

Data-Intensive Distributed Computing

CS 431/631 451/651 (Fall 2021)

Part 7: Analyzing Relational Data (1/3)

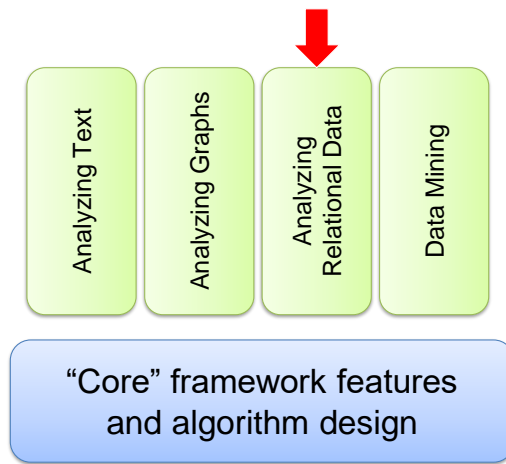
Ali Abedi

These slides are available at <https://www.student.cs.uwaterloo.ca/~cs451>



This work is licensed under a Creative Commons Attribution-NonCommercial-Share Alike 3.0 United States
See <http://creativecommons.org/licenses/by-nc-sa/3.0/us/> for details

Structure of the Course



Evolution of Enterprise Architectures

Next two sessions: techniques, algorithms, and
optimizations for relational processing

3

users

Monolithic
Application

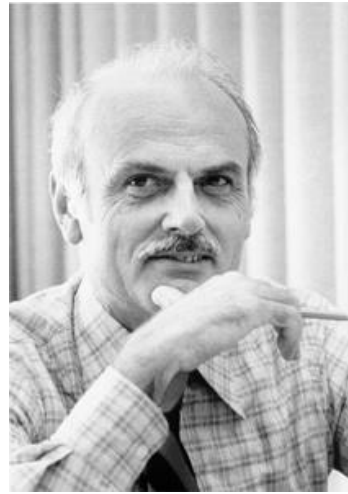
users

Frontend

Backend

Edgar F. Codd

- Inventor of the relational model for DBs
- SQL was created based on his work
- Turing award winner in 1981



users

Frontend

Backend

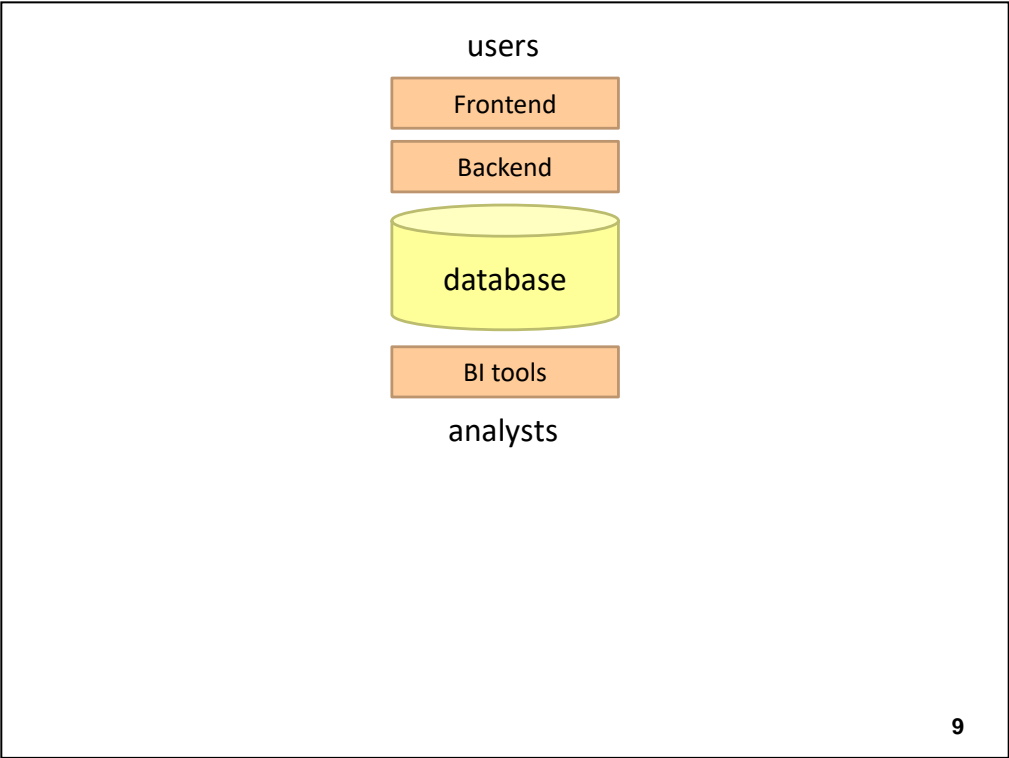
database

Why is this a good idea?

