

Fast, just-in-time queries on heterogeneous raw data

Anastasia Ailamaki

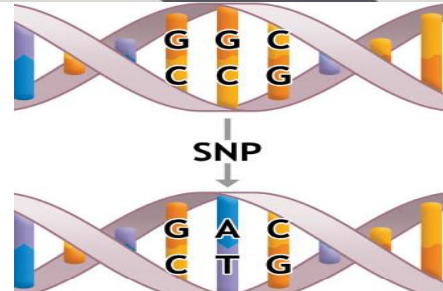
EPFL and RAW Labs SA

With **Manos Karpathiotakis, Stella Giannakopoulou, Matt Olma,**
and the **DIAS lab**

*most firms estimate that they are only analyzing
12% of the data that they already have*

Forrester, 2014

- growing data
- growing heterogeneity
- data movement restrictions



available data *impedes* business & scientific analytics

The Digital Universe: 50-fold Growth from the Beginning of 2010 to the End of 2020

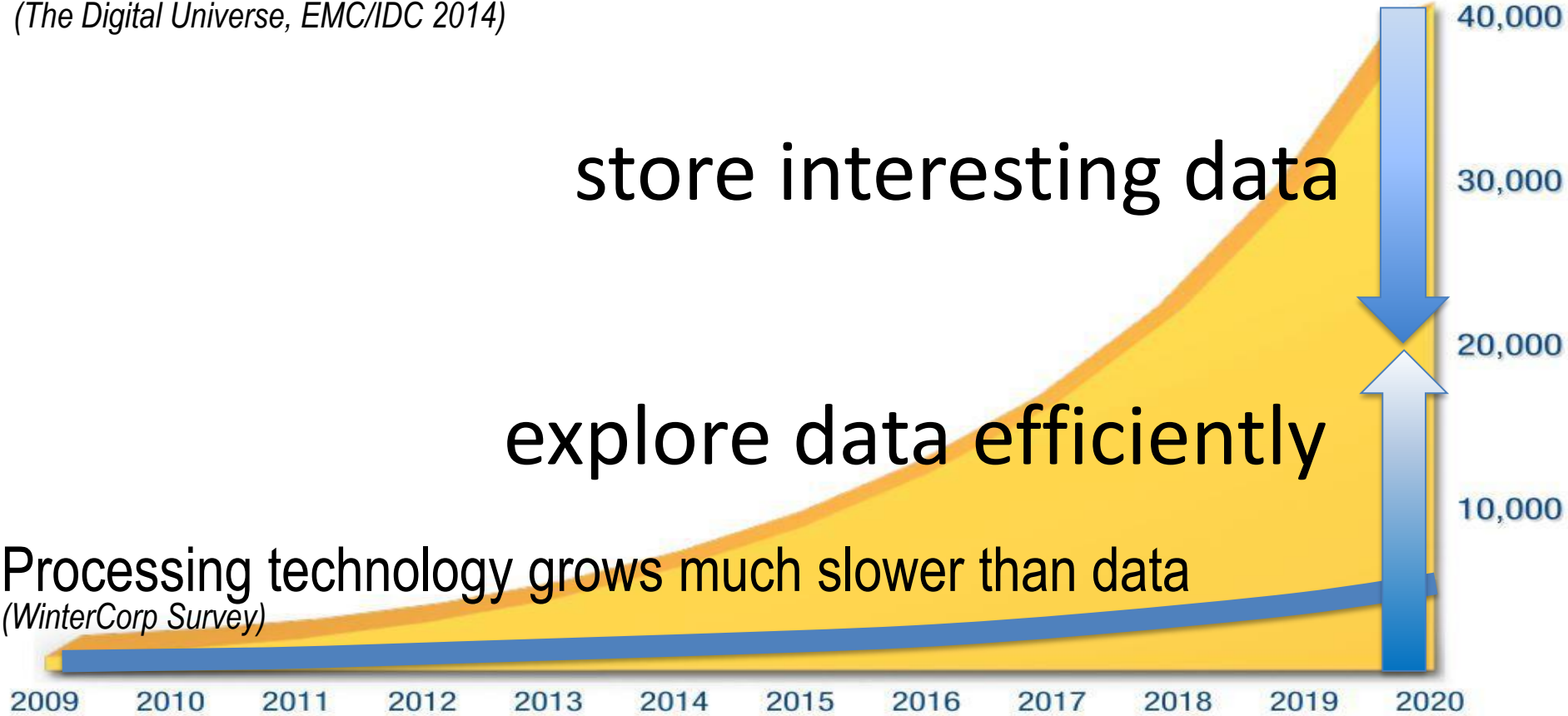
(The Digital Universe, EMC/IDC 2014)

