# Extract, transform, load (ETL) MySQL to MongoDB Atlas using Data Pipeline

> 📋 **Tldr**
>
> Basically hierarchical to nosql data migration which involves:
>
> - Queries to extract MySQL data into Azure Data Factory
> - Pipelines to transform SQL into NoSQL
> - Json data to be loaded into MongoDB
>
> Good:
>
> - Excellent Performance (4ms; sec.2.4)
>
> Bad:
>
> - Unwanted data duplication that also lead to significant update effort (sec. 3.1)
> - Lack data validation (sec. 3.3)
> - Lack aggregation for continuous update (sec. 3.3)
> - No reverse migration (sec. 3.3)

# Extract, transform, load (ETL)

## Overview



Source:

ETL MySQL to MongoDB Atlas.pptx

Estimate cost for 72 hours: US$16.82

https://azure.com/e/0d6d1f802abc478c9cedf6594e9a6731

## Source – MySQL

**Azure Database for MySQL** Burstable B1s
MySQL version 8.0
Resource location: Azure East US

# Configuration Procedures

1. Create a Azure Database for MySQL flexible server
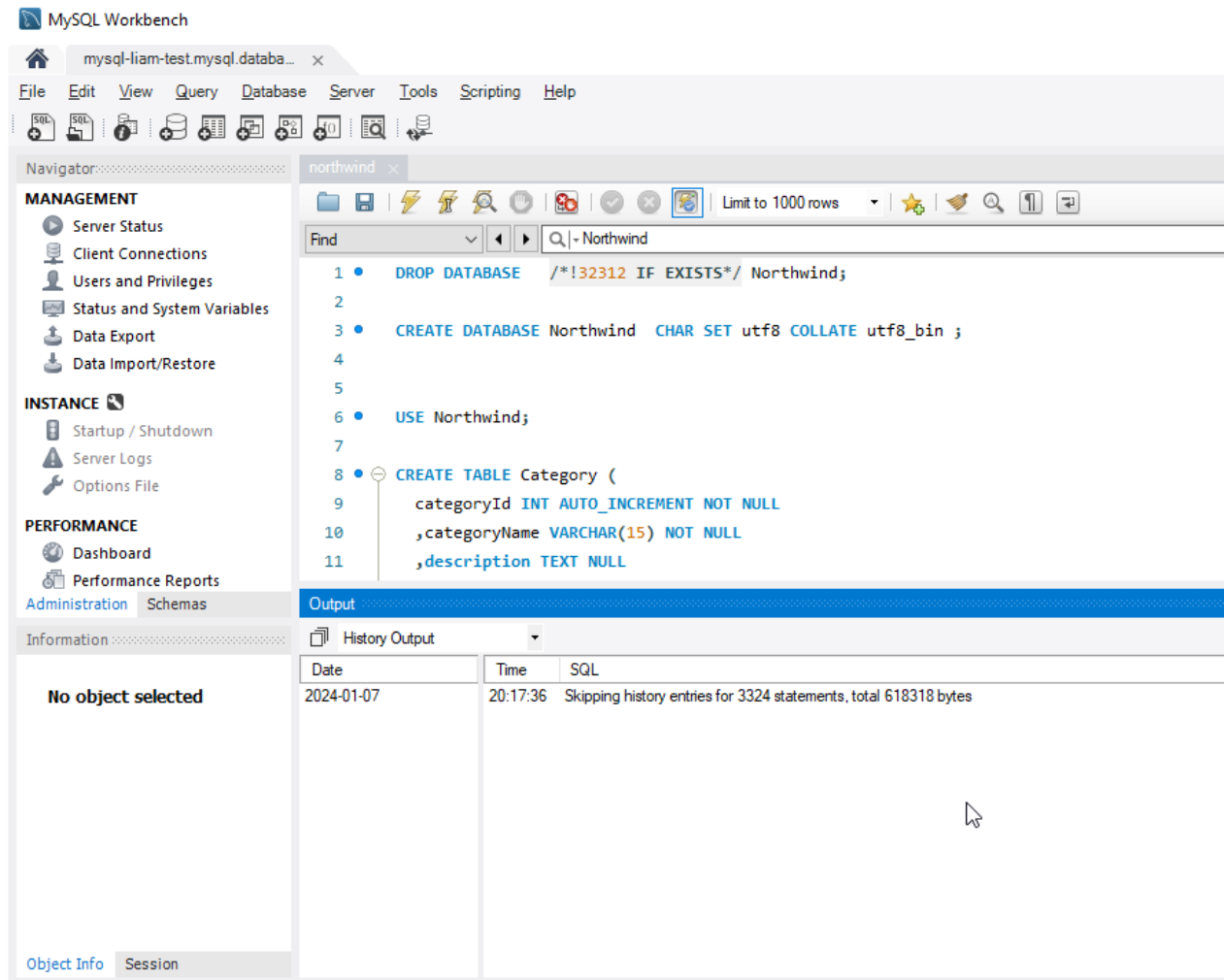2. Connect to Azure MySQL on MySQL Workbench
3. Import Northwind.sql