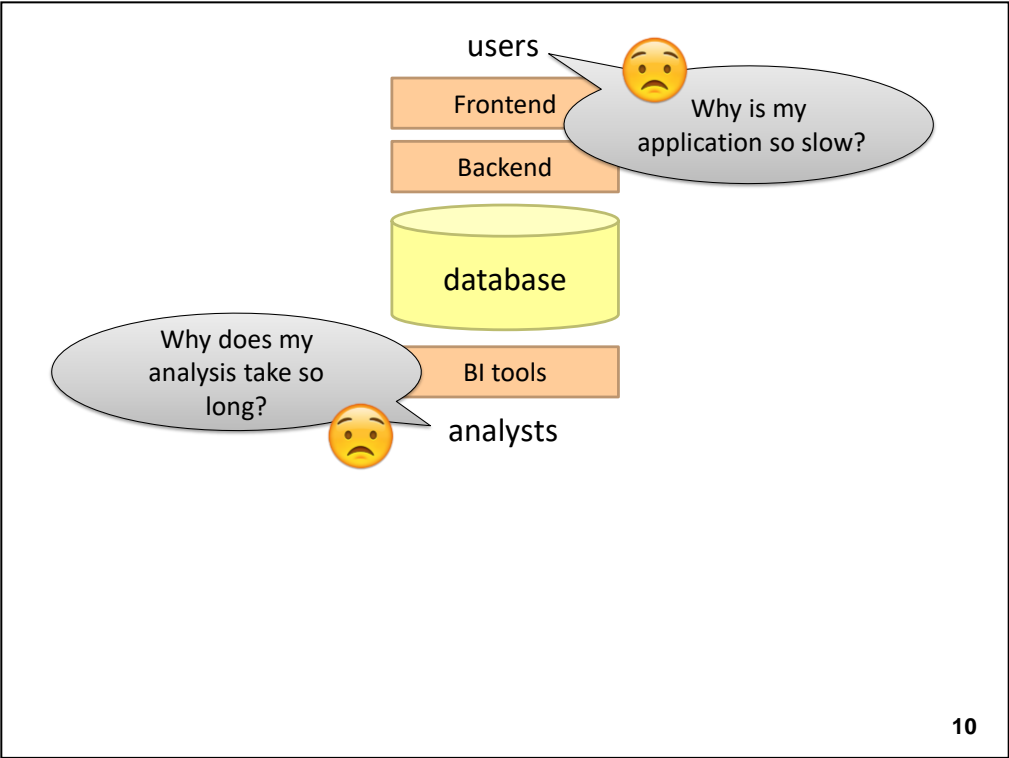
Business Intelligence

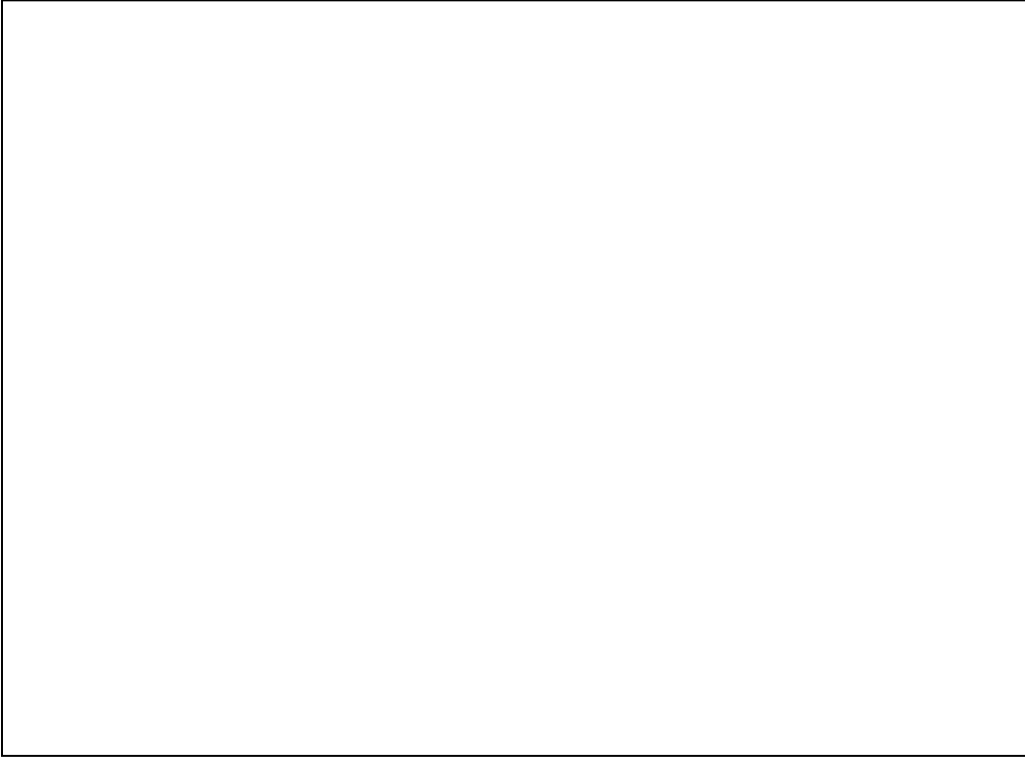
An organization should retain data that result from carrying out its mission and exploit those data to generate insights that benefit the organization, for example, market analysis, strategic planning, decision making, etc.

Duh!?



BI: Business intelligence





OLTP and OLAP Together?

Downsides of co-existing OLTP and OLAP workloads

Poor memory management
Conflicting data access patterns
Variable latency

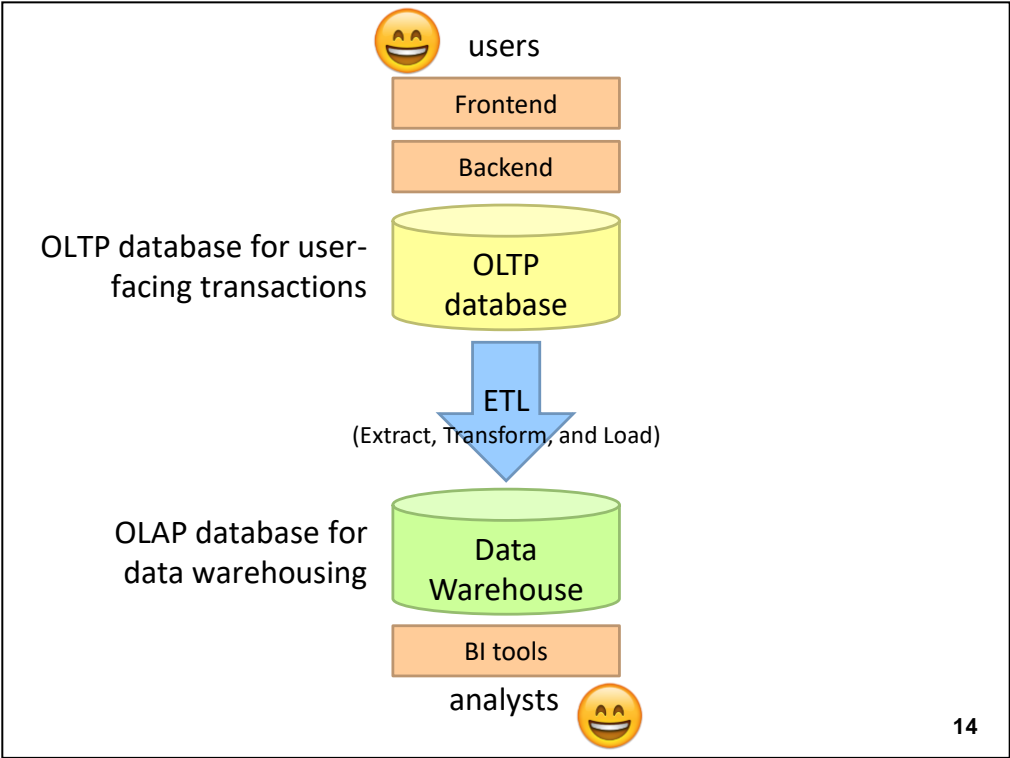


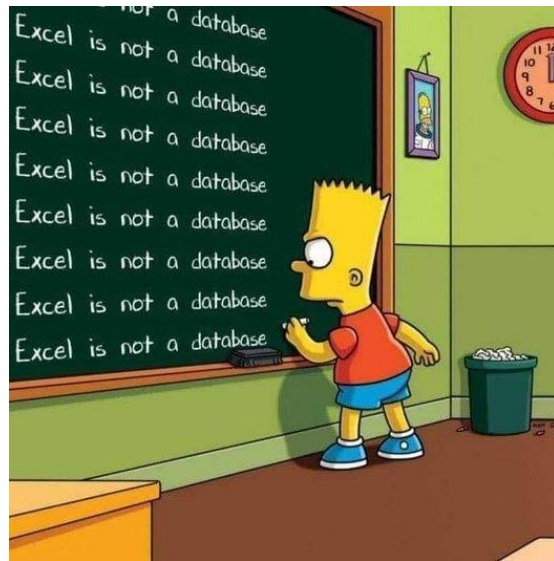
users and analysts

Solution?