store interesting data

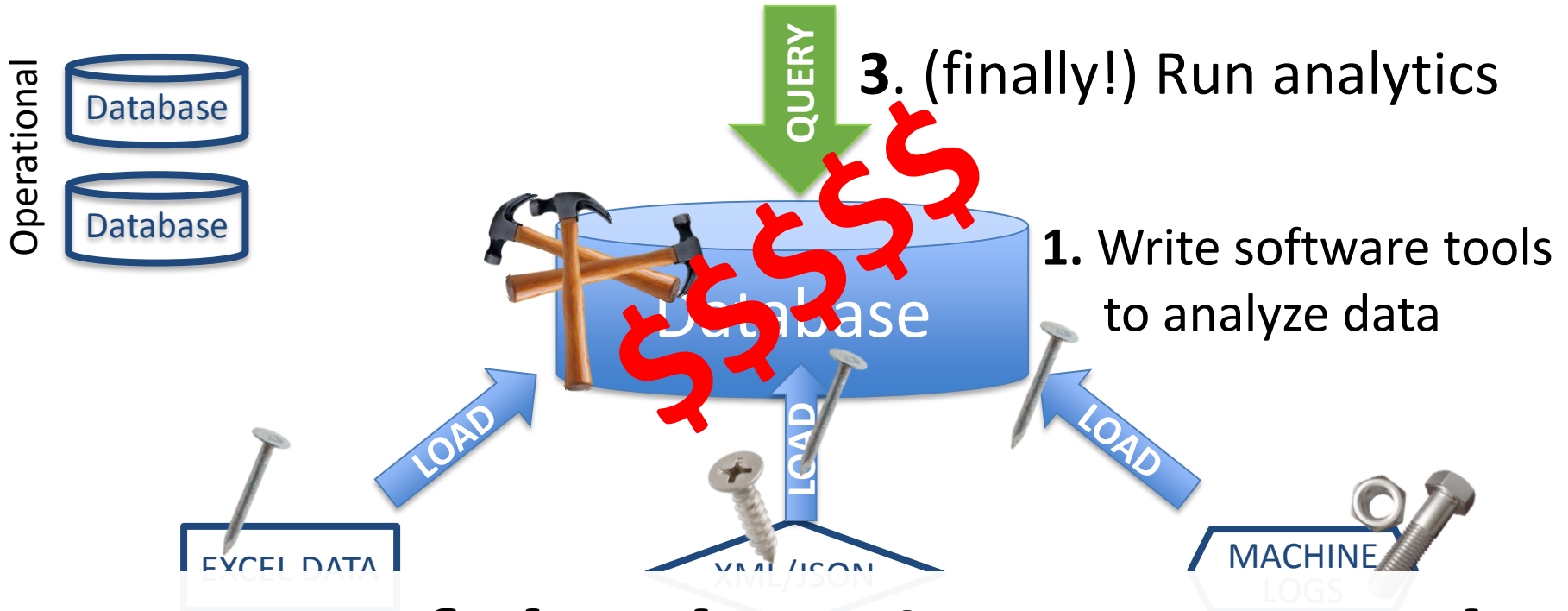
explore data efficiently

Processing technology grows much slower than data

(WinterCorp Survey)



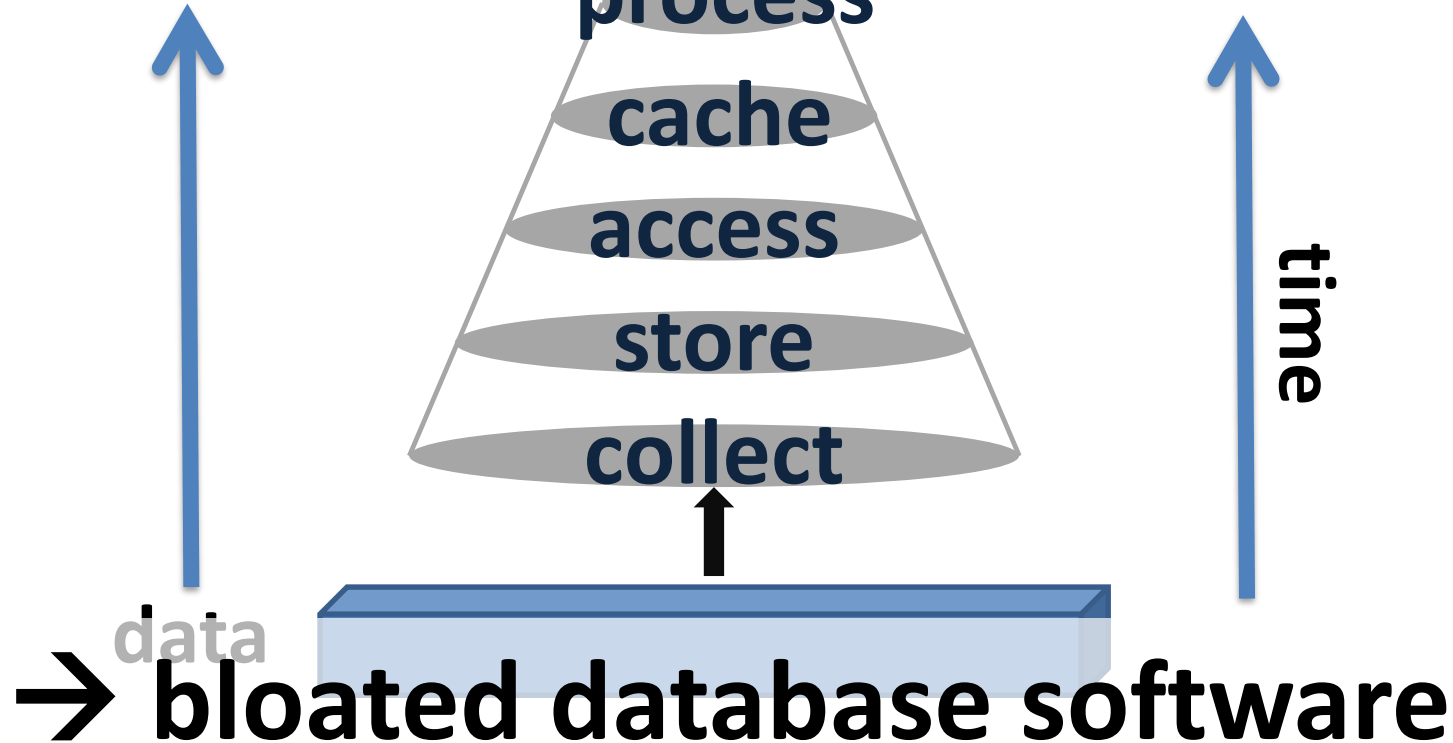
When you have a hammer...