Data Reference from https://github.com/harryho/db-samples

## Azure MySQL Credential

azureadmin
Z57v36VQ4sm3E2Lc5299U4mQK

## Source Data Structure - ER Diagram

Procedures: Create ER Diagram using Reverse Engineer



Only Primary Keys and Foreign Keys are shown.

# Destination - MongoDB

**MongoDB Atlas** Free tier M0 Sandbox cluster
MongoDB version 6.0
Resource location: AWS us-east-1

# MongoDB Credential

liamng
Gak2QAKFi0J0q2wA

sa
NmnXthqG341zT3IR

# Data API Secrets

`liam-mongo-api`
API Key: `xY6rwTPKLkmbSBsDYBE4X1JB8aC4hrJjHwStgVyznXmlJmXoQ5EOrOMaYxJ43dHN`

# MongoDB API Test on Postman

## Postman project variables



> ⚡ **Error**
>
> Enable `CONTENT_TYPE` may causes error `mime: unexpected content after media subtype`

1. Open Find Document
2. Update `filter` in `Body` to one of the records
3. Verify json content



> 🖉 **Future development for my PayPal API website (Node.js 18)**

# Document database schema designs

Consolidate into single Document using Tree pattern

**orderdetail**
- orderDetailId INT
- orderId INT
- productId INT
- unitPrice DECIMAL(10,2)
- quantity SMALLINT
- discount DECIMAL(10,2)

**salesorder**
- orderId INT
- custId INT
- employeeId INT
- orderDate DATETIME
- requiredDate DATETIME
- shippedDate DATETIME
- shipperid INT
- freight DECIMAL(10,2)
- shipName VARCHAR(40)
- 5 more...

**customer**
- custId INT
- companyName VARCHAR(...
- contactName VARCHAR(30)
- contactTitle VARCHAR(30)
- address VARCHAR(60)
- city VARCHAR(15)
- 7 more...

**custcustdemographics**
- custId INT
- customerTypeId INT

**customerdemographics**
- customerTypeId INT
- 1 more...

Separate with Document References to avoid data repetition

**shipper**
- shipperId INT
- companyName VARCHAR(40)
- phone VARCHAR(44)

**product**
- productId INT
- productName VARCHAR(40)
- supplierId INT
- categoryId INT
- quantityPerUnit VARCHAR(20)
- unitPrice DECIMAL(10,2)
- unitsInStock SMALLINT
- unitsOnOrder SMALLINT
- reorderLevel SMALLINT
- discontinued CHAR(1)

**supplier**
- supplierId INT
- companyName VARCHAR(...
- contactName VARCHAR(30)
- contactTitle VARCHAR(30)
- address VARCHAR(60)
- city VARCHAR(15)
- 7 more...

