

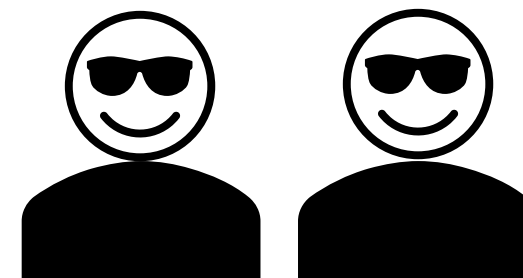
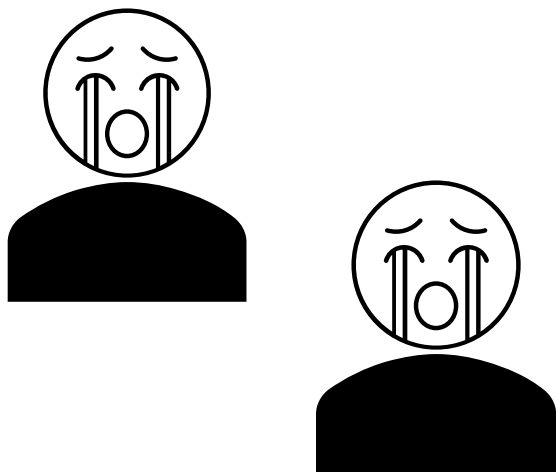
Federated Analysis of FAIR data

**César Bernabé and
Daphne Wijnbergen**

If you need to book a flight, how would you search for the best deal?

- **A)** Visit multiple airline websites separately
- **B)** Use a flight comparison tool like Skyscanner, Expedia
- **C)** Ask a travel agent
- **D)** Just book with the first airline you know

Example



Example

- **Several airlines: AMS – LCA**
 - Easyjet
 - Transavia
 - Sky Express
- Austrian Airways
- Aegean
- AirSerbia
- ...
- **Strategy one:** check one by one, download data and put it all in an Excel sheet

Example

Heenvlucht

Cyprus Larnaca
Momenteel bekijken
✈️

✈️ Outbound flight

Amsterdam (Schiphol) – Cyprus (Larnaca)

Monthly overview 📅

Apr 2025						
Mon 21	Tue 22	Wed 23	Thu 24	Fri 25	Sat 26	Sun 27
No flights available	From € 334	No flights available	From € 226	No flights available	From € 334	From € 334

◀ **vrijdag**
18 juli

Thu 24 Apr 2025

Geen vluchten
beschikbaar

🕒 14:45 ✈️ 20:00

Flight number
HV5313

10+ tickets available at this price

€ 226

Select

LAAGSTE TARIEF

€ 90,49 ✓

1 stoel over voor

beschikbaar

Example

- **Several airlines: AMS – LCA**
 - Easyjet
 - Transavia
 - Sky Express
 - Austrian Airways
 - Aegean
 - AirSerbia
 - ...
- **Strategy one:** check one by one, download data and put it all in an Excel sheet
- **Strategy two:** use a price comparison service (which performs federated queries!)

Example: Skyscanner



De uitstoot bij deze vlucht is **12% minder CO2e** dan bij een normale vlucht op deze route

easyJet	07:00 AMS	4u 15 Direct	→ 12:15 LCA	12 aanbiedingen van € 327
easyJet	12:50 LCA	4u 45 Direct	→ 16:35 AMS	Selecteren →

De uitstoot bij deze vlucht is **23% minder CO2e** dan bij een normale vlucht op deze route

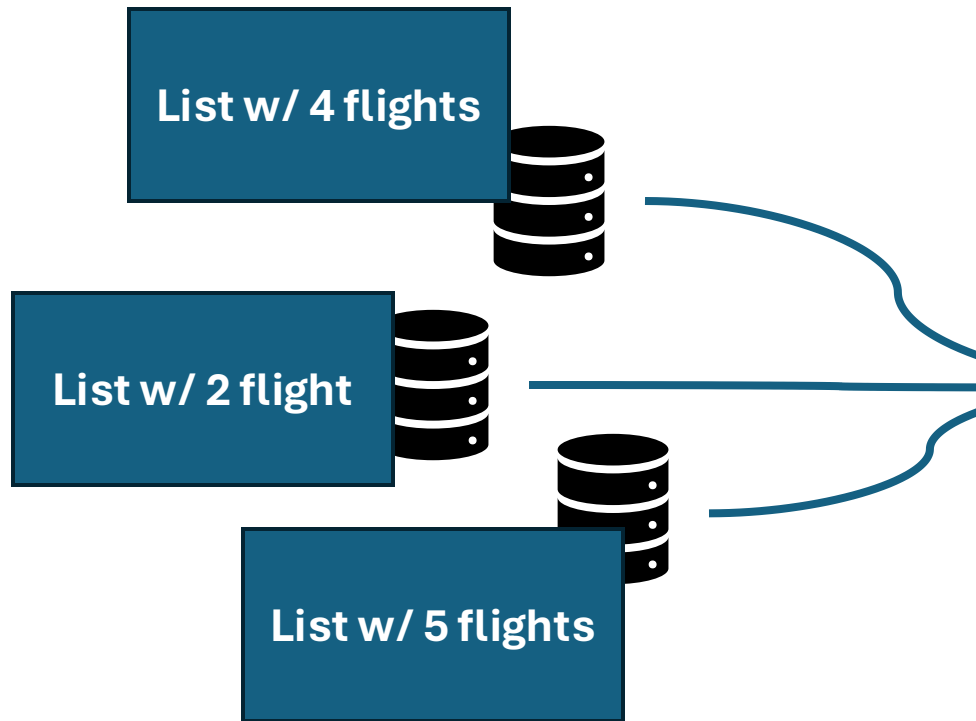
transavia	06:15 AMS	4u 15 Direct	→ 11:30 LCA	13 aanbiedingen van € 434
easyJet	12:50 LCA	4u 45 Direct	→ 16:35 AMS	Selecteren →