90% of the data is never used.

build database to run queries

information



data

bloated database software

new: one DB per app/data pair

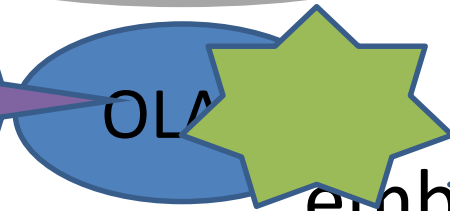
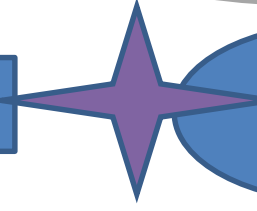
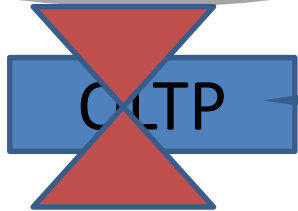
80% of analysts' time goes to data preparation and configuration

Main-memory
DBMS

Column
stores

NoSQL
systems

Stream
DBMS



embarrassingly parallel

databases will be extinct



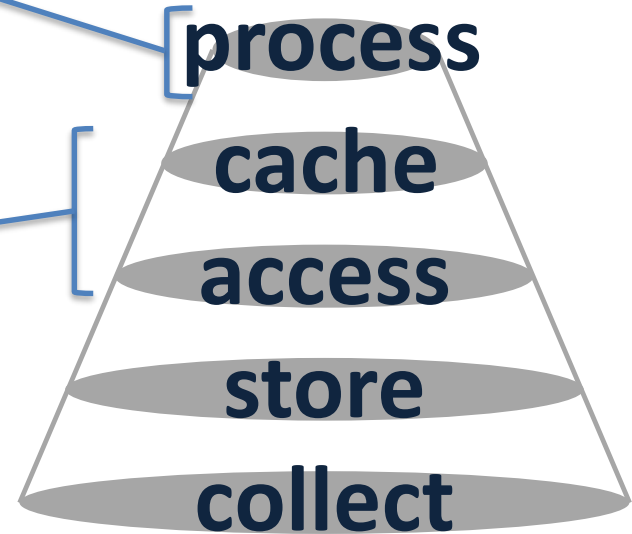
The way forward

- Data model:
 - Support variety (complex structured and unstructured data)
 - Col-store/Row-store are only two of many possible layouts
- Storage model:
 - Don't store!
 - Run in situ and cache based on actual needs/usage
- Execution model:
 - Generate engine based on query, available caches, history

Fundamentally rethink DB stack

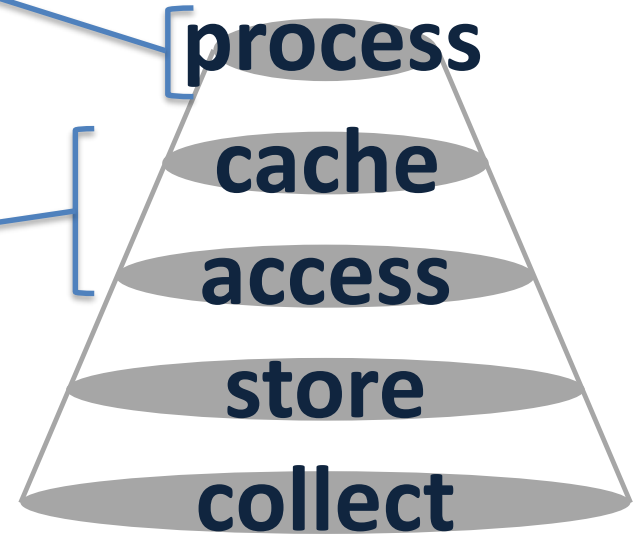
RAW — A lean and agile engine

- Adaptive Query Processing
 - A database per query and dataset
 - SQL++ to query and clean all data
- Adaptive data access
 - Tune database dynamically
- The Human Brain Project
 - An inspiring use case

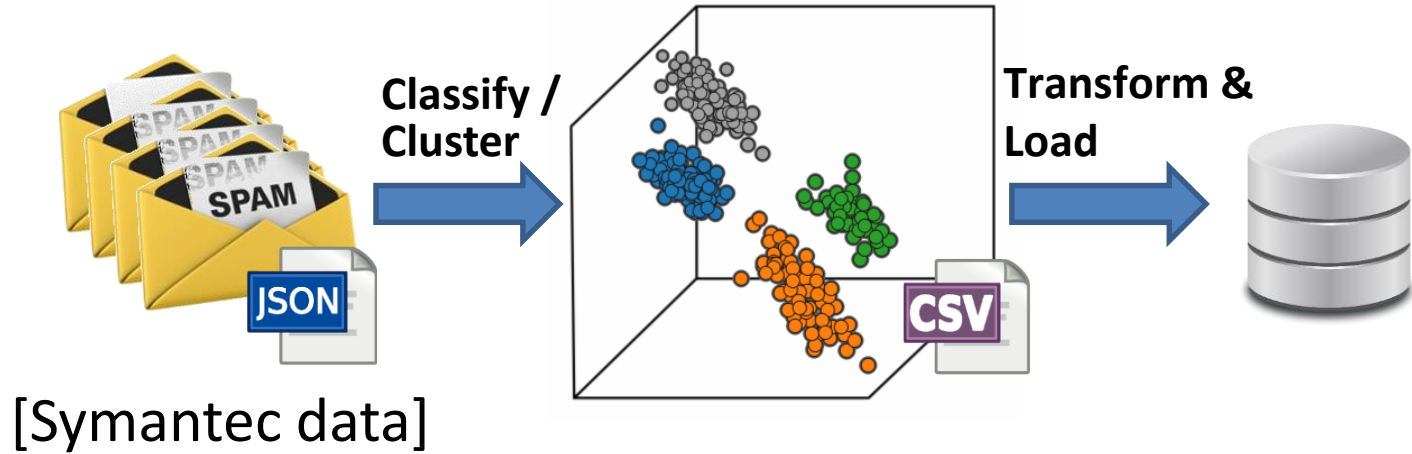


RAW — A lean and agile engine

- **Adaptive Query Processing**
 - A database per query and dataset
 - SQL++ to query and clean all data
- Adaptive data access
 - Tune database dynamically
- The Human Brain Project
 - An inspiring use case



Detecting active spambots



Flexibility

Ad-hoc queries
over diverse data formats



Performance

Fast queries
regardless of data format

fast queries on heterogeneous data

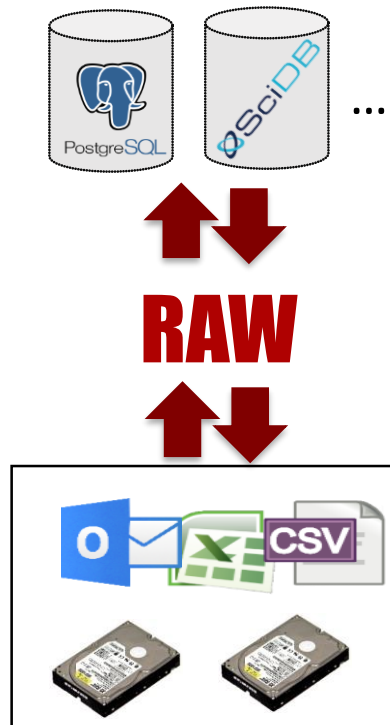
cannot load into a Database System!