**category**
- categoryId INT
- categoryName VARCHAR(15)
- description TEXT
- picture BLOB

Consolidate into single Document using Tree pattern

**employee**
- employeeId INT
- lastname VARCHAR(20)
- firstname VARCHAR(10)
- title VARCHAR(30)
- titleOfCourtesy VARCHAR(...
- birthDate DATETIME
- hireDate DATETIME
- address VARCHAR(60)
- 12 more...

**employeeterritory**
- employeeId INT
- territoryId VARCHAR(20)

**territory**
- territoryId VARCHAR(20)
- territorydescription VARCHAR(50)
- regionId INT

**region**
- regionId INT
- regiondescription VARCHAR(50)

**orderdetail**
- orderDetailId INT
- salesorder object
- product object
- unitPrice DECIMAL(10,2)
- quantity SMALLINT
- discount DECIMAL(10,2)

**salesorder**
- orderId INT
- customer object
- employeeId INT
- orderDate DATETIME
- requiredDate DATETIME
- shippedDate DATETIME
- shipper object
- freight DECIMAL(10,2)
- shipName VARCHAR(40)
- shipAddress VARCHAR(60)
- shipCity VARCHAR(15)
- shipRegion VARCHAR(15)
- shipPostalCode VARCHAR(...
- shipCountry VARCHAR(15)

**customer**
- custId INT
- companyName VARCHAR(...
- contactName VARCHAR(30)
- contactTitle VARCHAR(30)
- address VARCHAR(60)
- city VARCHAR(15)
- region VARCHAR(15)
- 6 more...

**shipper**
- shipperId INT
- companyName VARCHAR(40)
- phone VARCHAR(44)

**product**
- productId INT
- productName VARCHAR(40)
- supplier object
- category object
- quantityPerUnit VARCHAR(20)
- unitPrice DECIMAL(10,2)
- unitsInStock SMALLINT
- unitsOnOrder SMALLINT
- reorderLevel SMALLINT
- discontinued CHAR(1)

**category**
- categoryId INT
- categoryName VARCHAR(15)
- description TEXT
- picture BLOB

**supplier**
- supplierId INT
- companyName VARCHAR(...
- contactName VARCHAR(30)
- contactTitle VARCHAR(30)
- address VARCHAR(60)
- city VARCHAR(15)
- region VARCHAR(15)
- postalCode VARCHAR(10)
- 5 more...

To transform to nosql on Azure Data Factory:

1. "Derive" a new column based on existing column's value

2. "Select" the newly derived column to clean excess columns

3. "Join" the derived column to parent table

## Considerations

- Maximize number of Many-to-One (n:1) relationships to embed data in single document.
- Use references and keep the `customer demographics` information in a separate collection from the `customer` collection to avoid repetition of the demographics data,.
- Other patterns are not implemented in this project due to complexity and require more in-depth research on data modeling.

## SQL-MongoDB Queries Conversion Reference

| SQL query | Equivalent MongoDB query |
| --- | --- |
| `SELECT * FROM users` | `db.users.find({})` |

