Data, Data Storage, Data Collection Lecture 7: Data Storage

Romain Pascual

MICS, CentraleSupélec, Université Paris-Saclay

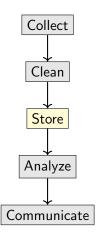
Recap

Recap: Data Wrangling

- 1 Start with initial assessment of the data quality
- 2 Handle missing data and outliers
- 3 Transform the data to prepare it for analysis
- Reshape the data for better readability or comparability with specific analysis techniques.
- 5 Keep the original data and document the cleaning steps

Introduction

Within the lifecycle



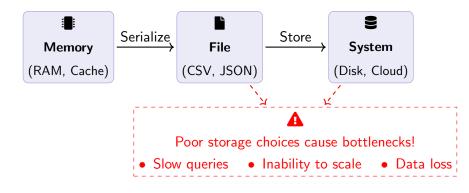
Session Objectives

At the end of this session, you should be able to:

- Distinguish between storage structures (flat files, tabular, hierarchical, columnar).
- Compare formats (CSV, JSON, Parquet) based on trade-offs (readability, size, speed, interoperability).
- Explain the differences between SQL, NoSQL, and cloud storage use cases.
- Apply a decision framework to select appropriate storage for a given scenario.
- Discuss the broader implications of storage choices (e.g., ethics, sustainability)

Why Storage Matters?

After wrangling, we need to store the data for analysis



We discuss the logic behind data storage, not just the technologies.

Storage Persistence, Strategies

In-Memory vs On-Disk





In-Memory (RAM) On-Disk (HDD/SSD)

Speed	10-100 ns	0.1-10 ms
Persistence	Volatile	Persistent
Cost (2025)	\$5-\$10 per GB	\$0.02-\$0.10 per GB
Capacity	GBs	TBs

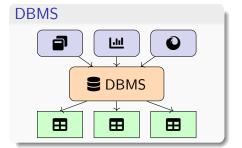
```
# In-Memory
data = {"key": "value"} # Stored in RAM

# On-Disk
import pickle
pickle.dump(data, open('data', 'wb')) # Save to disk
data=pickle.load(open('data', 'rb')) # Load from disk
```

Filesystems vs. Database Management Systems

FMS

- OS method to organize data files
- Simple folder/file hierarchy
- Keep redundant data
- Low complexity
- No built-in querying capabilities



- Software for accessing, and manipulating data
- Structured tables/collections
- Eliminates redundant data
- High complexity
- Supports complex queries and indexing

Storage Structures and Formats

Data Storage Models

Definition (Storage Model)

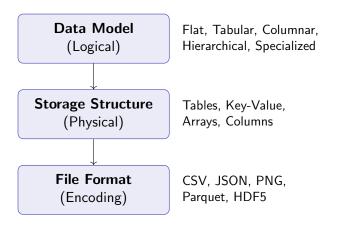
Conceptual structure or paradigm for organizing data, defining how it is logically stored and accessed.

- Flat files: Simple textual files
- Tabular: Organised in rows and columns
- Hierarchical: Nested key-value pairs

- Columnar: Data stored by columns
- Specialized: Domain-specific formats

This is about the structure of the data.

From Models to Formats



- Data model defines the logical structure.
- Storage Structure defines the physical organization.
- File format defines the encoding.

Excel (Tabular)

Characteristics:

- Widely used in business environments
- Proprietary format (.xlsx)
- 1 048 576 rows × 16 384 columns per sheet
- Supports formulas, charts, and multiple sheets
- User-friendly interface

Limitations:

- Not ideal for big data or automation
- Compatibility issues between versions

Tabular File Format Comparison

Excel

Characteristics

- Business standard (.xlsx)
- 1 048 576 rows × 16 384 columns
- Formulas, charts, multiple sheets
- User-friendly interface

Limitations

- Not for big data
- Version compatibility issues

CSV

Characteristics

- Universal, simple format
- Human-readable
- Works with almost any tool or programming language

Limitations

- No data types (everything is a string)
- Easy to break
- No complex structures

When Excel is not the right tool...