De uitstoot bij deze vlucht is **14% minder CO2e** dan bij een normale vlucht op deze route



SKY express	14:10 AMS	5u 45 1 tussenstop ATH	→ 20:55 LCA	15 aanbiedingen van € 382
easyJet	12:50 LCA	4u 45 Direct	→ 16:35 AMS	Selecteren →

easyJet	08:50 AMS	12u 50 1 tussenstop MXP	→ 22:40 LCA	4 aanbiedingen van € 341
easyJet	12:50 LCA	4u 45 Direct	→ 16:35 AMS	Selecteren → ! Overstap



Example: HemaSKY



De uitstoot bij deze vlucht is **23% minder CO₂e** dan bij een normale vlucht op deze route

	06:15 AMS	4u 15 Direct	11:30 LCA	13 aanbiedingen vanaf € 434 Selecteren →
	12:50 LCA	4u 45 Direct	16:35 AMS	

De uitstoot bij deze vlucht is **14% minder CO₂e** dan bij een normale vlucht op deze route

	14:10 AMS	5u 45 1 tussenstop ATH	20:55 LCA	15 aanbiedingen vanaf € 382 Selecteren →
	12:50 LCA	4u 45 Direct	16:35 AMS	

Example

- **The point is:** for César and Daphne, the whole search as perceived as an unique search
 - Prices were retrieved in real time
 - No need to copy data to a local repository
 - Results presented harmonised and unified
 - Results combined to be more valuable (e.g. flight in with Transavia and fly out with EasyJet)

Summary

1. A query is sent to multiple independent databases.

- Algorithms can be sent to databases instead of queries

2. Each database processes the request and returns results.

3. Results are combined and presented in a unified format.

4. Data remains in its original location, ensuring security.

What about rare diseases research?