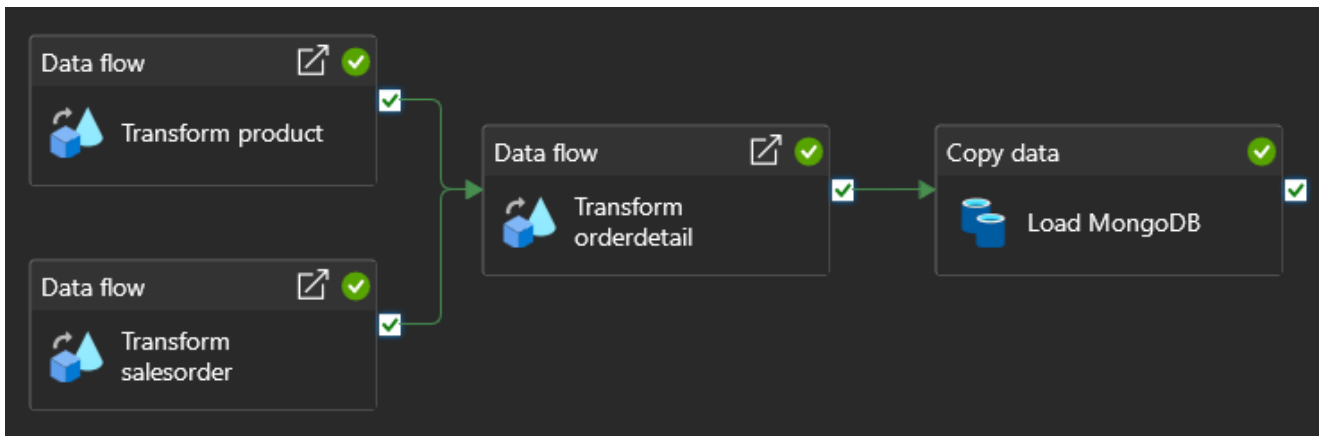More on

- Microsoft KB: Copy data from or to MongoDB using Azure Data Factory
- MongoDB KB: SQL to MongoDB Mapping Chart

# Azure Data Factory (ETL Tool)

## High-level ETL Steps

1. Query MySQL data into Azure Data Factory
2. Transform each table (**hierarchical**) into collection (Nosql, **Tree** pattern) in 3 steps:
   - Derive a new column based on existing column's value
   - Select the newly derived column to clean excess columns
   - Join the derived column to parent table
3. Store the transformed document into a sink (`.json`)
4. Copy the document to MongoDB

## Configuration Procedures

> ⚠️ **Pre-requisites**
>
> - Add client address to Azure Key Vault exception
> - Allow trusted Microsoft services to bypass this firewall on Key Vault
> - Add Azure RBAC (**Key Vault Administrator**) for yourself/ Access Policies in Azure Key Vault
> - Add Azure RBAC (**Key Vault Secrets User**) for Azure Data Factory Managed identity in Azure Key Vault
> - Add Azure RBAC (**Reader**) to Azure Data Factory Managed identity in Azure Database for MySQL
> - Add Azure RBAC (**Storage Blob Data Contributor**) to Azure Data Factory Managed identity in Storage Account
> - SSL in Azure Data Factory is not supported
> - You may need to add 0.0.0.0/0 into MongoDB `Network Access` > `IP Access List`, given Azure Data Factory (ADF) in this example is a public shared runtime, although it is possible to retrieve the public endpoints of ADF, it is unrealistic to add the extensive list of IP ranges which has over 50 entries and is also dynamic.

> 🖊️ **Note**
>
> | Secret Name | Value/Connection String |
> |---|---|
> | mysql-conn | `Server="mysql-liam-test.mysql.database.azure.com";UserID="azureadmin";Password="Z57v36VQ4sm3E2Lc5299U4mQK";Database="northwind";` |

| Secret Name | Value/Connection String |
|---|---|
| mongodb-conn | `mongodb+srv://liamng:Gak2QAKFi0Joq2wA@cluster0.f67f2bk.mongodb.net/` |

1. Create Azure Data Factory
2. Configure Azure Data Factory Linked Services
3. Extract and Transform tables into document-like format and repeat for all tables
4. Load data into MongoDB
5. Connect to MongoDB with VS code to verify result

> 🔥 **Important**
>
> For Step 3, to prevent potential data loss, `left join` (outer) is preferred in queries.

> 🖊 **Note**
>
> For step 3, as an alternative, it is possible to prepare a **sample schema** (with Tree pattern) for data sink, the columns will be mapped automatically. However, upon testing, if there is missing column, errors are expected and no fixes were found for the mapping afterwards. Further study is required.