- diverse formats
- legacy software
- privacy limitations
- data “owned” by one database

RAW: interface to raw data

With extended SQL

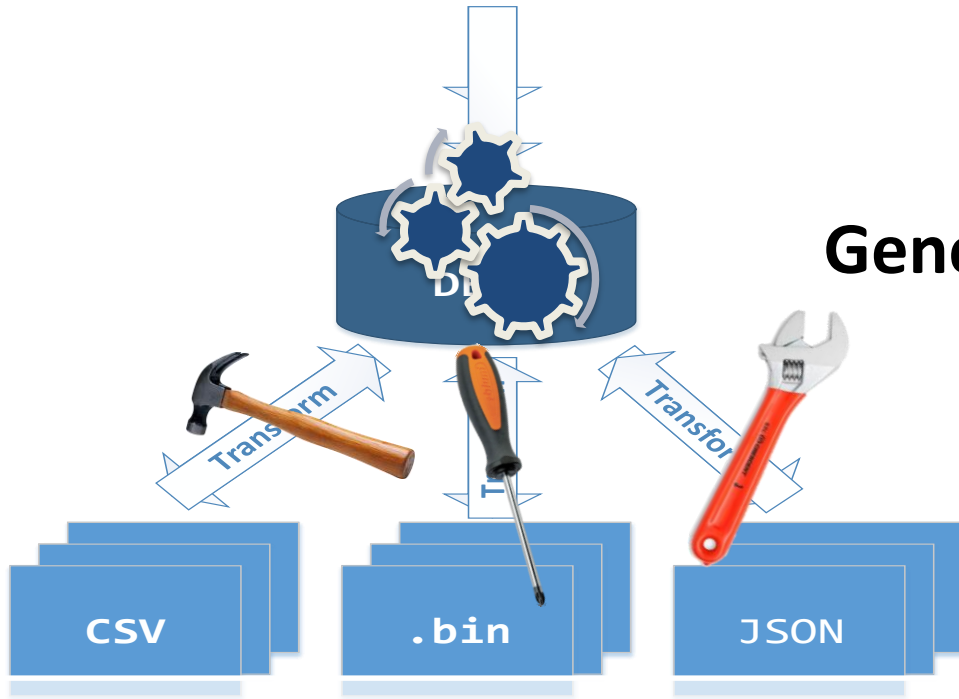
code-generated engine



key: data virtualization

Adapting a query engine to data

Query



Generate plug-in per data source



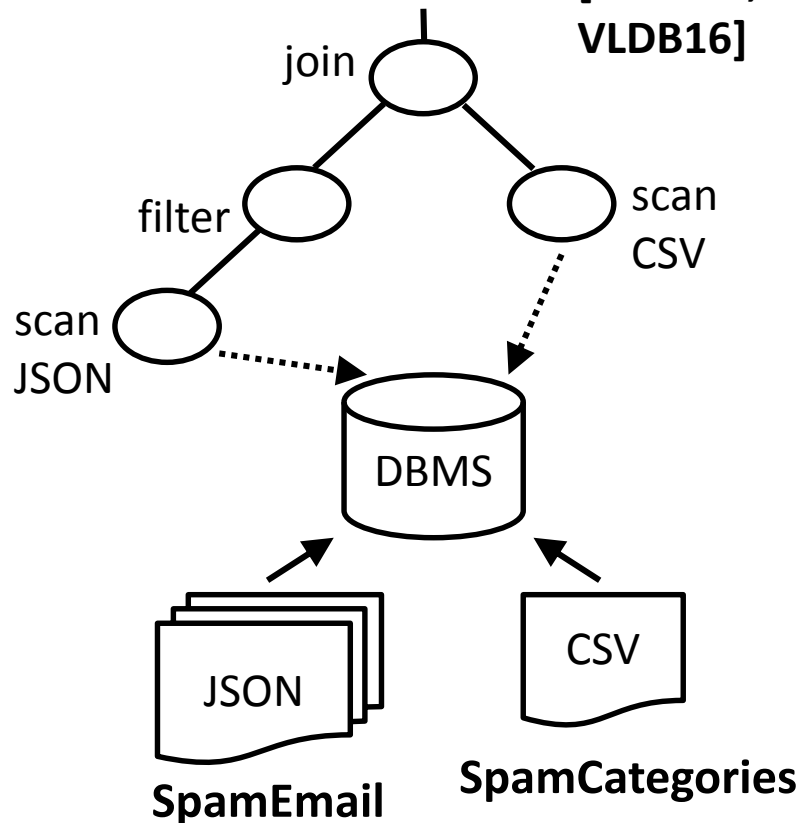
Treat each source as
native storage format

Query original data formats, files, and scripts

How to build a just-in-time data base

```
SELECT bot, country, ...  
FROM SpamEmail e, SpamCategories c  
WHERE e.id == c.id AND  
      e.lang = 'English' AND ...
```

[VLDB14,CIDR15,
VLDB16]