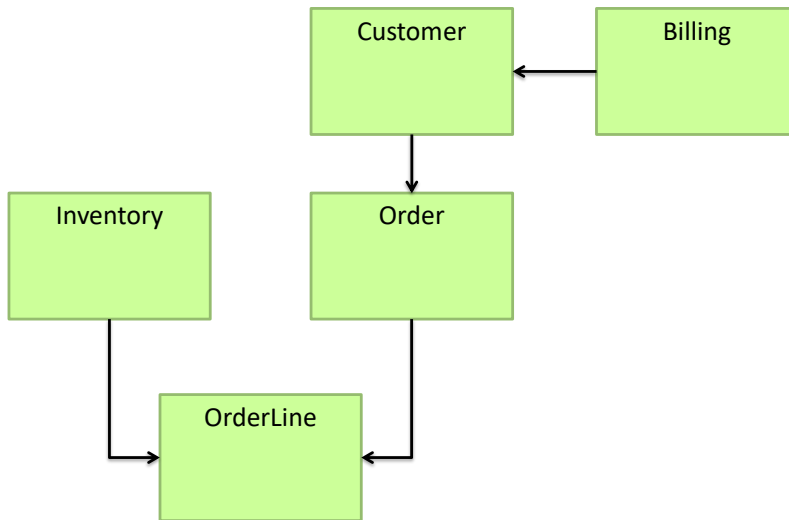


Source: Wikipedia (Warehouse)