Grille d'auto-évaluation par compétence				
étudiant :				
date				
activité pédagogique : Stage de fin d'études				
				PASS
COMPETENCE Ci	QUESTIONS CLES D'EVALUATION PAR Ci : pour chaque Ci : Est-ce que l'élève (ou l'équipe) a su	Points forts	Points d'amélioration	/FAII
C1-complexité	Analyser, concevoir et réaliser des systèmes complexes			
	* analyser le système dans sa globalité, identifier les disciplines mises en jeu, caractériser les interactions internes du système, et ses dimensions extérieures ?			
Analyser un système dans sa globalité, identifier ses dimensions scientifiques, économiques, humaines, etc., modéliser avec l'échelle et les hypothèses pertinentes, résondre le problème, et concevoir tout ou partie d'un système complexe.	unneassons exerteures? ** unifore un modèle dapté, avec une échelle de modélisation adéquate, des hypothèses pertinentes, et argumenter le choix de son modèle? ** résondre un problème avec une pratique de l'approximation, simulation, observation et expérimentation, en sachant identifier les incertitudes et recorder du recur sul se arésultat sobteurs?			
	* concevoir, spécifier un cahier des charges, réaliser, tester et valider tout ou partie d'un système complexe, en sachant adapter sa démarche de imanière ilétative selon les résultats intermédiaires? * globalement promère du recuel et faire un retour d'expérience sur sa démarche ?			
C2-métier ingénieur	Développer ses compétences dans un domaine d'ingénieur et dans un métiers			
L'approfondissement disciplinaire ou sectoriel doit permettre au diplônd d'ailler au fond d'un domaine, de réfléchir de manière profonde, de comprendre les difficultés ets subtilisé d'un sigle. Cette approche de lyes expertise », permet de s'approprier des modes de pensée complexe qui poverent être ensuite transposés à d'autres socteurs, elle dome d'arrea plora probre dérice ficiencement de first apprensisages, 2 axes sont visés : le secteur d'activité, et la ligne métier transverse.	*A developer ous experites dans un domaine des extences de l'application ? **Leiseffier les domaines moments au guit de not result, en exploré les constanantes, métiger les contraintes spécifiques, et ainsi enricht l'administration de l'application de l'application de l'application de l'application par le production de l'application par le production de l'application par le production de l'application de l'application par le production de l'application de l'ap			
C3-innover entreprendre	Agir, entreprendre, innover en environnement scientifique et technologique			
La compétence entrepressur se comprend dans un sens large (entrepressur de soi-même, capacité d'action, intrapressur, cristeur d'entreprés, ext.) o boserver et s'autorier à critiquer le monde tel qu'il est, remettre en cause ses hypothèses de départ, proposer des alternatives indignant triaques el meritande, mettre en œuvre concrètement des idées novatrices, industrialisation pour délivrer des résultats tangibles.	*questionner la formulation du problème posé pour comprende le besoin sous-jacent dans un contract plus large ? "propose de sides nouvelles pour réponde na problème et/ou quant la démanche à autre pour en melliourer l'efficacié ? "mounteur, une ces démoneturels, autres qu'et duite l'14 liababble à le solution perchaige, chomaige, fantaire, échique) 2º l'infort side pour le bénéficiaire? "request se ca problème ne del résolut ailleurs, et de quelle masière ? (ne pas réinventer la noue et ne positionner pur rapport a d'autres solution)			
C4-création valeur	Avoir le sens de la création de valeur pour son entreprise et ses clients			
Partant d'une définition pour l'entreprise d'augmentation de la productivité et de la ternabilité, la Création de Valeur est élargie in à l'apport de progrès et d'amélioration dans des domaines comme l'efficacité d'un processos, la performance ou la fabilité u'une solution technique, ou l'impact entrommentant la ouscière ; de telle sorte que le (ou les bénéficiaires), individués ou collectifs (entreprises, organisations, société) missent no constater et ne meure les effect en meure les effet.	*congrencie, per un dialogue portions avec le demodere (en clarit, les nejons de la problématique, identifier les éventrelles autres parties remantes et les actue de major didit oirées et ventre.? **referentes insplaement la cristium de values attendues par le demondere, obtenir son approblematique ser entre conveille formalistics et convenir les de la façon de nameur ? **propose et rejonser de manier consciuntes aux obtains confirme à la reformalistius de besoin du demodere (voire plusieurs solutions) ? **elevate constantium en formation de montantes aux obtains confirme à la reformalistius de besoin du demodere (voire plusieurs solutions) ? **elevate constantium en formation de montantes aux obtains confirme à la reformalistius de besoin du demodere (voire plusieurs solutions) ? **elevate constantium en formation de miscentes de recitories de values convenir.			