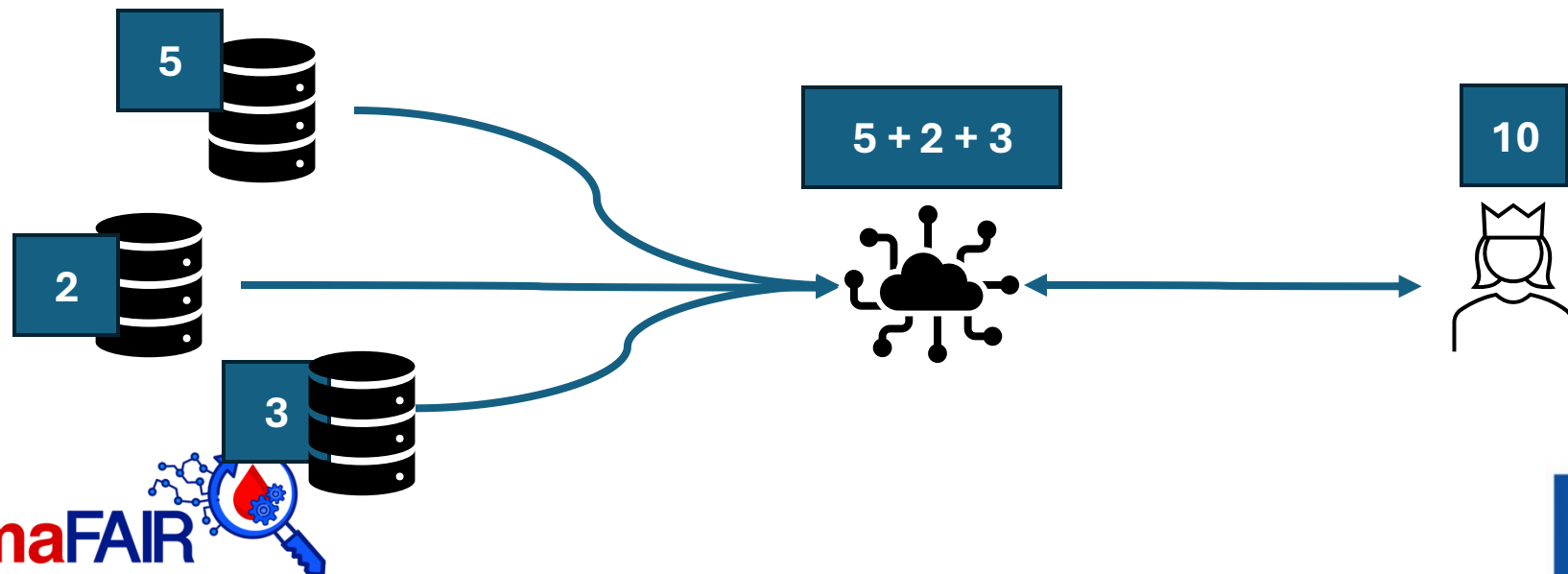
- **Patient Cohort Identification**

“Retrieve anonymized counts of patients diagnosed with Duchenne Muscular Dystrophy (DMD) across multiple hospitals, filtering by age group and presence of a specific genetic mutation.”

What about rare diseases research?

- **Patient Cohort Identification**

*“Retrieve **anonymized counts** of patients diagnosed with Duchenne Muscular Dystrophy (DMD) across multiple hospitals, filtering by age group and presence of a specific genetic mutation.”*



What are the implications?

- There must be an agreement on minimal infrastructure, common data terms and legal interfaces to allow for federated querying
- For instance, to query for “**Duchenne muscular dystrophy**”
 - Different languages: Distrofia Muscular de Duchenne, Μυϊκή Δυστροφία Duchenne
 - Different means to capture information: “X-linked muscular dystrophy with abnormal dystrophin”, “Duchenne and Becker muscular dystrophy”, SNOMED:240048000, ORPHA:262

Does FAIR data enable federated?

- **FAIR Data** facilitates federated analysis, as:
 - Relevant datasets are **found** by machines
 - **Access conditions** are explicit
 - Machines now if they can send and receive queries
 - Results can be **interoperated** as data is harmonised
 - Results are presented using relevant formats for **reuse**