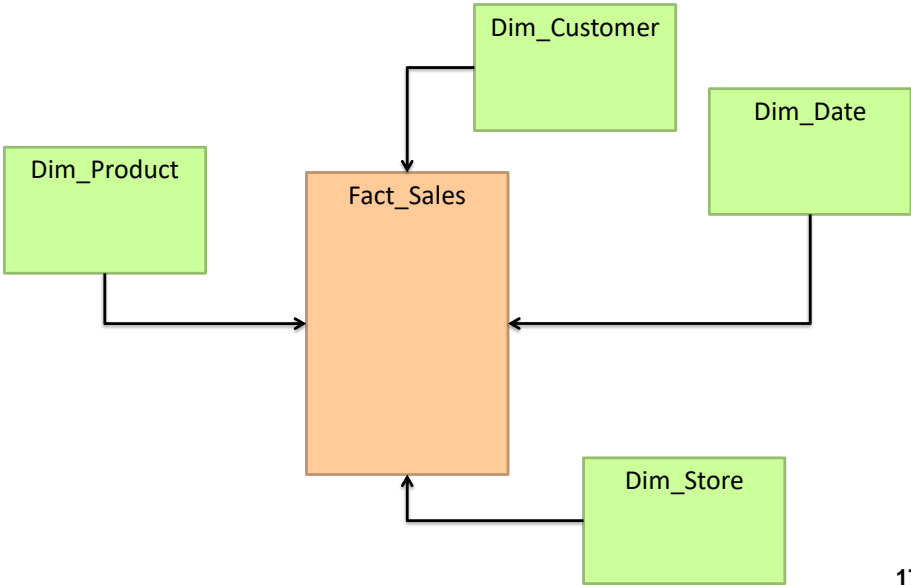


A Simple OLTP Schema



16

A Simple OLAP Schema



ETL

Extract

Transform

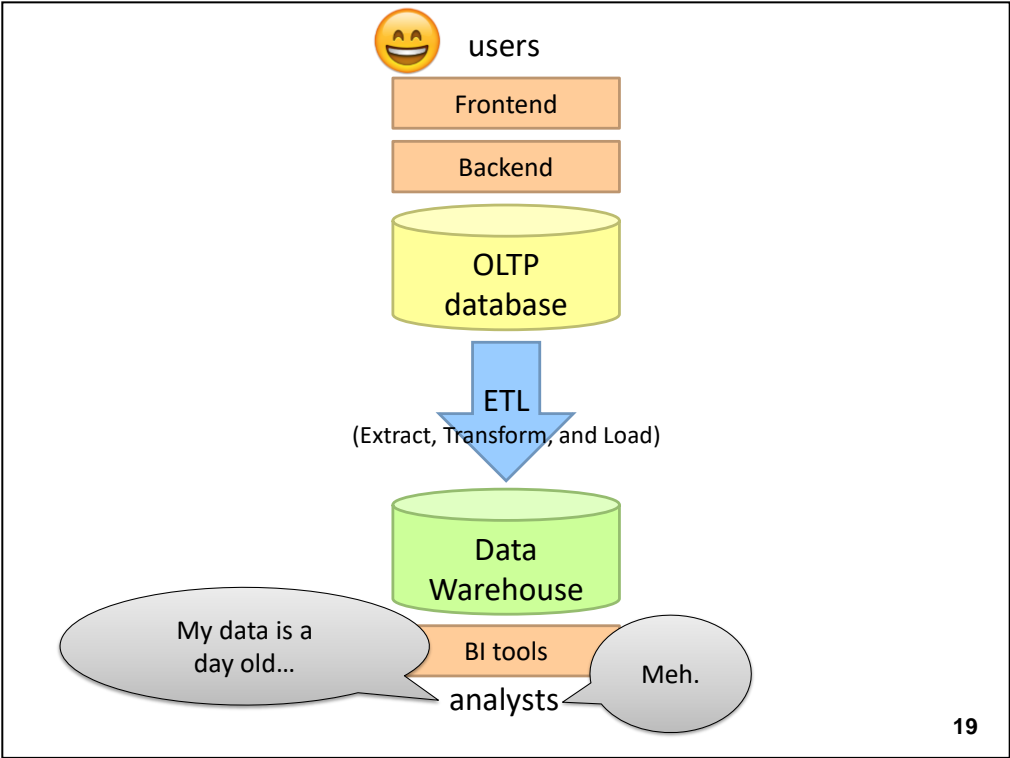
Data cleaning and integrity checking

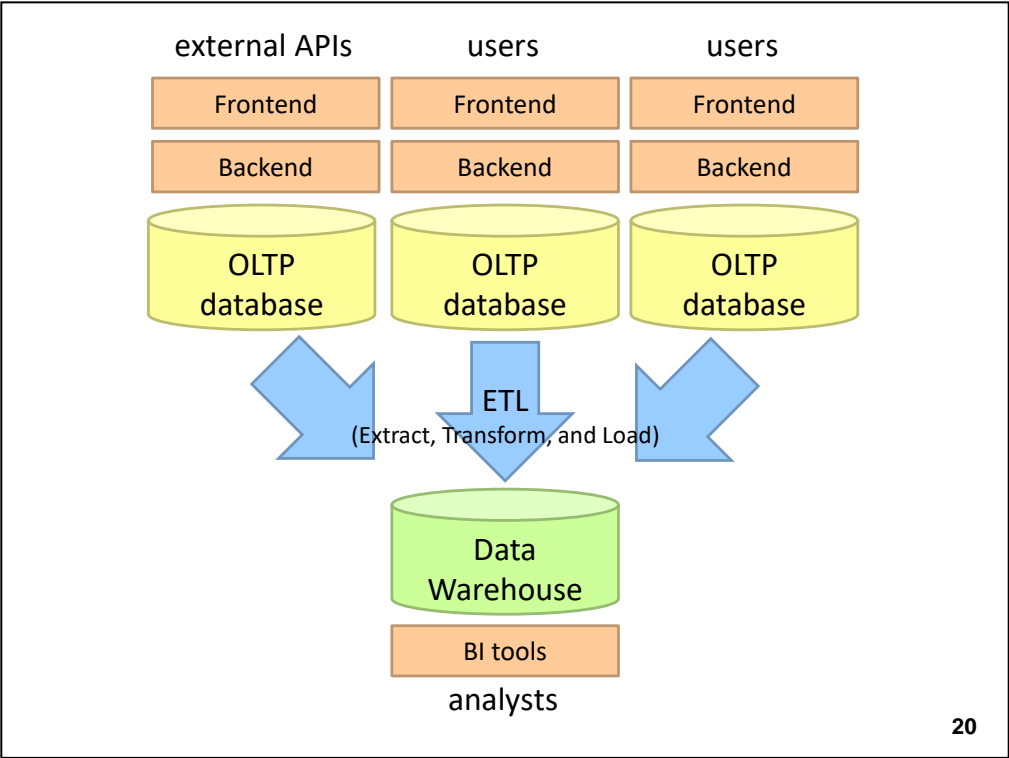
Schema conversion

Field transformations

Load

When does ETL happen?





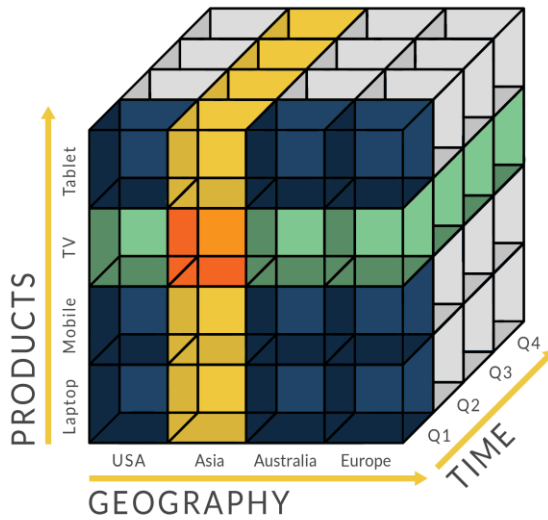
What do you actually do?

Report generation

Dashboards

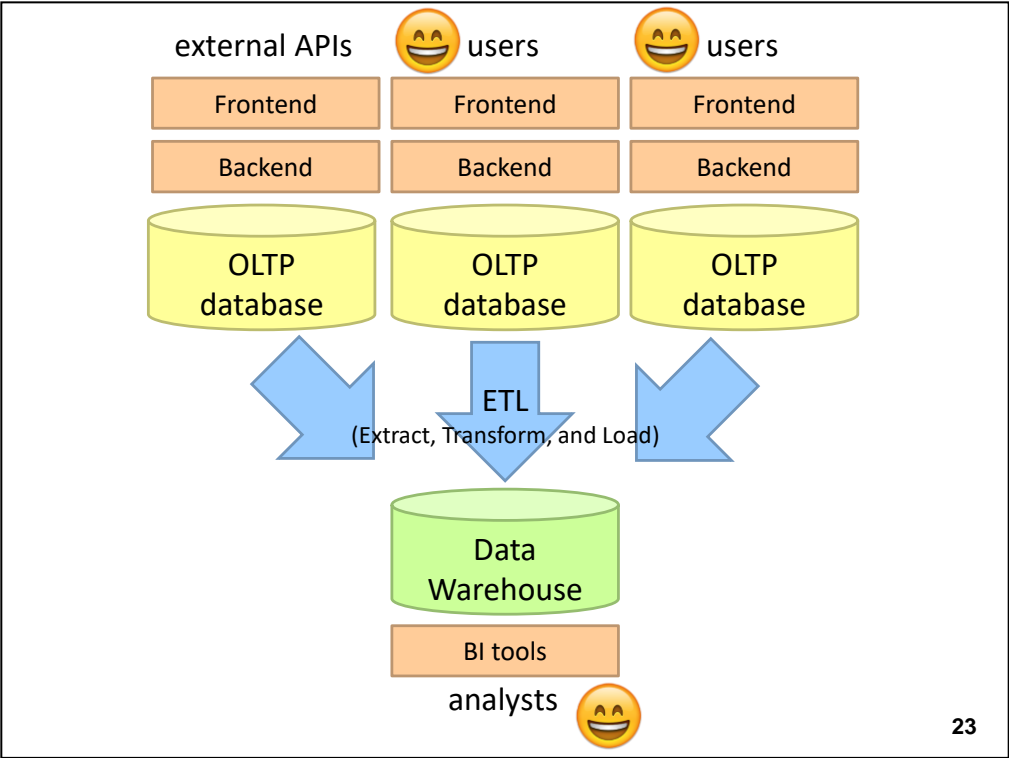
Ad hoc analyses

OLAP Cubes



Common operations
slice and dice
roll up/drill down
pivot

<https://youtu.be/LRdsZqrwOrc>



Fast forward...