> ⚡ **Error**
>
> - Clean `.json` file before running pipeline
> - **Output to single file** must be selected for flat namespace. As a result, only Single partition can be set and this will make data flow execution longer.
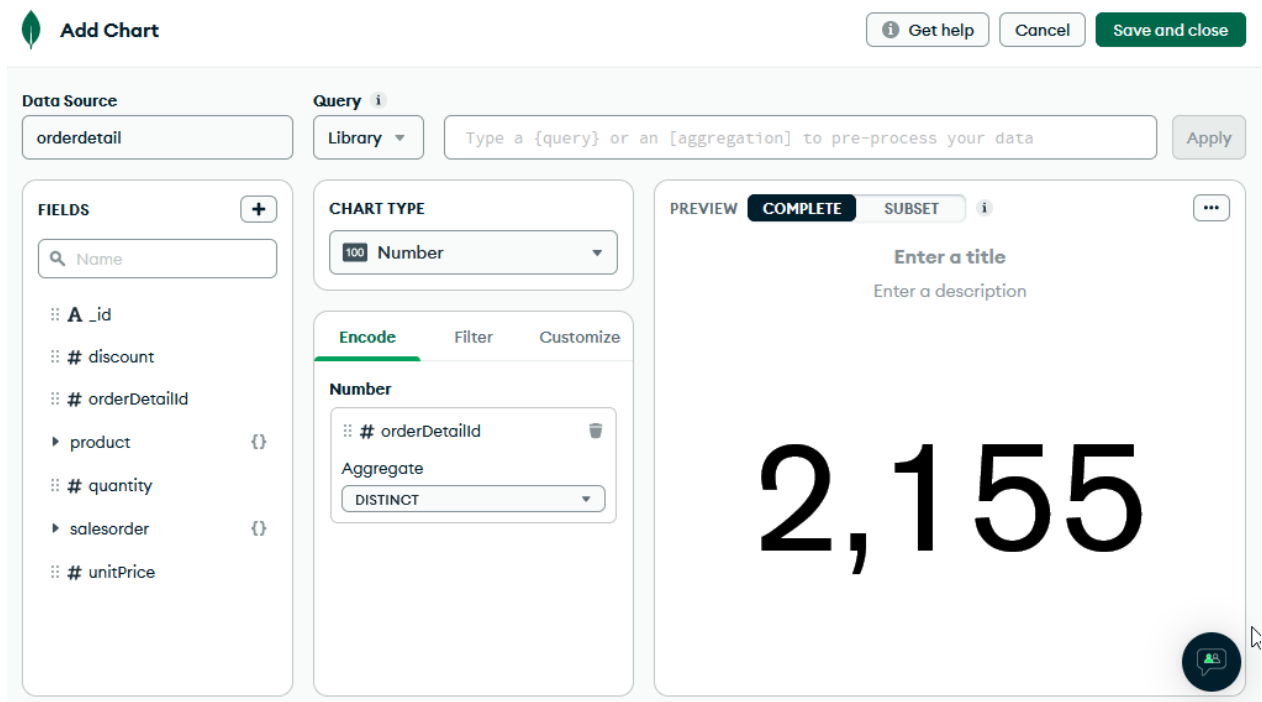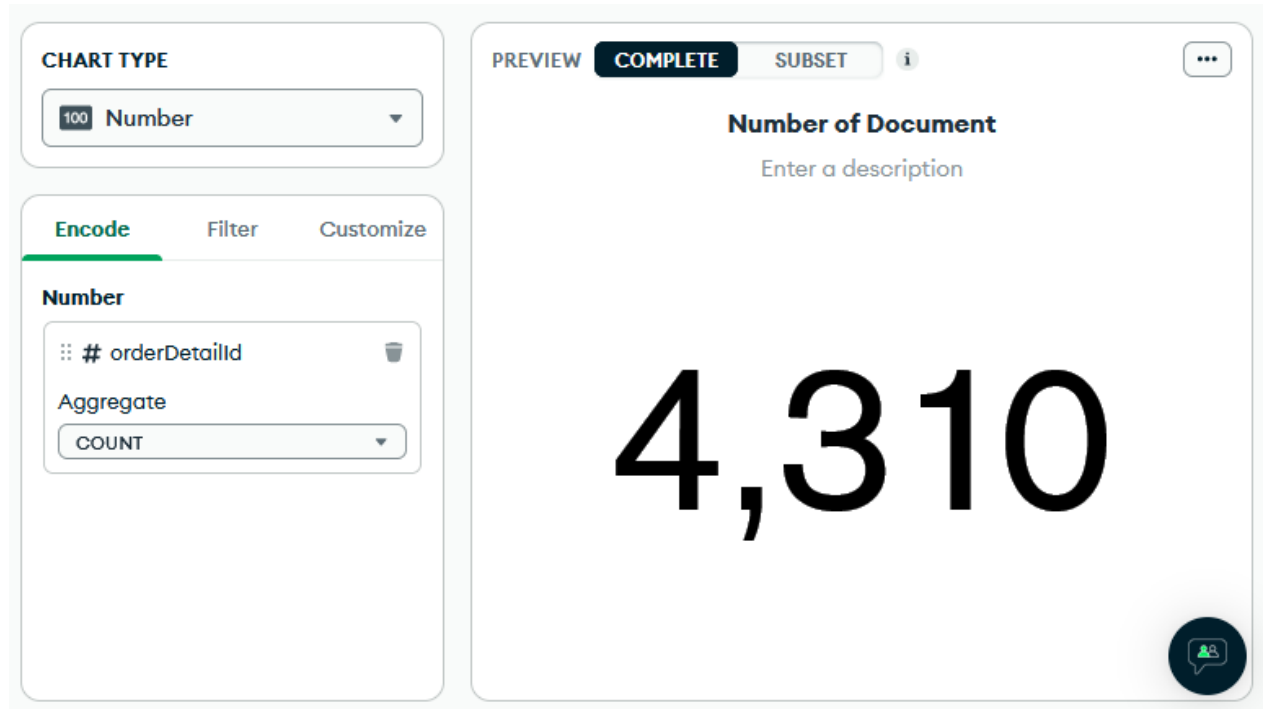
# Test Case