JSON: JavaScript Object Notation (Hierarchical)

Characteristics

- Standard for API and config.
- Nested key-value structure
- Human-readable format

Limitations

- Verbose syntax
- No support for comments
- Inefficient for large datasets

```
1 {
2    "title": "The Hitchhiker's Guide to the Galaxy",
3    "author": "Douglas Adams",
4    "year": 1979,
5    "characters": [
6    { "name": "Arthur Dent", "species": "Human", "role": "Protagonist"},
7    { "name": "Ford Prefect", "species": "Betelgeusian", "role": "Researcher"},
8    { "name": "Marvin", "species": "Robot", "role": "Android"}
9    ],
10    "notable_quote": "Don't Panic.",
11    "answer_to_life": 42,
12    "adaptations": ["radio", "TV", "film"]
13 }
```

Parquet: (Columns)

Characteristics:

- Columnar format optimized for big data
- High performance compression and encoding schemes
- Efficient for analytical queries
- Works well with Spark, Hadoop, and other big data tools

Limitations:

Binary format (not human-readable)

Row vs. Column Storage

Name	Age	City	
Alice	28	Paris	
Bob	32	London	
Charlie	25	Berlin	



	(Alice		
Row 1	28		
	Paris		
	Bob		
Row 2	32		
	London		
Row 3	Charlie ¦		
	25		
	Berlin		
Row-based access			
reads entire row			

Row Storage

Column Storage Alice Name Bob Charlie **í** 28 32 Age 25 **Paris** City London Berlin

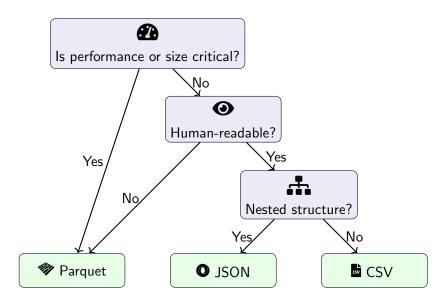
Column-based access

Summary: Comparing Common File Formats

Criterion	Excel	CSV	JSON	Parquet
Fast to read/write	X	•	•	✓
Space-efficient		X	X	✓
Compatibility	×	✓	✓	X
Human-readable	✓	✓	✓	X
Supports nested data		X	✓	X
Good for sharing		✓	✓	•
Good for analysis	✓	•	X	✓

Excel is user-friendly but proprietary, **CSV** and **JSON** are widely compatible and flexible, **Parquet** is optimized for large-scale analytics.

Choosing the Right File Format



Databases

SQL: Relational Databases

Key Concepts

- Tables with rows and columns
- Keys (primary, foreign)
- Atomicity, Consistency, Isolation, Durability
- Underlying formalization: Relational algebra

users SELECT name, email Input id email name age Alice 37 alice@example.com 1 ORDER BY name: Frank frank@example.com 3 7oe 21 zoe@example.com Result David 20 david@example.com 30 Eve eve@example.com email name Bob 19 bob@example.com Alice alice@example.com Grace 27 grace@example.com Charlie charlie@example.com 32 charlie@example.com Charlie eve@example.com Eve

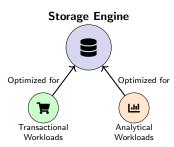
Understanding Storage Engines

You will not build a storage engine from scratch . . .

