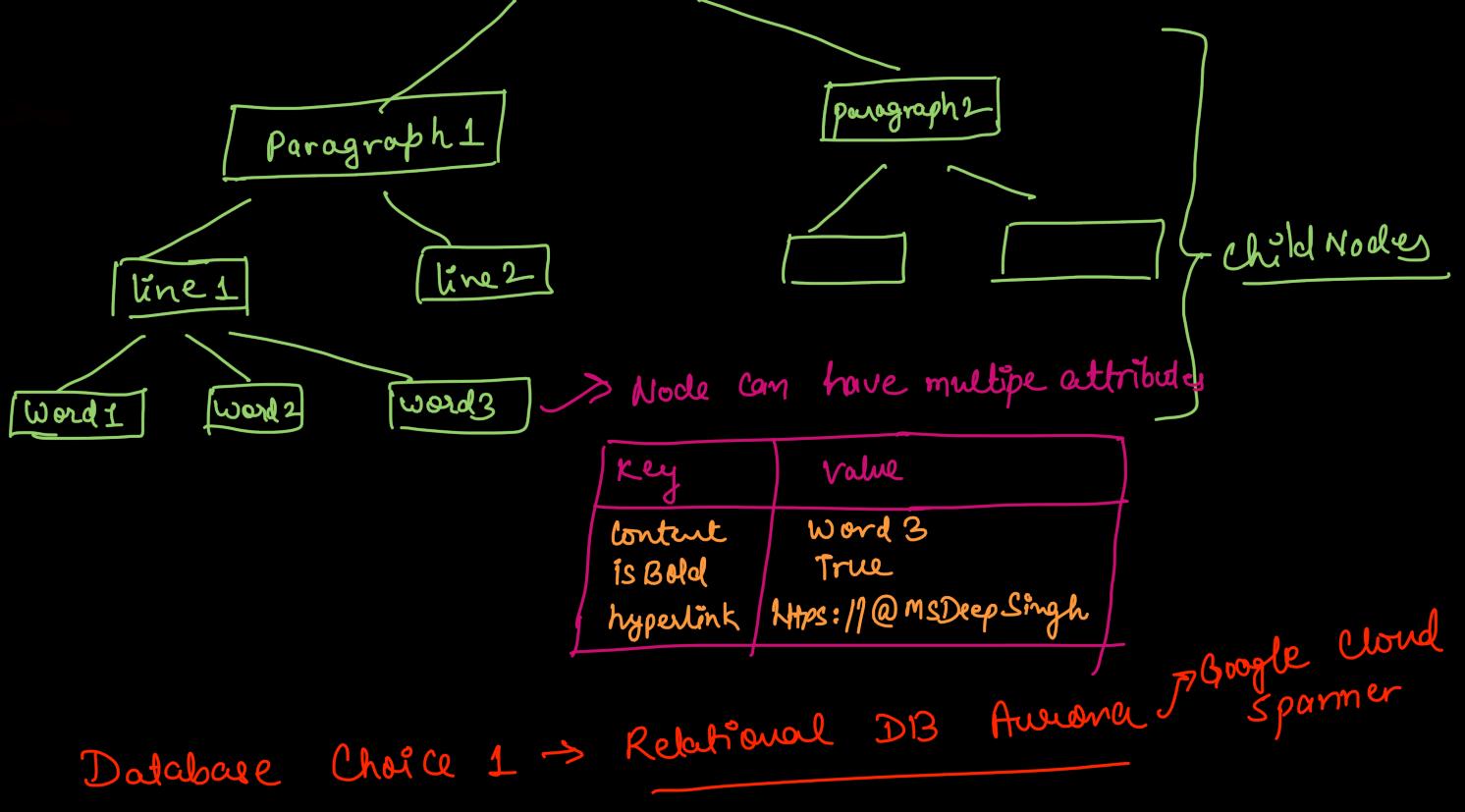
- ① Store document content
- ② User information like credentials, preferences
- ③ Document metadata like documentId, createdBy, etc.
- ④ User sessions database for websocket connection management
- ⑤ Document collaboration data - user actions or comments.

Note → I'm talking in terms of database choices from AWS context, Google Cloud has similar services - check link in video description.