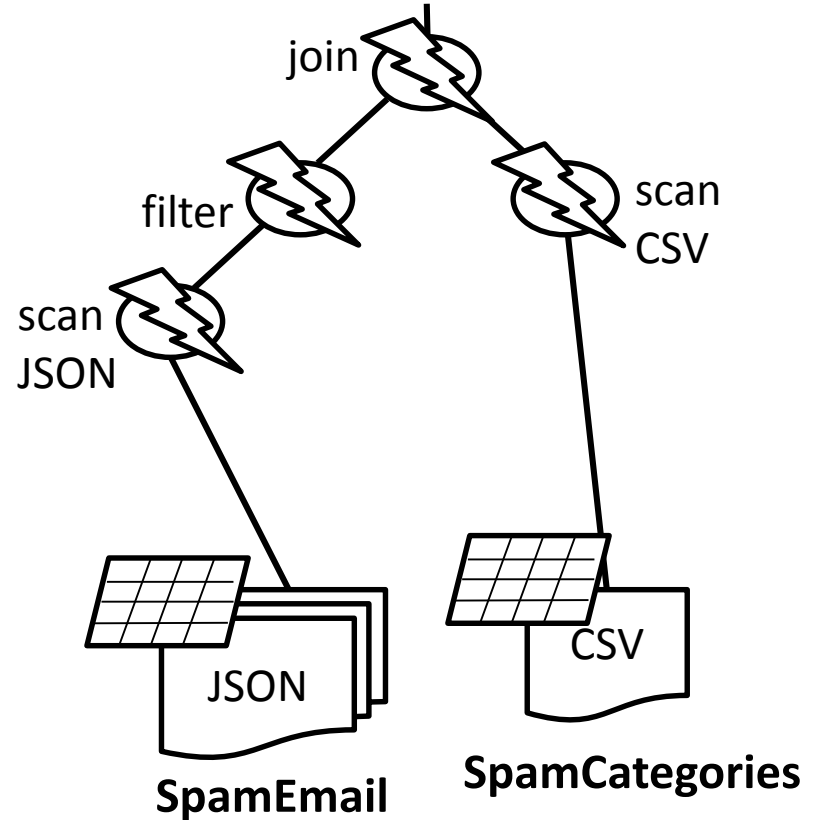
How to build a just-in-time data base

```
SELECT bot, country, ...  
FROM SpamEmail e, SpamCategories c  
WHERE e.id == c.id AND  
      e.lang = 'English' AND ...
```

 Code Generate the Access Paths

 Code Generate the Query

 Build Position and Data Caches



Queries → Monoid comprehensions

Monoids:

- Abstraction for “aggregates” computation

*Fegaras
[SIGMOD95,
TODS 2000,...]

Monoid Comprehensions*:

- Operations between monoids

```
for {  
  p <- Patients, r <- BrainRegions,  
  p.id = r.id, r.amygdala.Vol > 0.2  
} yield bag p.age
```

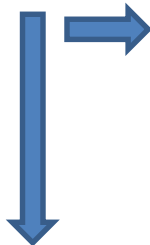
Sum/Bag/List/Set/Top-K/...

Support multiple data models as input & output

“SQL++” → Comprehensions → Algebra

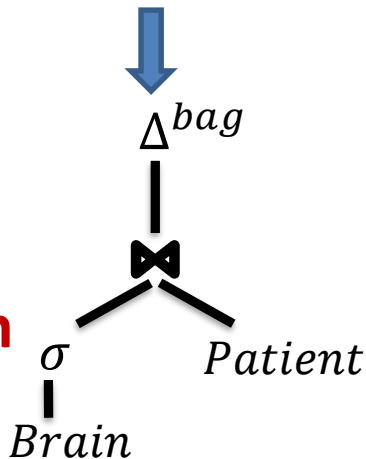
```
SELECT r.age
FROM Patients p
JOIN BrainRegions r
ON (p.id = r.id)
WHERE r.amygdala.Vol > 0.2
```

Internal Calculus



if-else
record construction
function application
(nested) comprehension
...

```
for {
  p <- Patients,
  r <- BrainRegions,
  p.id = r.id,
  r.amygdala.Vol > 0.2
} yield bag r.age
```



Algebraic form amenable to relational optimizations

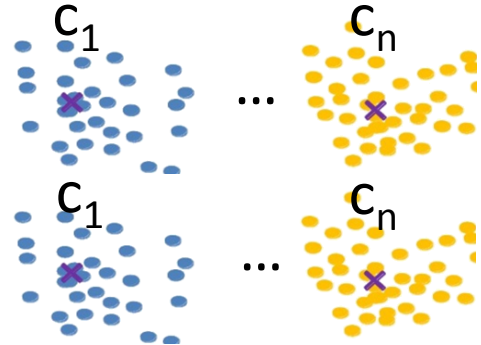
Data cleaning using monoid comprehensions

```
for(o←orders) yield list split(o.ship_date,"/")
```

```
dataGroup := for (o←orders)  
  yield cluster(o.item,kmeans)
```

```
dictGroup := for (d←dict)  
  yield cluster(d.item,kmeans)
```

```
for(d1←dataGroup,  
  d2←dictGroup,  
  d1.center = d2.center,  
  similar(metric,d1.item,d2.item,θ))  
yield group (d1.item)
```



Optimize cleaning operations holistically

SQL-like extensions for data cleaning

Functional Dependencies:

orderno, item \rightarrow quantity

```
SELECT o.orderno, o.item, *
```

```
FROM Orders o
```

```
FD((o.orderno, o.item), o.quantity)
```