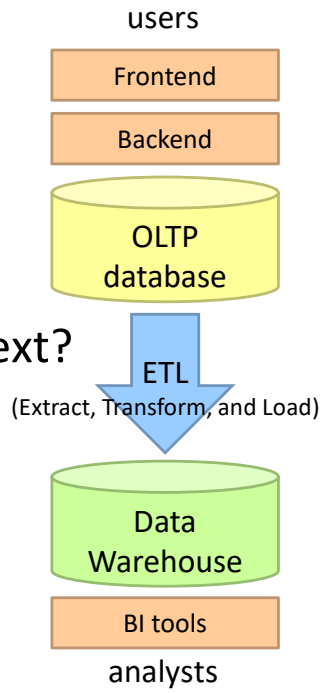
facebook®

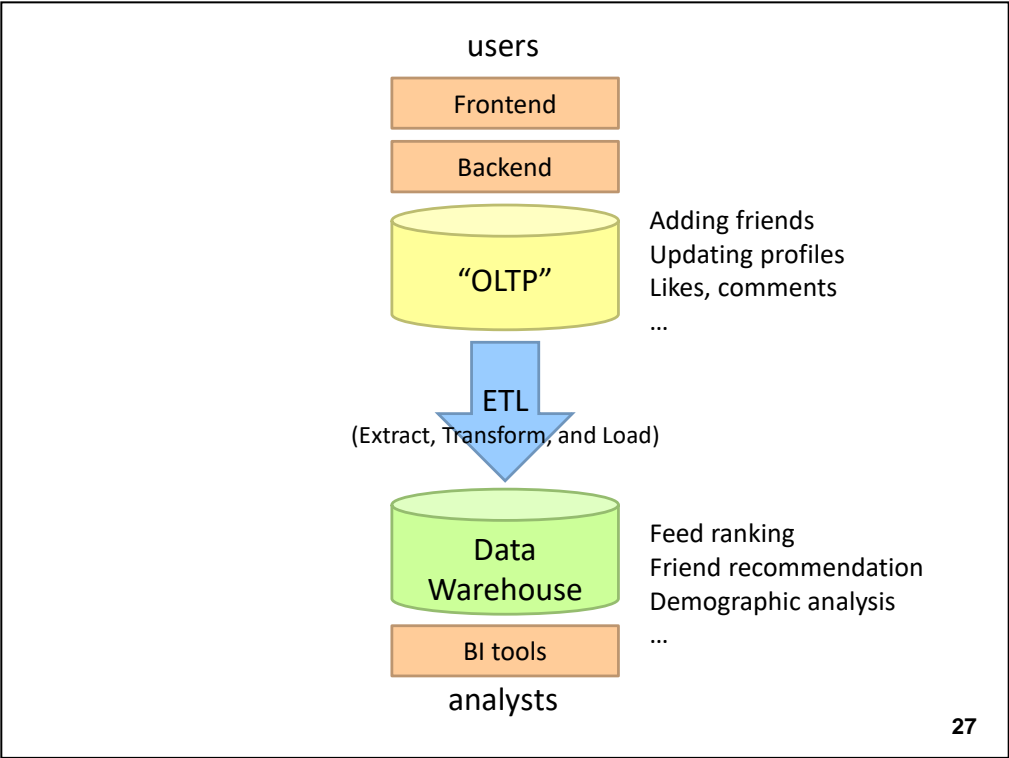
Jeff Hammerbacher, Information Platforms and the Rise of the Data Scientist.
In, *Beautiful Data*, O'Reilly, 2009.

“On the first day of logging the Facebook clickstream, more than 400 gigabytes of data was collected. The load, index, and aggregation processes for this data set really taxed the Oracle data warehouse. Even after significant tuning, we were unable to aggregate a day of clickstream data in less than 24 hours.”

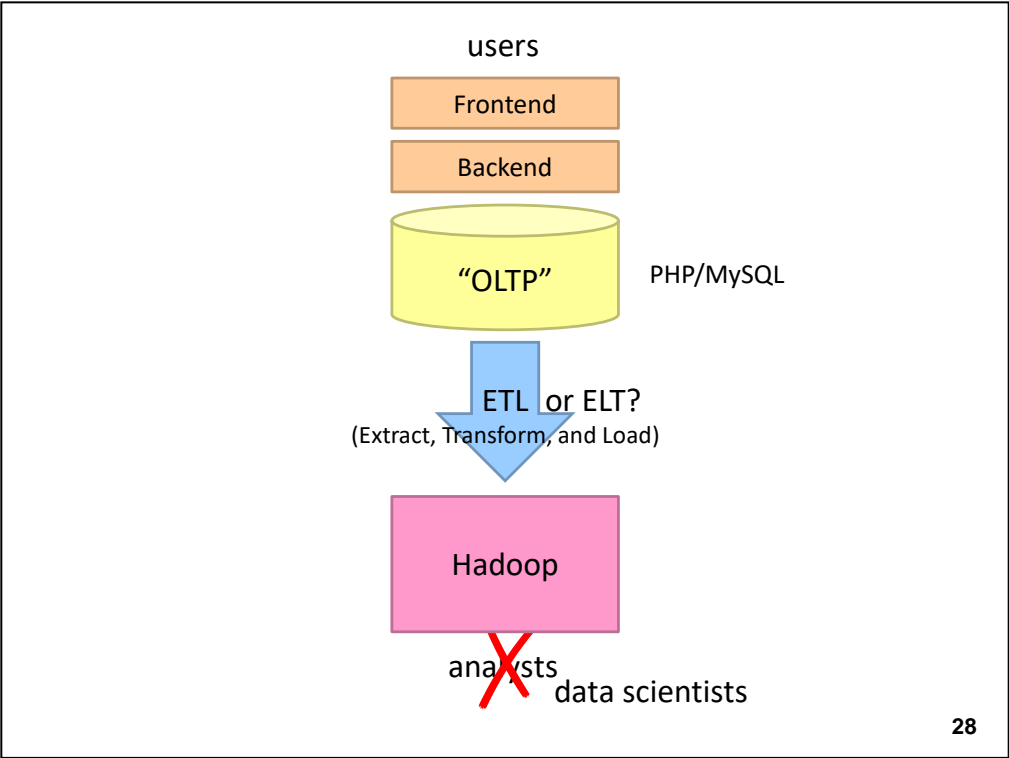
25

Facebook context?