...But you will need to

- Select the right engine.
- Tune it for your workload.

So you need to understand what is under the hood!



In particular, storage engines are optimized for different workloads

Databases: OLTP vs. OLAP

OLTP (Online **Transaction** Processing)

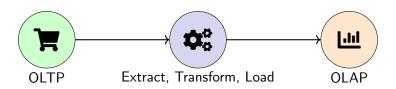
- Fast, frequent reads/writes
- Small, simple transactions
- Indexes for quick lookups

Example: Banking transactions

OLAP (Online **Analytical** Processing)

- Complex queries and aggregations
- Large datasets, few columns
- Calculates statistics

Example: Business intelligence



NoSQL: Beyond Relational

• Volume: Massive data growth

Velocity: High-speed data generation

• Volume: Massive data growth

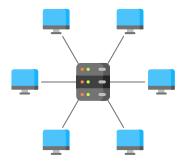
Velocity: High-speed data generation



Centralized

• Volume: Massive data growth

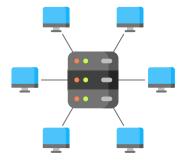
Velocity: High-speed data generation



Centralized

• Volume: Massive data growth

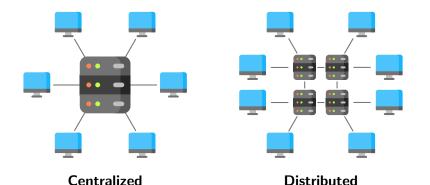
Velocity: High-speed data generation



Centralized

• Volume: Massive data growth

Velocity: High-speed data generation



The Evolution of Database Systems: From SQL to NoSQL

Historical Context

- Pre-SQL Era: Early databases (1960s-70s) were often hierarchical or network models without standardized schema
- 1970: Edgar F. Codd proposed the relational model
- early 1970's: SEQUEL (Structured English Query Language), designed by IBM
- 1979: First SQL database (Oracle) commercialized
- 1980s-90s: SQL became the dominant standard for structured data

The Rise of "Not Only SQL"

As a response to the three 3 V's

- Architecture Shift: From centralized to distributed systems
- Data Shift: Unstructured or rapidly changing data

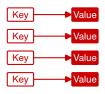
BASE Principles (vs. ACID)

- Basically Available: System appears to work most of the time
- Soft state: Data may change over time without explicit updates
- Eventual consistency: Data will be consistent... eventually

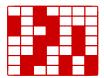
Why NoSQL?

- Better suited for large-scale, distributed data
- More flexible schema design
- Horizontal scalability
- Optimized for specific data models and access patterns

Four families of NoSQL



Key-Value Simple, fast lookups



ColumnOptimized for analytics

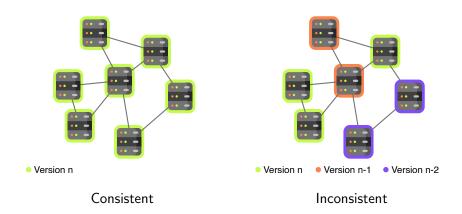


DocumentNested KV pairs



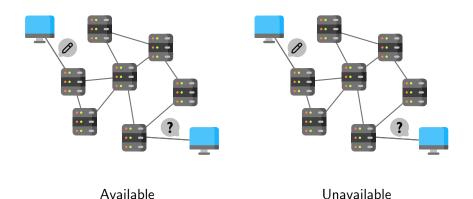
Graph Relationships

The CAP Theorem and its Implications



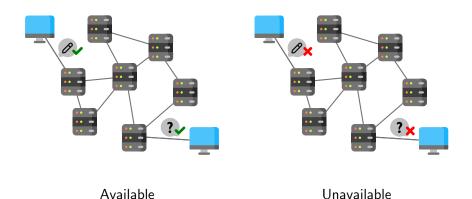
Consistency (C):

- Every read receives the most recent write (or an error).
- All nodes see the same data at the same time.