Data Deduplication:

```
SELECT <projections>
```

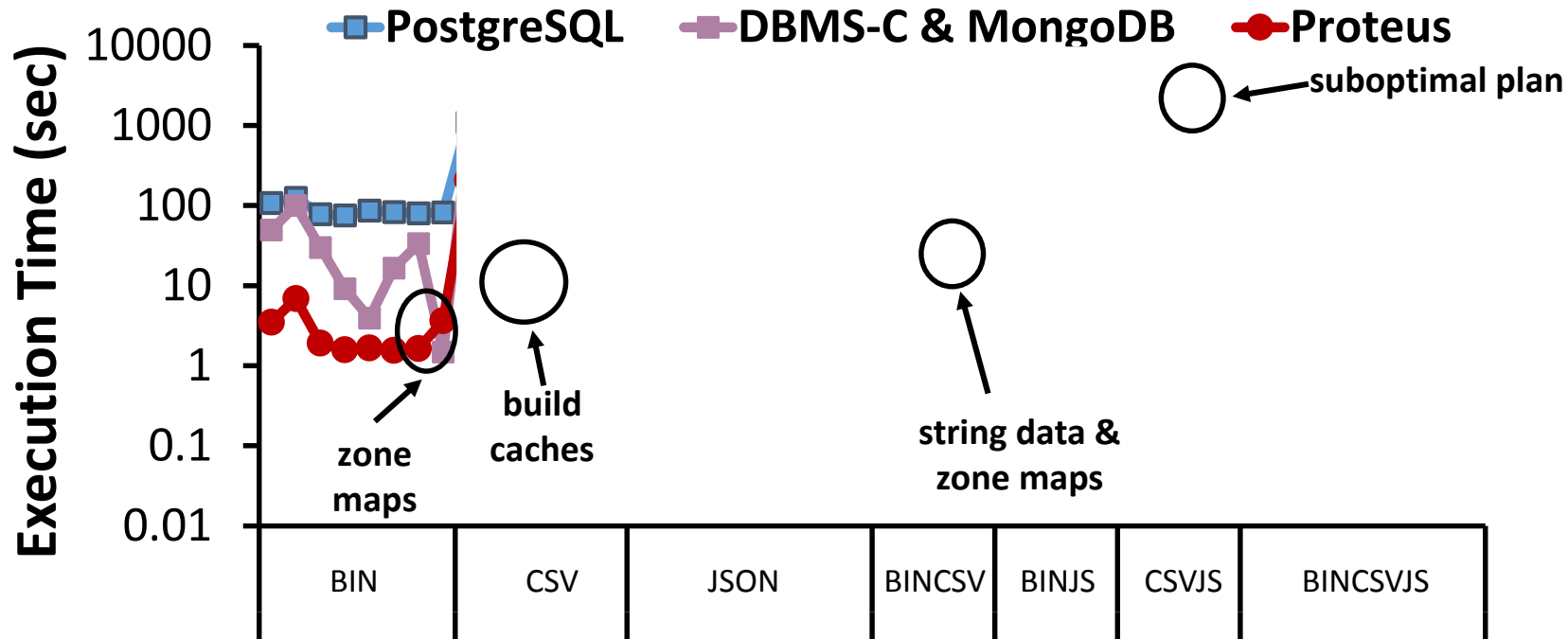
```
FROM <dataset>
```

```
DEDUP([<metric>,) [<theta>,) <attributes>)
```

Mask complex comprehension syntax

Symantec Spam Email Analysis

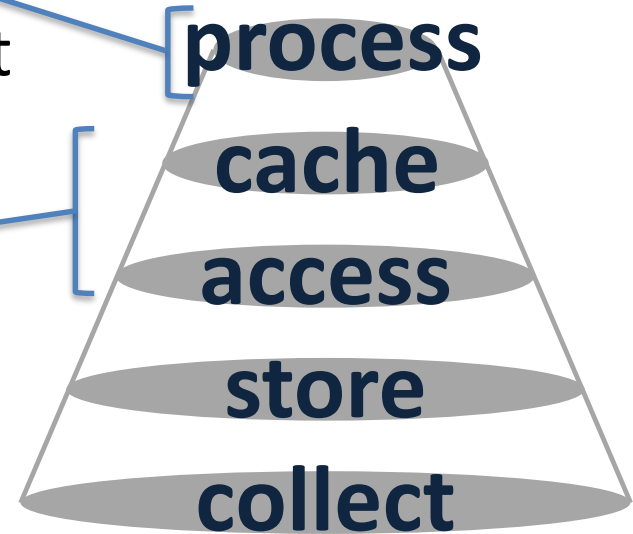
95GB Binary - 22GB JSON - 22GB CSV
50 queries