But we have tools to deal with this, right?



What's changed?

Dropping cost of disks

Cheaper to store everything than to figure out what to throw away



5 MB hard drive in 1956

What's changed?

Dropping cost of disks

Cheaper to store everything than to figure out what to throw away

Types of data collected

From data that's *obviously* valuable to data whose value is less apparent

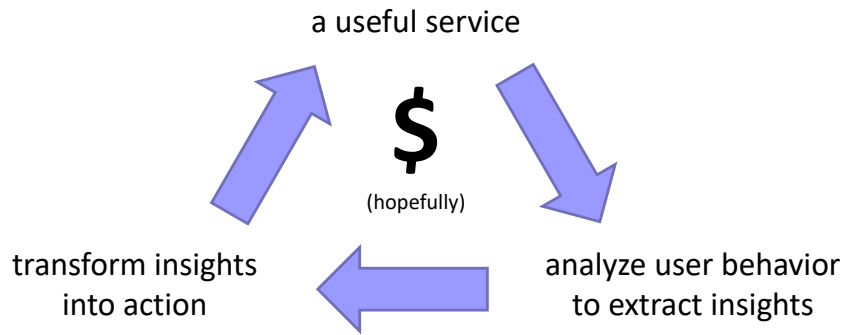
Rise of social media and user-generated content

Large increase in data volume

Growing maturity of data mining techniques

Demonstrates value of data analytics

Virtuous Product Cycle



Google. Facebook. Twitter. Amazon. Uber.

What do you actually do?

Report generation

Dashboards

Ad hoc analyses

“Descriptive”

“Predictive”

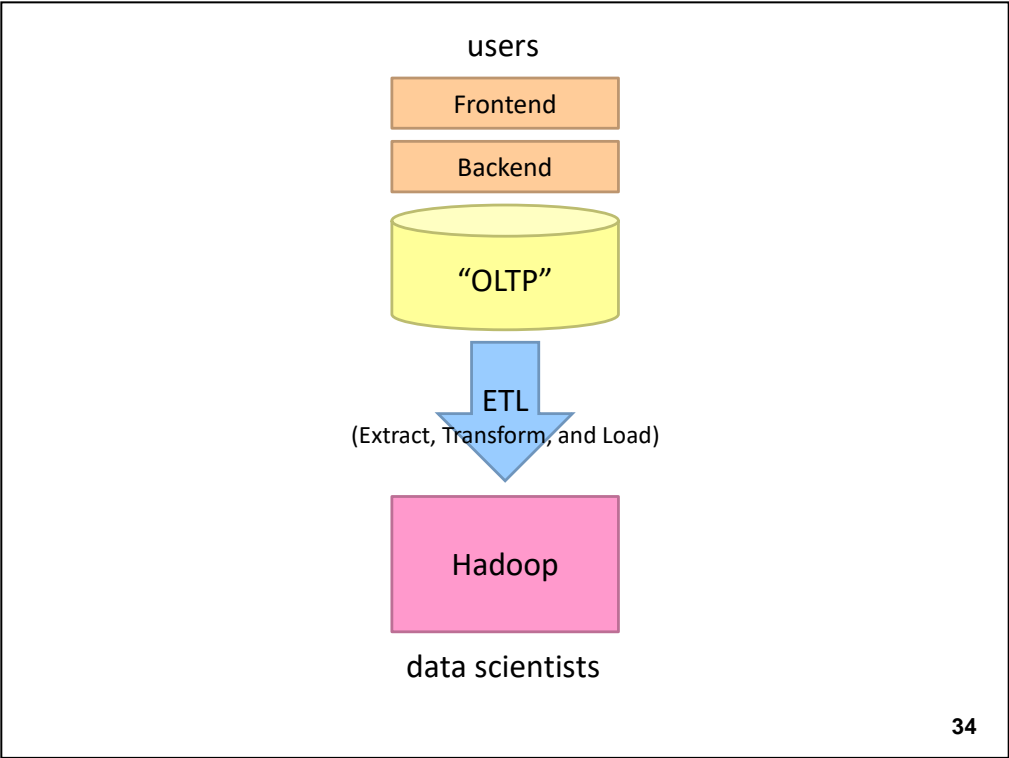
Data products

facebook®

Jeff Hammerbacher, Information Platforms and the Rise of the Data Scientist.
In, *Beautiful Data*, O'Reilly, 2009.

“On the first day of logging the Facebook clickstream, more than 400 gigabytes of data was collected. The load, index, and aggregation processes for this data set really taxed the Oracle data warehouse. Even after significant tuning, we were unable to aggregate a day of clickstream data in less than 24 hours.”

33



users

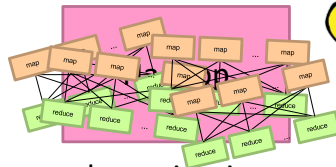
Frontend

Backend

"OLTP"

ETL

(Extract, Transform, and Load)



data scientists

Wait, so why not use a
database to begin with?

35

Why not just use a database?

SQL is awesome

Scalability. Cost.

Databases are great...

If your data has structure (and you know what the structure is)

If your data is reasonably clean

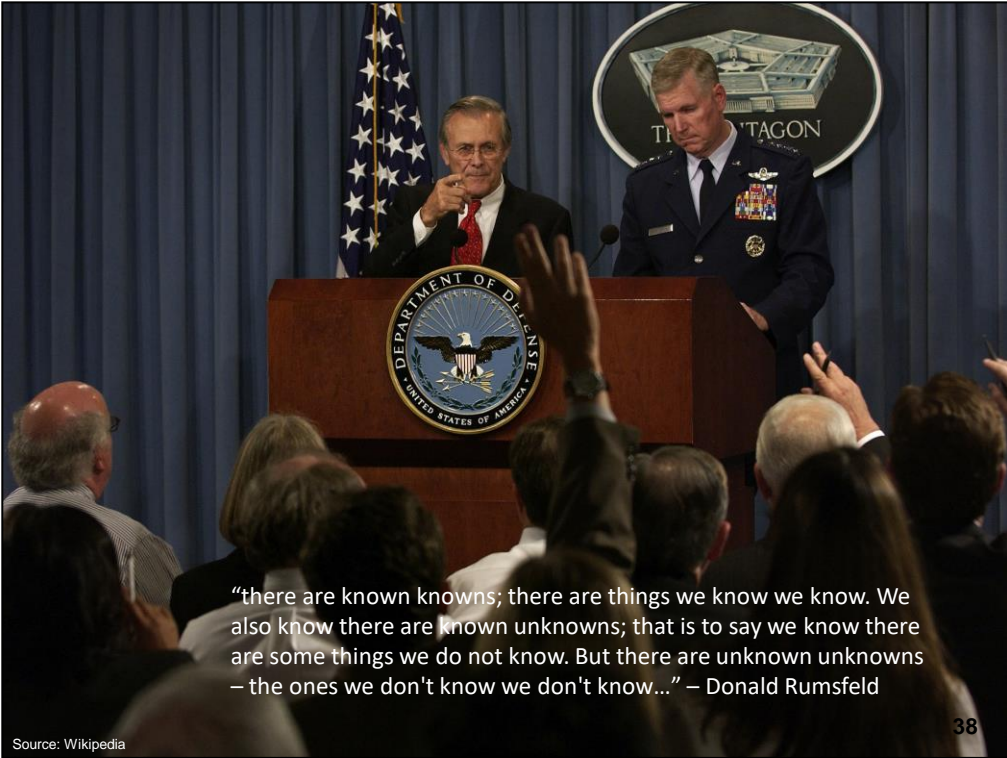
If you know what queries you're going to run ahead of time