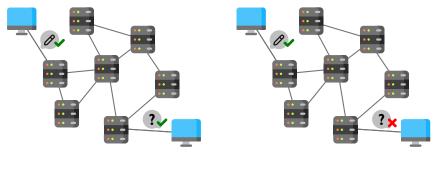
Availability (A):

 Every request received by a non-failing node in the system must result in a response



Availability (A):

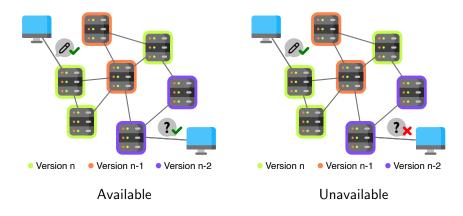
 Every request received by a non-failing node in the system must result in a response



Available Unavailable

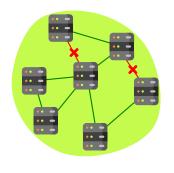
Availability (A):

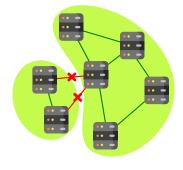
 Every request received by a non-failing node in the system must result in a response



Availability (A):

 Every request received by a non-failing node in the system must result in a response, even if it is not with the latest data.





No partition

Partition

Partition Tolerance (P):

- The system continues to function despite network partitions.
- Tolerates arbitrarily many message losses.

Theorem (CAP¹)

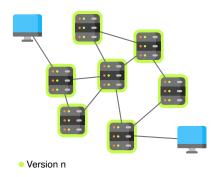
Any distributed system can have at most two of the following three properties:

- Consistency
- Availability
- Partition Tolerance



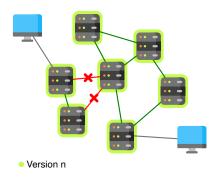
Eric Brewer

¹Armando Fox and Eric Brewer. "Harvest, Yield, and Scalable Tolerant Systems". In: Proceedings of the Seventh Workshop on Hot Topics in Operating Systems. Mar. 1999, pp. 174–178



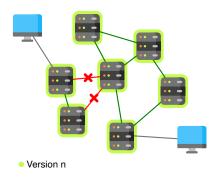
• By contradiction, consider a system *S* that fulfills the three properties.

²Seth Gilbert and Nancy Lynch. "Brewer's Conjecture and the Feasibility of Consistent, Available, Partition-Tolerant Web Services". In: **SIGACT News** 33.2 (2002), pp. 51–59.



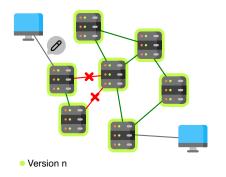
2 Partition the system into two subsystems S_1 and S_2 .

²Seth Gilbert and Nancy Lynch. "Brewer's Conjecture and the Feasibility of Consistent, Available, Partition-Tolerant Web Services". In: **SIGACT News** 33.2 (2002), pp. 51–59.



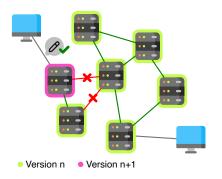
3 By **partition tolerance**, the system continues to function.

²Seth Gilbert and Nancy Lynch. "Brewer's Conjecture and the Feasibility of Consistent, Available, Partition-Tolerant Web Services". In: **SIGACT News** 33.2 (2002), pp. 51–59.



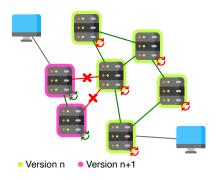
4 Client c_1 requests to write new data to a node in S_1 .

²Seth Gilbert and Nancy Lynch. "Brewer's Conjecture and the Feasibility of Consistent, Available, Partition-Tolerant Web Services". In: **SIGACT News** 33.2 (2002), pp. 51–59.



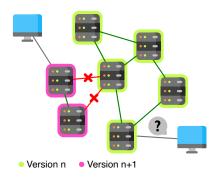
5 By **availability**, the system handles the request and writes the new data.