Flexible and fast by specializing & adapting

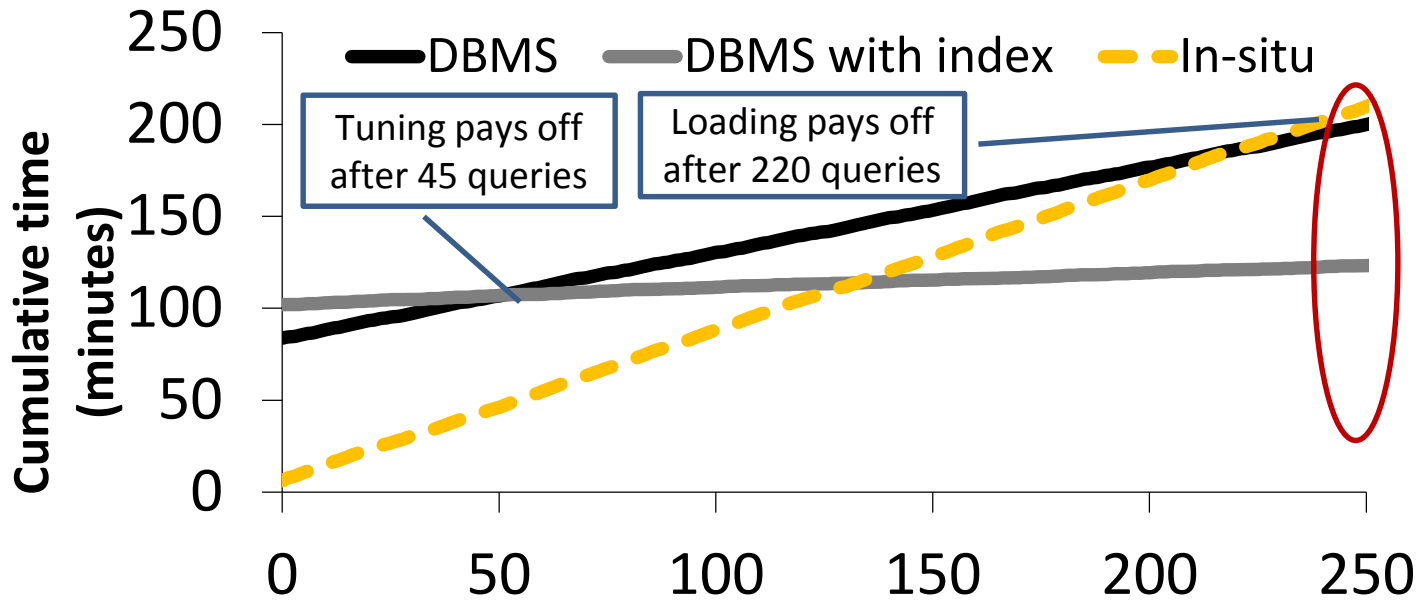
RAW — A lean and agile engine

- Adaptive Query Processing
 - A database per query and dataset
 - SQL++ to query and clean all data
- **Adaptive data access**
 - **Tune database dynamically**
- The Human Brain Project
 - An inspiring use case



Querying raw data w/o index: Diminishing returns

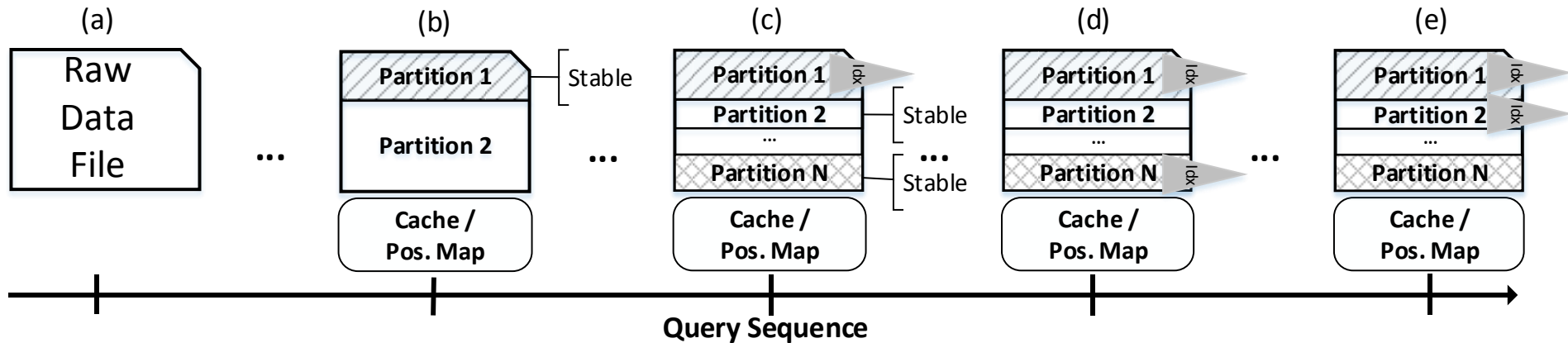
60GB smart meter data,
selectivity 1%,
128GB RAM, 1 thread