## Verify total number of records

- MySQL



- MongoDB Chart

Number of *Distinct* Document
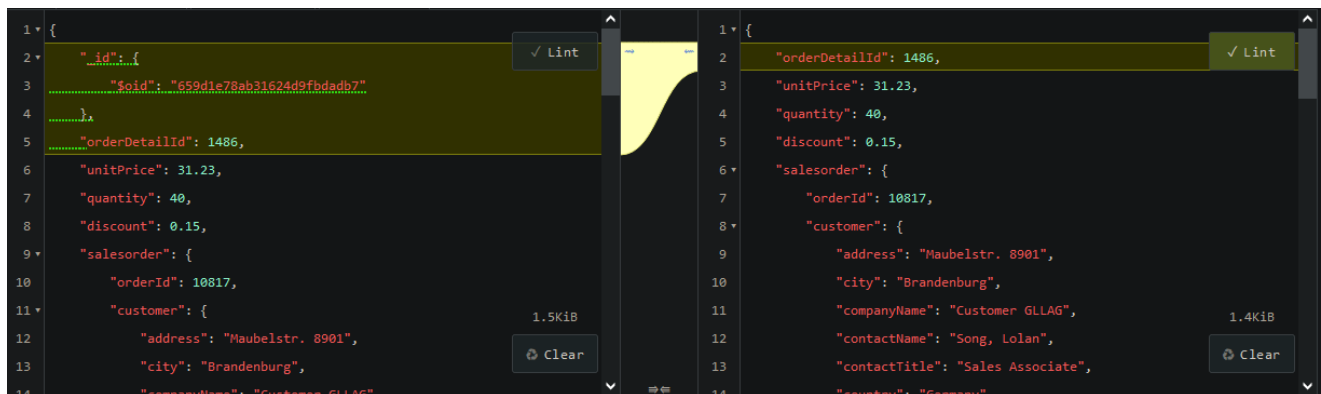
Number of Document



Since the pipeline has been run multiple times (although `Upsert` is selected), there are duplicated records in MongoDB.