²Seth Gilbert and Nancy Lynch. "Brewer's Conjecture and the Feasibility of Consistent, Available, Partition-Tolerant Web Services". In: **SIGACT News** 33.2 (2002), pp. 51–59.



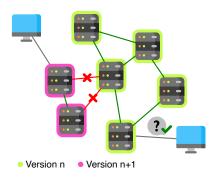
6 Since the network is partitioned, the data cannot be replicated to the nodes in S_2 .

²Seth Gilbert and Nancy Lynch. "Brewer's Conjecture and the Feasibility of Consistent, Available, Partition-Tolerant Web Services". In: **SIGACT News** 33.2 (2002), pp. 51–59.



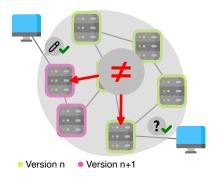
7 Client c_2 requests to read the data to a node in S_2 .

²Seth Gilbert and Nancy Lynch. "Brewer's Conjecture and the Feasibility of Consistent, Available, Partition-Tolerant Web Services". In: **SIGACT News** 33.2 (2002), pp. 51–59.



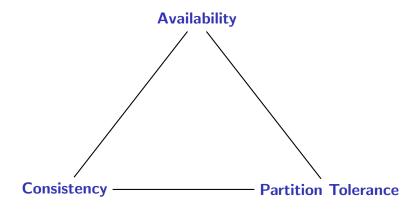
8 By availability, the system handles the request and outputs the old data.

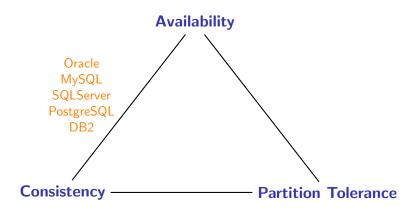
²Seth Gilbert and Nancy Lynch. "Brewer's Conjecture and the Feasibility of Consistent, Available, Partition-Tolerant Web Services". In: **SIGACT News** 33.2 (2002), pp. 51–59.



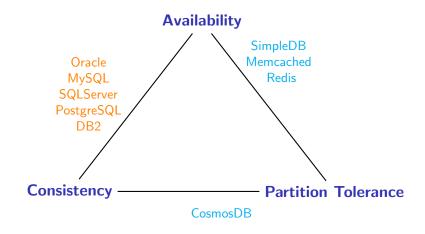
The second request yields an inconsistent answer.

²Seth Gilbert and Nancy Lynch. "Brewer's Conjecture and the Feasibility of Consistent, Available, Partition-Tolerant Web Services". In: **SIGACT News** 33.2 (2002), pp. 51–59.