Part II

Infrastructure for federated analysis

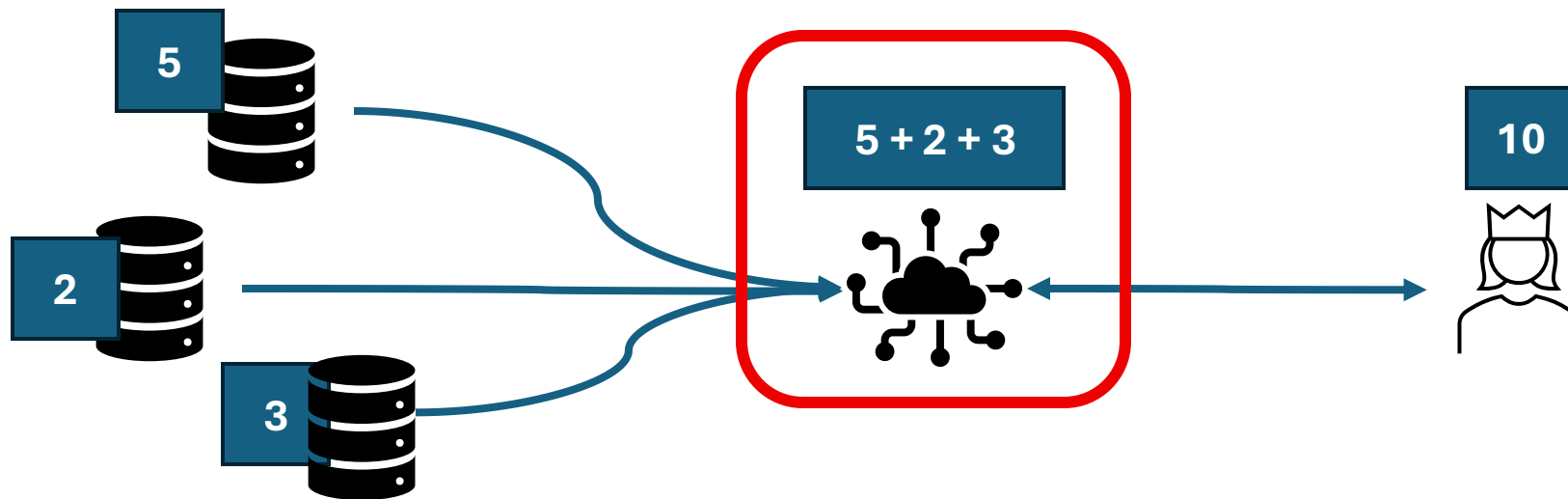


Funded by
the European Union



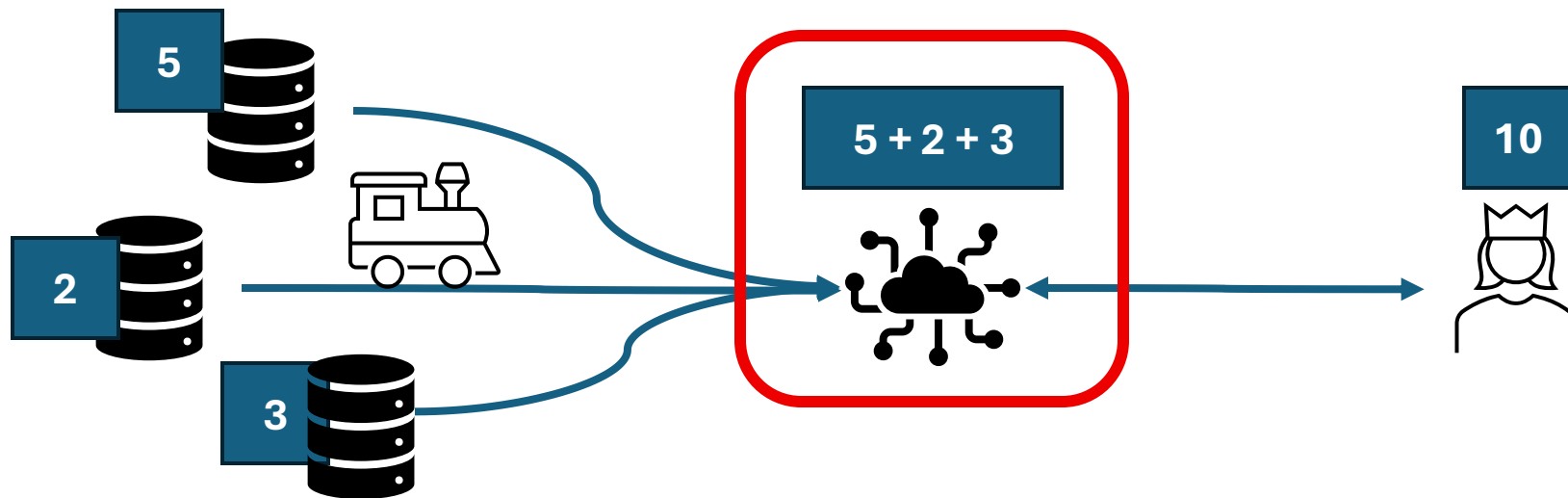
Infrastructure for federated analysis

- What is happening here?
- What components are necessary to achieve this?