To **aggregate** these records, see Aggregation Commands details at MongoDB KB

## Schema Validation

Compare Azure Data Factory output and MongoDB record



The number of column in MySQL also matches number of fields in MongoDB (except `_id`). To prevent loss of null columns in Data Factory, `'null'` text is used (also due to product limitation).

But still, due to the transformation, comparing the schema of both DB is difficult.

## Azure Data Factory ETL Process Performance

Enabled **incremental column** and **change data capture**:

First-time:



Subsequent:



Improvement in ETL process time: 40+%

However, the scale and number of pipeline execution are trivial at this moment.

# Query Performance Difference between MySQL and MongoDB

- MySQL
  - Query takes 47ms

```sql
Select *

From orderdetail natural left outer join salesorder natural left outer join customer
natural left outer join product natural left outer join supplier

Where companyName = "Customer GLLAG" AND freight > 50 AND supplierId in (3,5,10)
```

- MongoDB
  - 4ms to search through 4.3k documents/ records
  - 10 times faster than MySQL

Query string `{"salesorder.customer.companyName": "Customer GLLAG", "salesorder.freight": {"$gt": 50}, "product.supplier.country": {"$in": ["Japan", "Germany", "USA"]}}`

# Discussion

## Data Model

For One-to-Many relationship, embedding is preferred BUT putting all into single document is very likely to be considered anti-pattern, in terms of data duplication, administrative burden from updating documents, and data privacy/ security due to lack of segregation;

e.g. Every time we need to update information about customer, we'll need to update the document for every orderDetail

This is unlikely for these information to be frequently displayed or updated together every time.

**Good**:

- Excellent Performance (4ms)

**Bad**:

- Data duplication
- Large update effort

There are 3 approaches to this problem:

1. Embed (single document)
2. Reference (separate collections)
3. **Extended Reference** (separate collections but with duplicated data)

**Fix**:
**Balance of storage cost and performance**

- Separate `product`, `customer`, `shipper` from `orderdetail`
- Keep frequently access information in `orderdetail` document, such as, `productName`, `companyName`, `contactName`
- Reference new collections with `$lookup` operations

## Data Transformation

- Tradeoff between simplicity and performance.
    - In the case of the tree pattern, you get better performance by avoiding multiple joins, however, you will need to manage the updates to your graph. Plus, excess/ duplicated data across many documents.

Combined with conclusion in previous section, these are the rule of thumbs:

- 1: Favor embedding unless there is a compelling reason not to.
- 2: Needing to access an object on its own is a compelling reason not to embed it.
- 3: Avoid joins and lookups if possible, but don't be afraid if they can provide a better schema design.

- **4**: Arrays should not grow without bound. If there are more than a couple of hundred documents on the many side, don't embed them; if there are more than a few thousand documents on the many side, don't use an array of ObjectID references. High-cardinality arrays are a compelling reason not to embed.

As always, with MongoDB, how you model your data depends entirely on your particular application's data access patterns. You want to structure your data to match the ways that your application queries and updates it.