Databases are not so great...

If your data has little structure (or you don't know the structure)

If your data is messy and noisy

If you don't know what you're looking for



“there are known knowns; there are things we know we know. We also know there are known unknowns; that is to say we know there are some things we do not know. But there are unknown unknowns – the ones we don't know we don't know...” – Donald Rumsfeld

Source: Wikipedia

38

One who knows and knows that he knows
His horse of wisdom will reach the skies

One who doesn't know, but knows that he doesn't know
His limping mule will eventually get him home

One who doesn't know and doesn't know that he doesn't know
He will be eternally lost in his hopeless ignorance!

Ibn Yamin (1286-1368)

Databases are great...

If your data has structure (and you know what the structure is)

If your data is reasonably clean

If you know what queries you're going to run ahead of time

Known unknowns!

Databases are not so great...

If your data has little structure (or you don't know the structure)

If your data is messy and noisy

If you don't know what you're looking for

Unknown unknowns!

What do you actually do?

Report generation

Dashboards

Ad hoc analyses

“Descriptive”

“Predictive”

Data products

Which are known unknowns and
unknown unknowns?

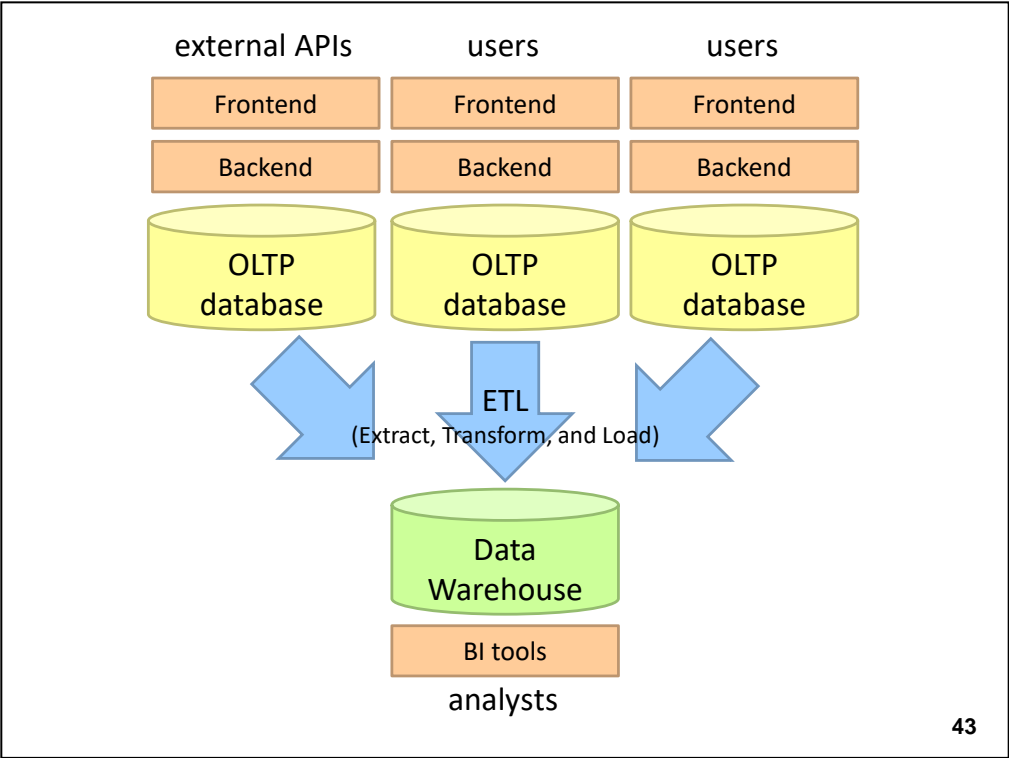
Advantages of Hadoop dataflow languages

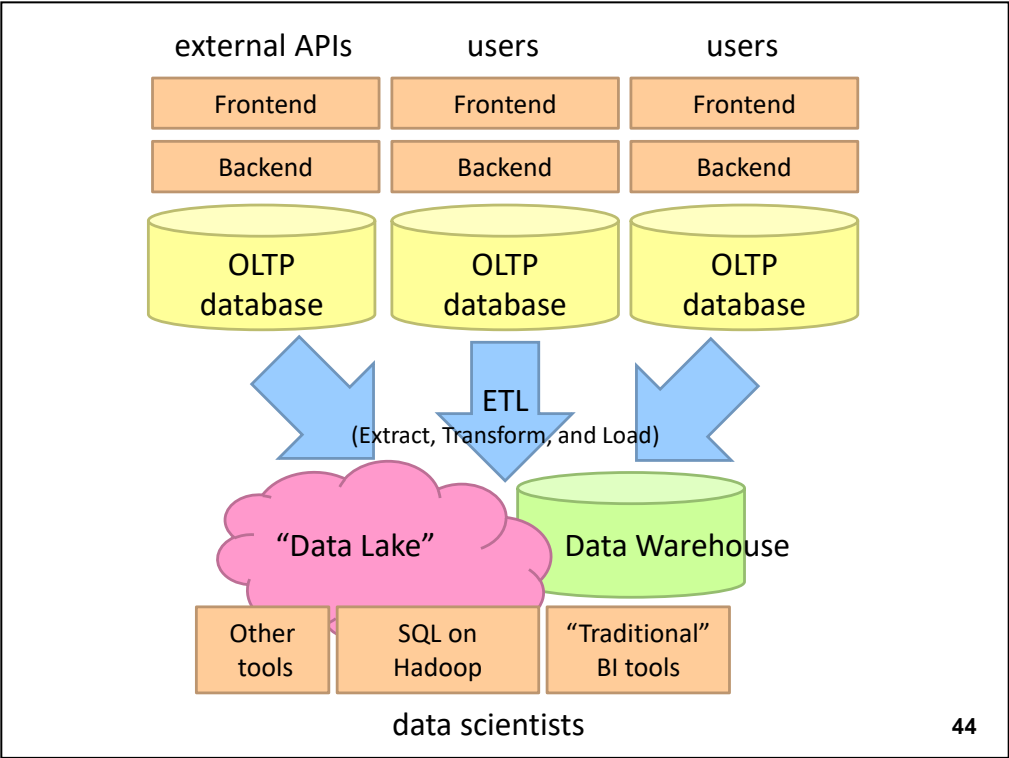
Don't need to know the schema ahead of time