Data Train Analogy

- Questions and answers being moved around will be represented as trains

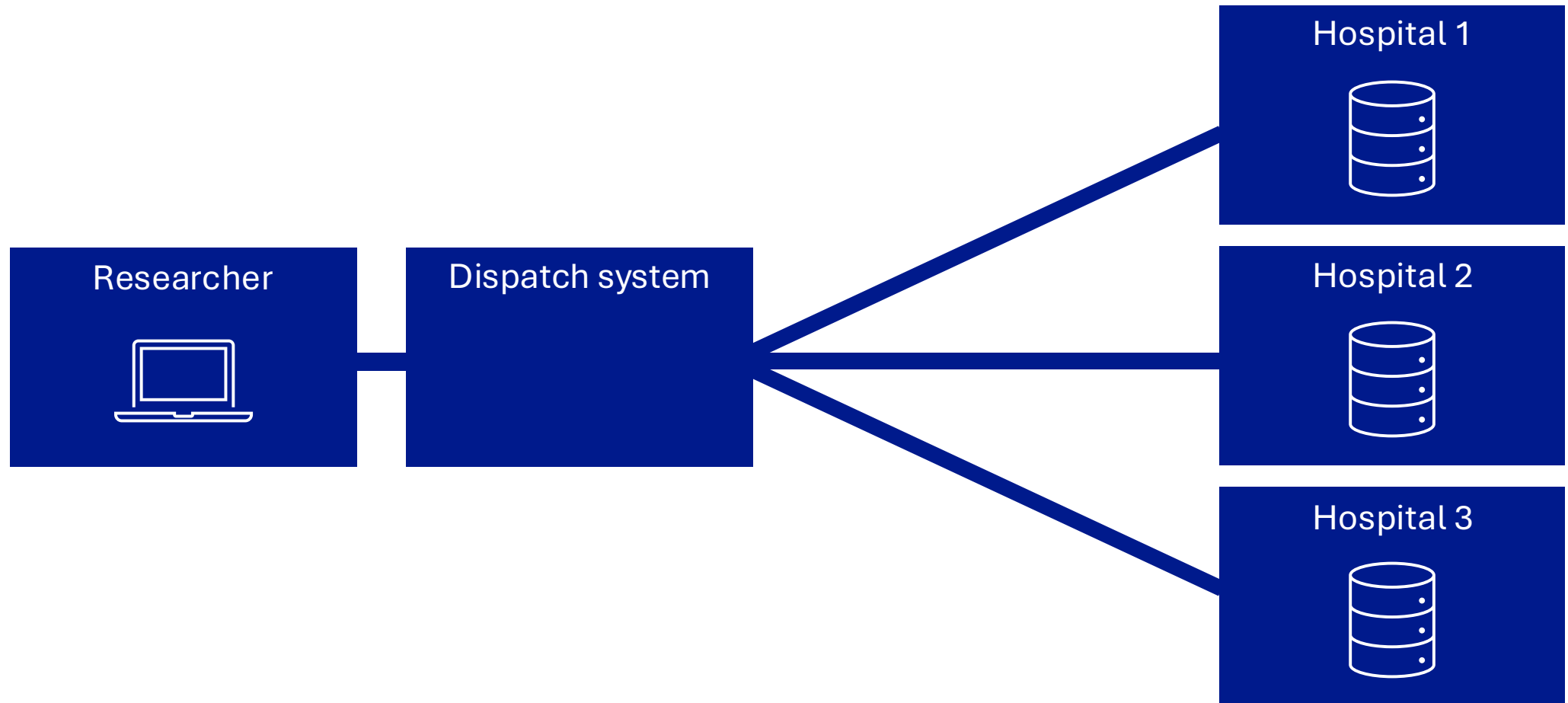


How many patients with DMD?

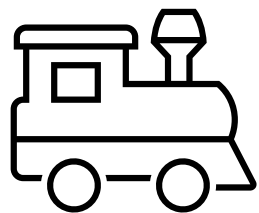
Starting point

- Three hospitals have data about patients with DMD
- We want to ask the question: how many patients are there with DMD?

Starting point

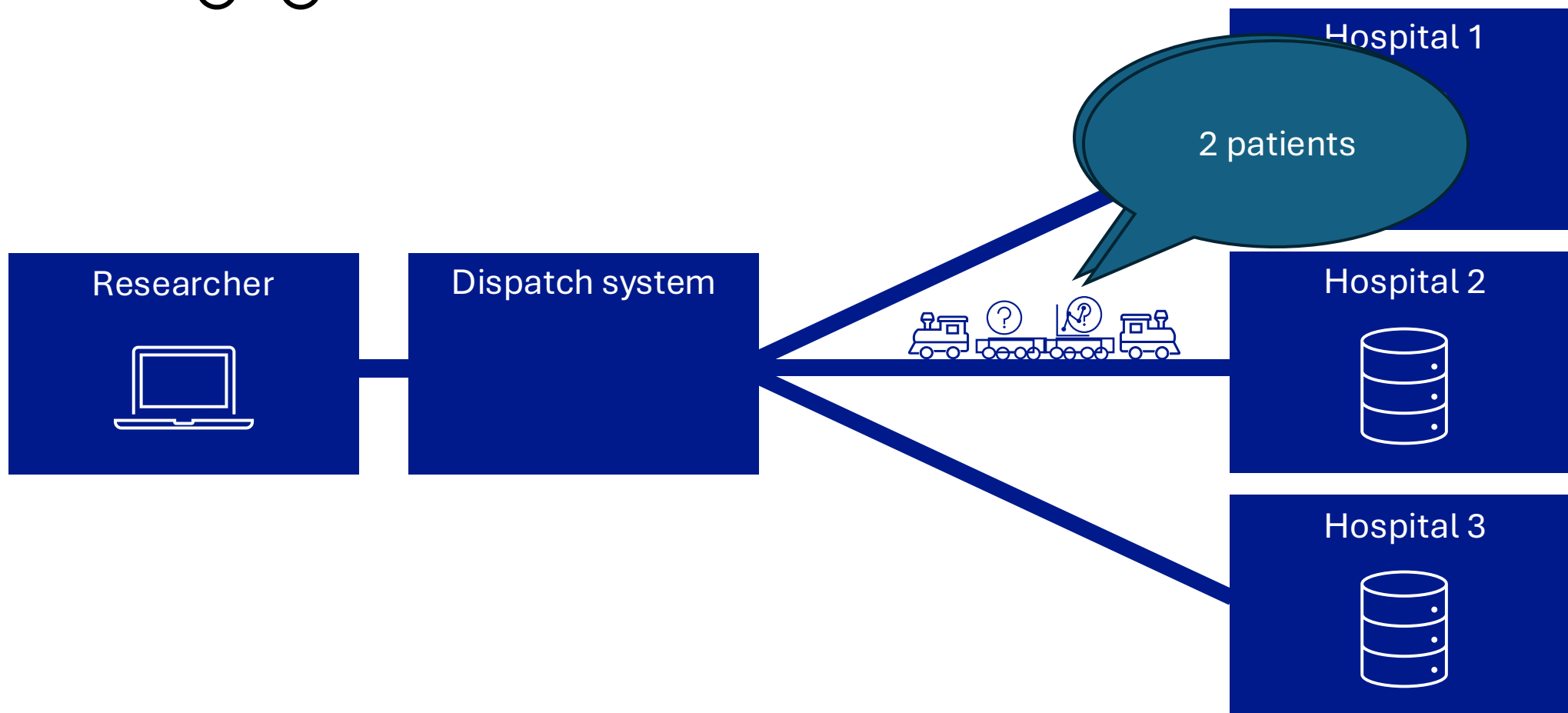
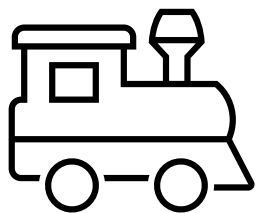