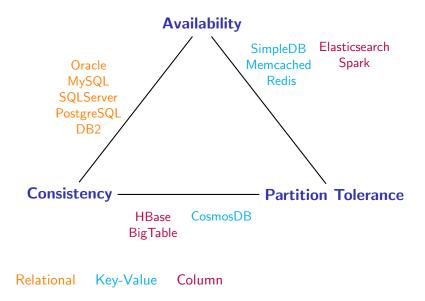


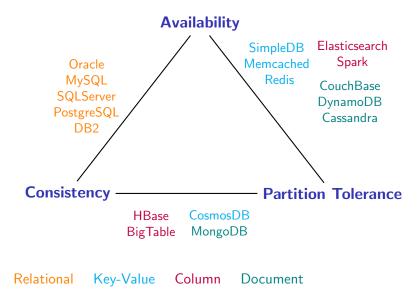


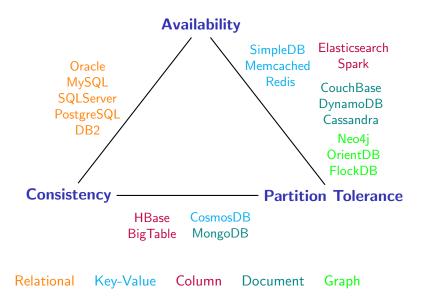
Relational

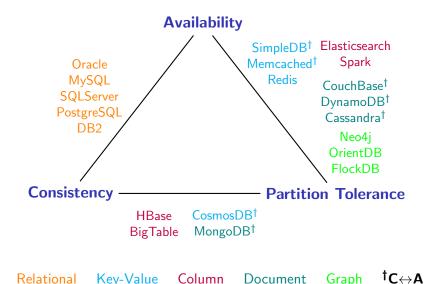


Relational Key-Value









Example: Bank Database

A financial database stores each client's bank balance.

Clients interact via an ATM to:

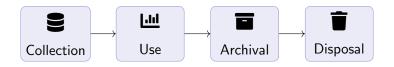
- Check balance
- Withdraw money
- Deposit money



Given these operations, which aspect of the CAP theorem would you prioritize: consistency, availability, or partition tolerance?

Long Time Storage

Long-Term Storage in the Data Science Lifecycle



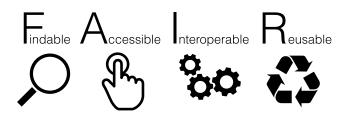
Why Is Data Archive Part of Data Science?

- Reproducibility of analyses and results
- Compliance with data retention policies
- Future reuse for new analyses or model training
- Preservation of data provenance and metadata

Challenges and Principles of Long-Term Storage

Technical Challenges:

- Format Compatibility: Ensuring readability by future tools
- Metadata Preservation: Maintaining reproducibility
- Data Integrity: Preventing corruption over time
- Storage Costs: Balancing accessibility and expense



Data Retention and Disposal





Formats for Archival

Criteria for Archival Formats:

- Open standards (e.g., CSV, JSON, Parquet)
- Self-describing (e.g., JSON with schema, Parquet with metadata)

```
"metadata": {
    "format": "JSON",
    "created": "2025-10-12"
},
    "data": "..."
"]
```

Example: Archiving a Data Science Project

```
1 project/
2 |-- data/
3 | |-- raw/
                            # Original immutable data
 | |-- processed/
                            # Cleaned/processed data
 | \-- README.md
                            # Data documentation
6 |-- models/
  |-- model.onnx
                       # Serialized model
 | \-- metadata.json  # Model metadata
 |-- notebooks/
                            # Analysis notebooks
  -- requirements.txt
                            # Python dependencies
 \-- README.md
                            # Project documentation
```

Key Files to Archive

- Raw and processed data (in open formats)
- Trained models and their metadata
- Analysis code and notebooks
- Environment specifications
- Documentation and data dictionaries

Conclusion

Ethical and Practical Considerations

Ethical Implications

- A Privacy: Who has access to the data?
- **Accessibility**: Is the data available to those who need it?
- Sustainability: What's the environmental impact?

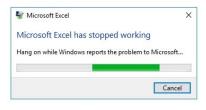
Practical Considerations

- U Long-term storage: Will you need this data in 5 years?
- \$ Cost: What's the total cost of ownership?
- **\$\pi** Maintenance: Who will maintain this storage system?

Common Pitfalls and Best Practices

Pitfalls

- Using Excel for big data
- Over-engineering (e.g., graph DB for simple data)
- Ignoring costs (e.g., cloud storage fees)



Best Practices

- 1 Start simple (CSV/SQL), scale up as needed
- 2 Document storage schema and access patterns
- 3 Test with real data early

Takeaways: Data Storage

- Understand the difference between storage models (flat, tabular, hierarchical, columnar) and formats (CSV, JSON, Parquet)
- 2 Choose the right format based on your use case: Excel for business, CSV/JSON for sharing, Parquet for analytics
- 3 Consider the trade-offs between readability, size, speed, and interoperability when selecting a format
- For long-term storage, prioritize open standards and self-describing formats with proper metadata
- 5 Document your storage decisions and maintain data provenance for reproducibility
- Consider ethical implications including privacy, accessibility, and sustainability