① Document Content Storage

- Strong consistency
- High availability and scalability
- Content stored in chunks for easy change management.
 - ↓
 - One feasible way is to maintain a tree structure.

Complete doc - root node



Database Choice 1 → Relational DB Awareness ↗ Google Cloud Spanner

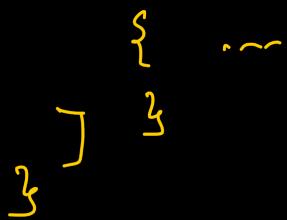
- ① Document Table → documentId, title, document rootNode. ↗ paragraph, line etc
- ② Nodes Table → nodeId, parentNodeId, NodeType, AttributeMap
- ③ Collab Table → documentId, collaboratorEmail, permissions

Database Choice 2 → DynamoDB ↗ Google BigTable
partition Key - documentId

Sort Key - creation Time

```
{
  "documentId": "12345",
  "title": "Subscribe to YT channel",
  "creationTime": "2020",
  "author": "@msDeep Singh",
  "content": [
    {
      "ContentSectionId": "1",
      "type": "text",
      "content": "Hello"
    }
  ]
}
```

As document size grows, it will become difficult in continuous update for real-time collaborations.



Other DB schema →

Node Id → Partition Key
 parentNode Id → Sort Key
 child Node Ids → Node Ids along with attributes

Database choice 3 - Document DB

Similar kind of structure

→ More complex query support.

Node Id	Parent	Child Nodes
doc1	root	[section1, section2]
section1	doc1	[para1, para2]
para1	section1	[line1, line2]
line1	para1	[word1, word2]
word1	line1	[]

Key consideration → RDBMS can provide just specific node entry update instead of querying for entire document.



the information around figuring out node Id when multiple users are contributing can be derived with combination of user's input & conflict handling by system.

② User Information → Aurora or RDS.

(a) User Email

(b) Credentials

(c) preferences

③ Document Metadata → Dynamo or Aurora both are feasible.

(a) document Id

(b) Created By

(c) creation Date

d) Access Permissions = {userId, email, permissions. each

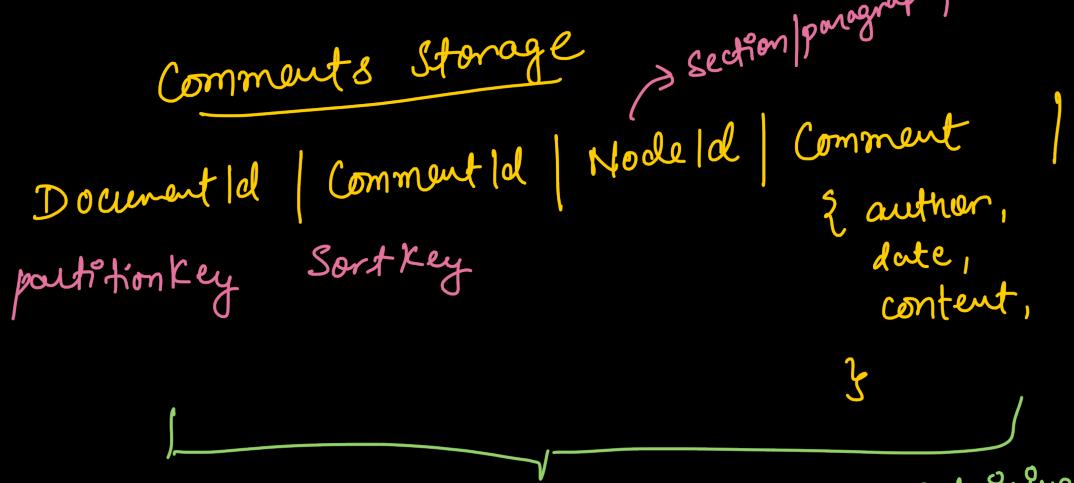
④ User sessions → Redis with persistent storage.

⑤ Collab Data → This will work in sync with DocumentStorage DB

① Comments

② Suggestions

③ Document Update Info → who updated specific document Node



Will this scale for maintaining comment replies?

↳ Each comment should have parent as well

Share your thoughts on how you'll think of
databases and their schema in comments.

Subscribe for more - @msDeep Singh
@msdeep14

Happy Learning ☺