Trains

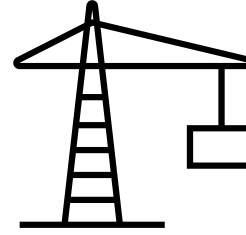


- We can use the analogy of trains
- These trains carry a question to a hospital, and carry a result back to the researcher

Trains

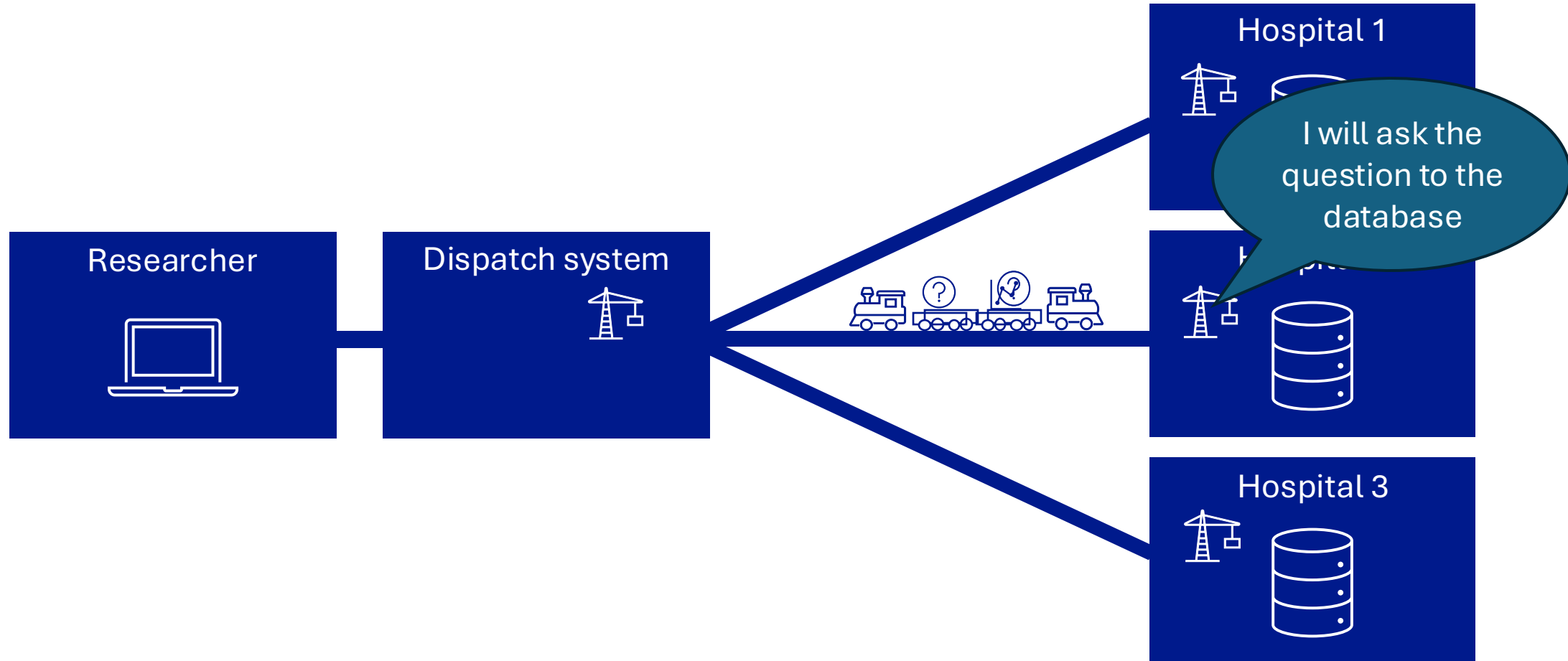
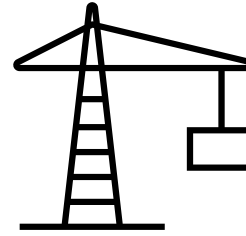


Interaction mechanism

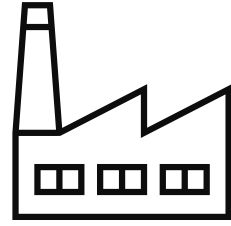


- A train can't connect to a database directly
- An interaction mechanism should exist between the train and the data
- This mechanisms needs to perform several actions:
 - Obtain the question from the train
 - Run this question on the database/file
 - Return the answer to the train
 - And more..