Raw scans are the most common operations

Many analyses are better formulated imperatively

Much faster data ingest rate







users

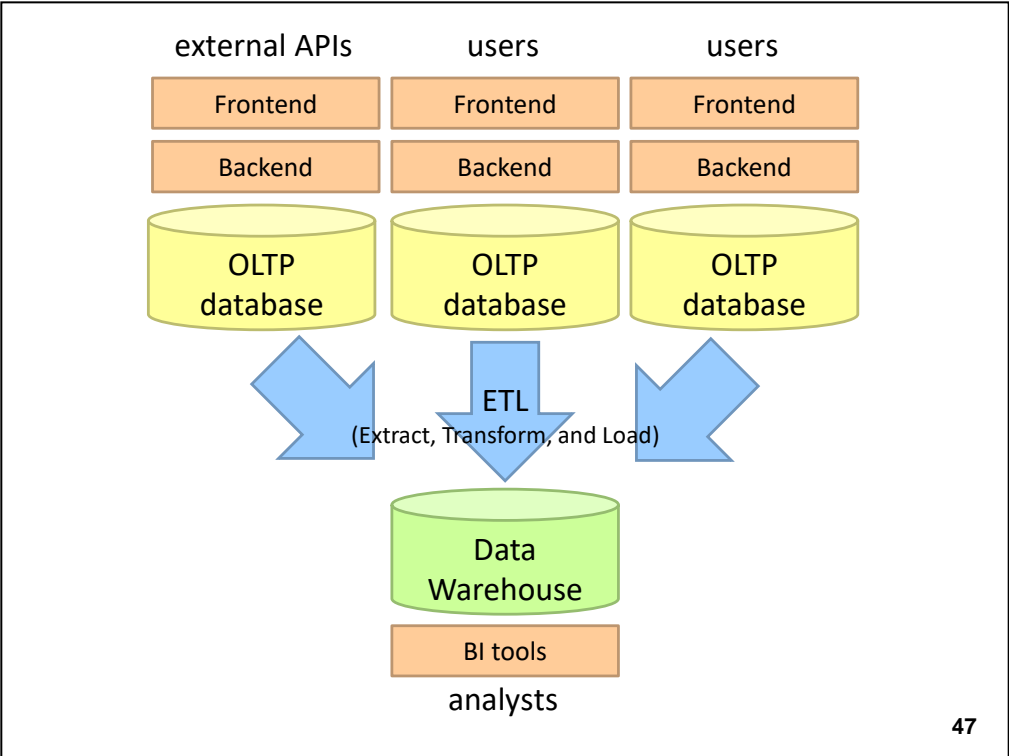
Frontend

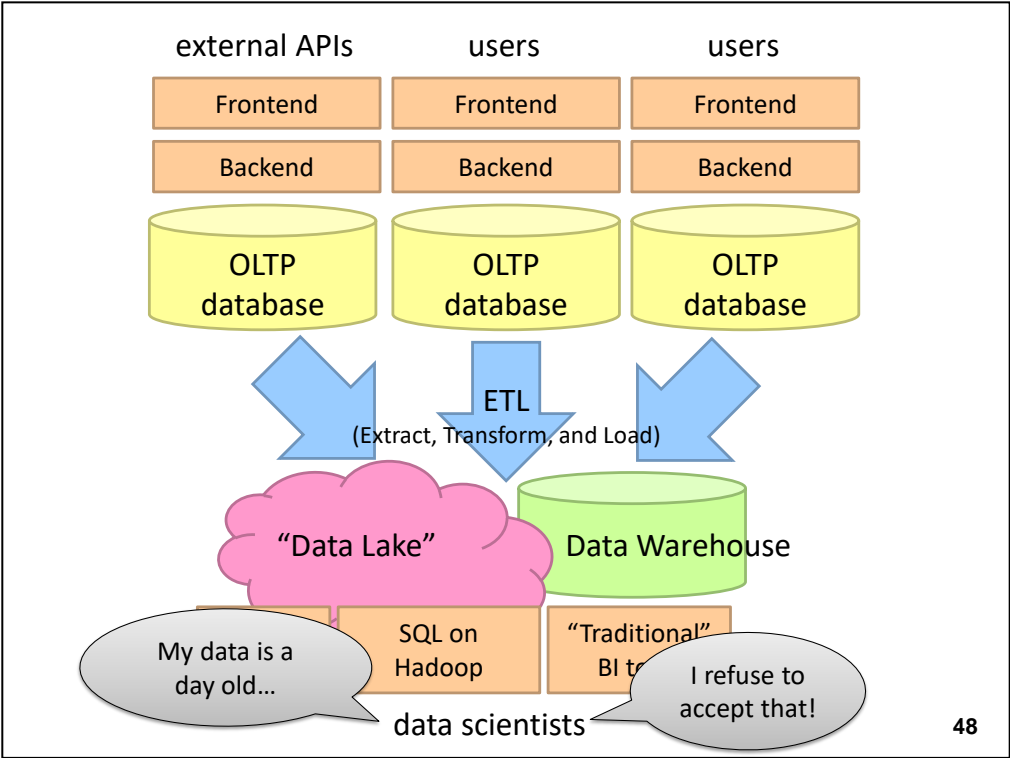
Backend

database

BI tools

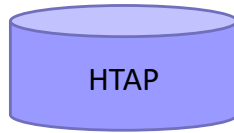
analysts







What if you didn't have to do this?



Hybrid Transactional/Analytical Processing (HTAP)

Coming back full circle?