## Data Consistency and Replication

- Lack Data Validation, very difficult to compare apple-to-apple... need a solution
- Lack Aggregation on MongoDB side? every time the pipeline run (copy operation) will duplicate data set. No continuous/ incremental update to MongoDB.
- Potential improvement with `Trigger` for pipeline execution when data / schema update is performed on source MySQL DB. (ADF CDC is not yet supported for MongoDB)
    - Timer trigger for pipeline execution in regular basis
    - Event-driven trigger, and if Data Lake is used, could benefit streaming data
- No Reverse Migration native function on Azure Data Factory (from MongoDB back to MySQL)

## Others

1. Scalability and Performance
    1. Hierarchical namespace would allows optimized partition + parallel process and make data flow execution faster
2. Security Compliance and Privacy

---

# Appendix

## Source Files

In chronological build order

### Key Vault

📄 ExportedTemplate-kv-liam-test.zip

### Storage Account (Blob)

📄 ExportedTemplate-stliamtestdf.zip

### Azure Database for MySQL

ExportedTemplate-mysql-liam-test.zip

## Azure Data Factory

arm_template_adf-liam-test.zip

**Factory Resources**  ⟱  «

🔽 Filter resources by name   +

◢ **Pipelines**   1

   🎞 ETL from MySQL to MongoDB

▷ **Change Data Capture (preview)**   0

◢ **Datasets**   3

   ▦ collection_MongoDbAtlas

   ▦ Json_blob

   ▦ table_northwind

◢ **Data flows**   3

   🗘 Prepare orderdetail

   🗘 Prepare product

   🗘 Prepare salesorder

▷ **Power Query**   0

dataset_table_northwind.json

dataset_Json_blob.json

dataset_collection_MongoDbAtlas.json

dataflow_Prepare product.json

dataflow_Prepare salesorder.json

dataflow_Prepare orderdetail.json

ETL from MySQL to MongoDB.json

# MongoDB

Sample Schema

```json
{
  "_id": {
    "$oid": "659d1e78ab31624d9fbdadb7"
  },
  "orderDetailId": {
    "$numberInt": "1486"
  },
  "unitPrice": {
    "$numberDouble": "31.23"
  },
  "quantity": {
    "$numberInt": "40"
  },
  "discount": {
    "$numberDouble": "0.15"
  },
  "salesorder": {
    "orderId": {
      "$numberInt": "10817"
    },
    "customer": {
      "address": "Maubelstr. 8901",
      "city": "Brandenburg",
      "companyName": "Customer GLLAG",
      "contactName": "Song, Lolan",
      "contactTitle": "Sales Associate",
      "country": "Germany",
      "custId": {
        "$numberInt": "39"
      },
      "phone": "0555-34567",
      "postalCode": {
        "$numberInt": "10060"
      }
    },
    "employeeId": {
      "$numberInt": "3"
    },
    "orderDate": "2008-01-06 00:00:00",
    "requiredDate": "2008-01-20 00:00:00",
    "shippedDate": "2008-01-13 00:00:00",
    "shipper": {
```

```json
      "companyName": "Customer GLLAG",
      "phone": "0555-34567",
      "shipperid": {
        "$numberInt": "2"
      }
    },
    "freight": {
      "$numberDouble": "306.07"
    },
    "shipName": "Destination RMBHM",
    "shipAddress": "Maubelstr. 1234",
    "shipCity": "Brandenburg",
    "shipRegion": "null",
    "shipPostalCode": {
      "$numberInt": "10209"
    }
  },
  "product": {
    "productId": {
      "$numberInt": "26"
    },
    "productName": "Product HLGZA",
    "supplier": {
      "address": "Tiergartenstraße 3456",
      "city": "Berlin",
      "companyName": "Supplier ZPYVS",
      "contactName": "Jain, Mukesh",
      "contactTitle": "Sales Manager",
      "country": "Germany",
      "phone": "(010) 3456789",
      "postalCode": {
        "$numberInt": "10016"
      },
      "supplierId": {
        "$numberInt": "11"
      }
    },
    "category": {
      "categoryId": {
        "$numberInt": "3"
      },
      "categoryName": "Confections",
      "desc": "Desserts, candies, and sweet breads"
    },
    "quantityPerUnit": "null",
    "unitPrice": {
```

```
      "$numberDouble": "31.23"
    },
    "unitsInStock": "null",
    "unitsOnOrder": "null",
    "reorderLevel": "null",
    "discontinued": "null"
  }
}
```