Interaction mechanism

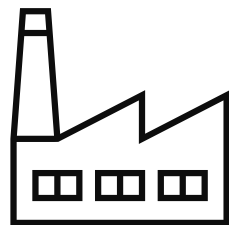


Train depot



- Trains with algorithms that can be send out should be stored somewhere

Train depot



I want to use
the DMD count
train

Researcher



Dispatch system



Hospital 1

How many
patients with
DMD?

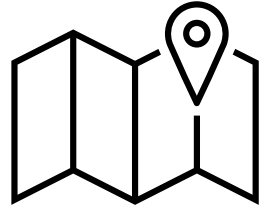
Hospital 2



Hospital 3

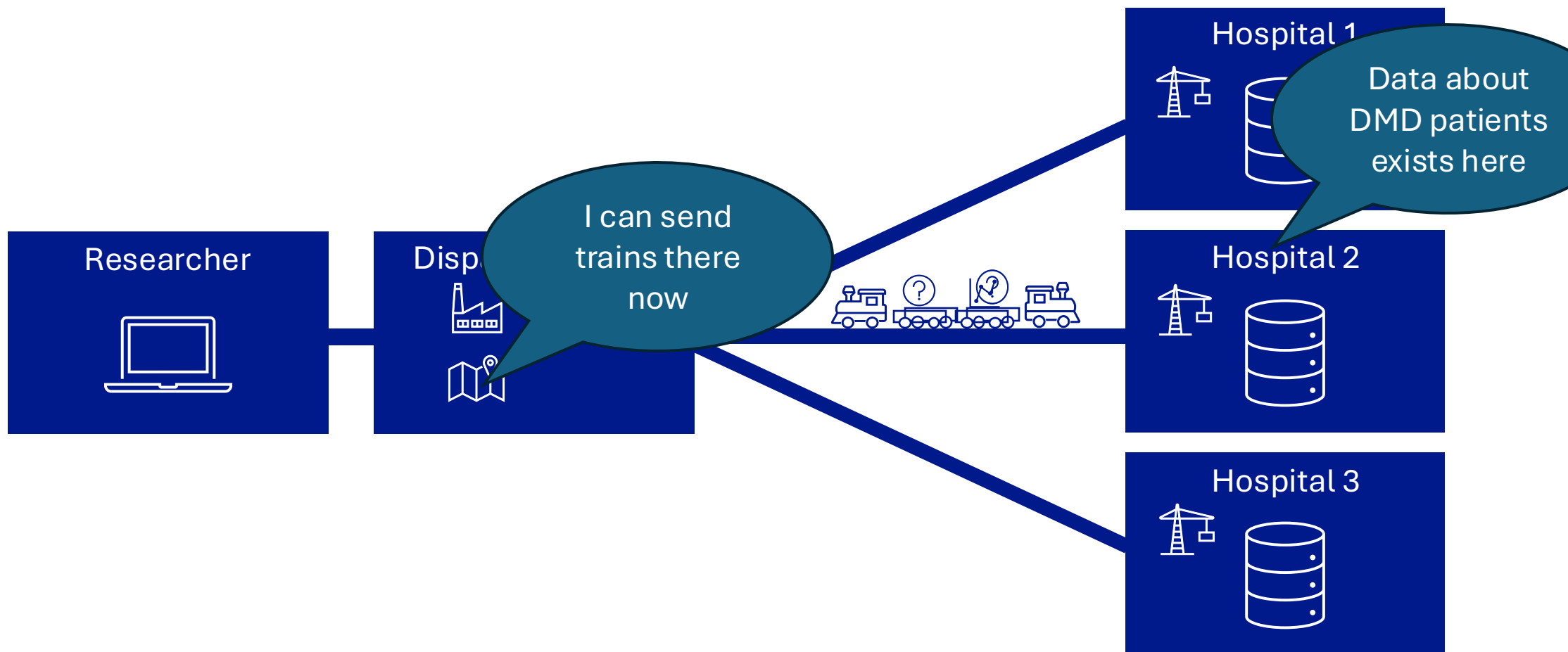
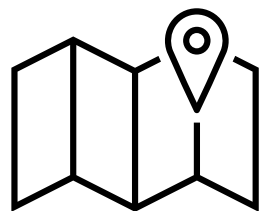


Railway map

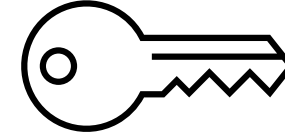


- The trains need to know which hospitals can be visited
 - Know that the data exists.
 - Know that it can ask a question about the data
- This is one part where FAIR is important, as FAIR Data Points can point towards the data and describe it with metadata

Railway map

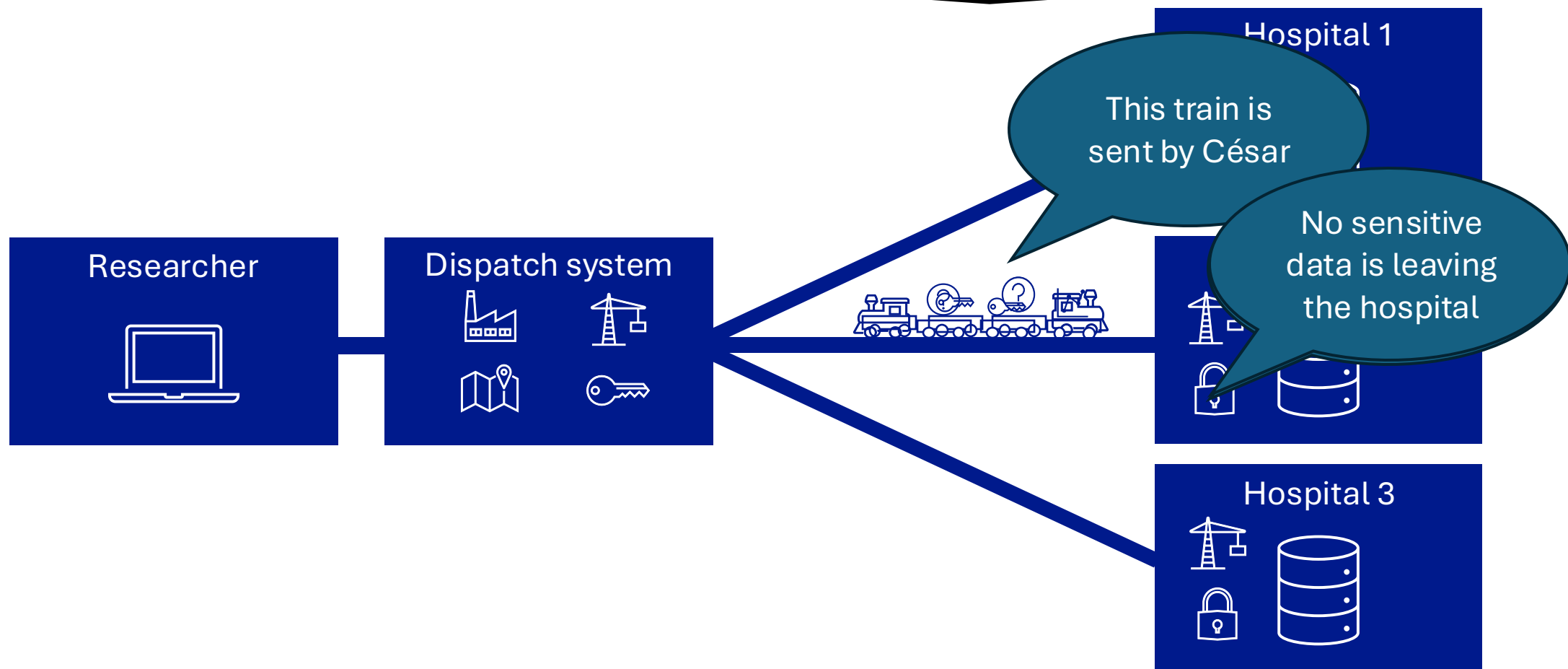
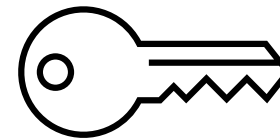


Authentication and validation



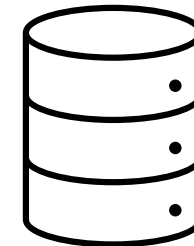
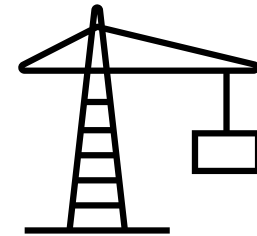
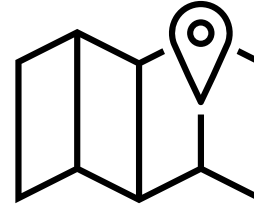
- Privacy is the main concern, which is why we do federated analysis
 - The hospital should verify the identity of a researcher
 - The hospital should verify that no sensitive data leaves the hospital

Authentication and validation



What can FAIR metadata do for federated analysis?

- Advertise the data
 - E.g. There is data for x here
 - E.g. This data has y access conditions
- Explain to the hospital how it can interact with trains
 - E.g. The train contains a SPARQL/FHIR/SQL query
 - E.g. The train needs x GB of memory to run
- Enable harmonization of data so that the same train works for multiple hospitals



Conclusion

- Several components are needed for federated analysis
- The FAIR principles can enable several of these components