Scanning all data is slow

Indexing/tuning non-trivial for ad hoc data

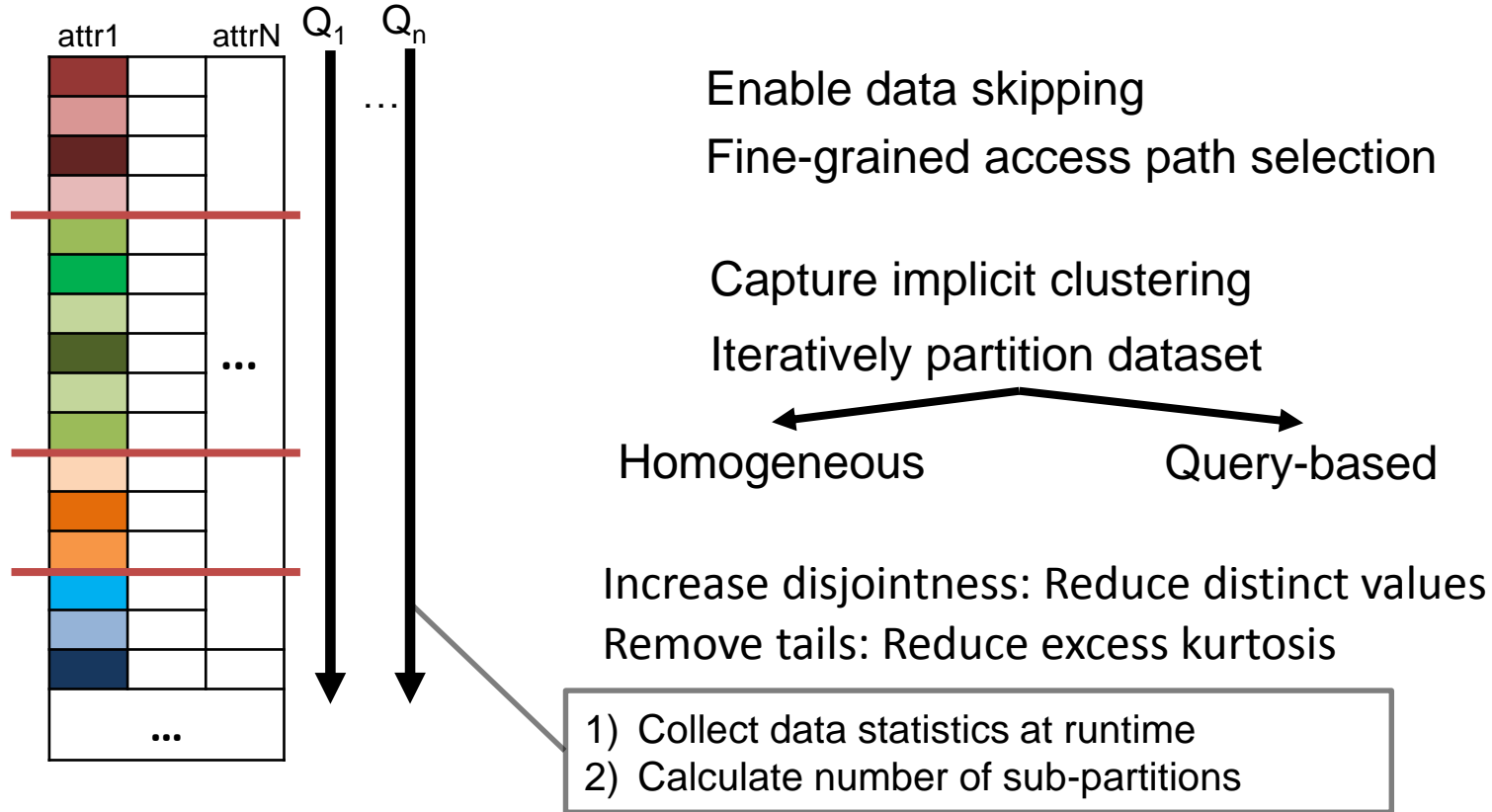
Invest in popular data subsets



Refine partitions over the data => Skip if useless for query

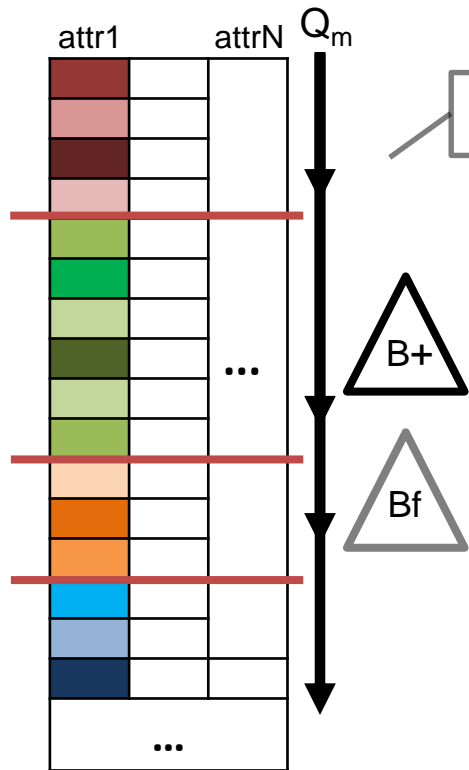
Tune indexes over popular partitions => Minimize data accesses

Adapt to data: Logical partitioning



Set the “ground” for reducing data access

Adapt to queries: index tuning



Index tuning on partition level

costs vs. gains
Should I build or not?

Choose what & when to build

What

- Value-Existence (i.e., Bloom filters)
- Value-Position (i.e., B+ Trees)

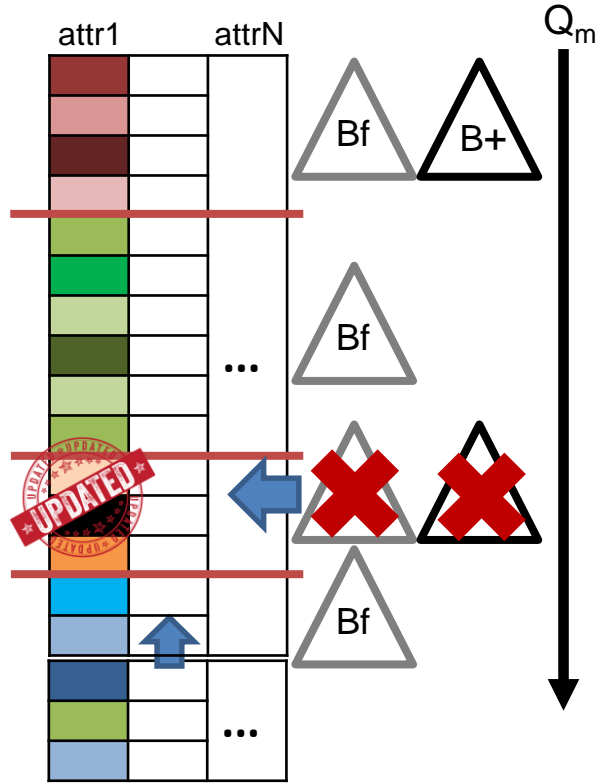
When

- Based on randomized algorithm
- Cost of scan vs. cost of build + gain

Build and drop based on budget

Maximize gain: build cost vs performance

Append & in-place updates



Store partition state

- Calculate hash value (MD5)

Monitor file for modifications

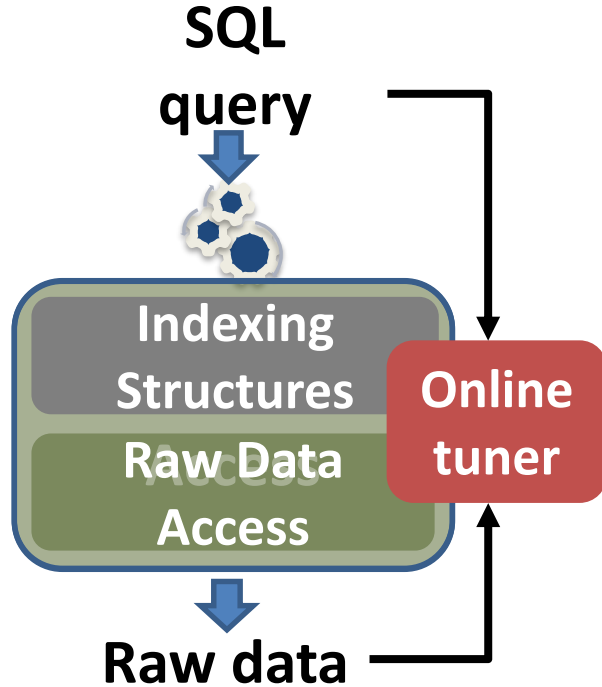
Recognize updated partitions

Fix modified partitions

- Drop/Re-build cache/index

Minimize update overhead

Slalom: Adaptive indexing over raw data



Incremental logical partitioning

- Based on data distribution

Adaptive partition indexing

- Based on access patterns

Monitors data for updates

- Updates data structures

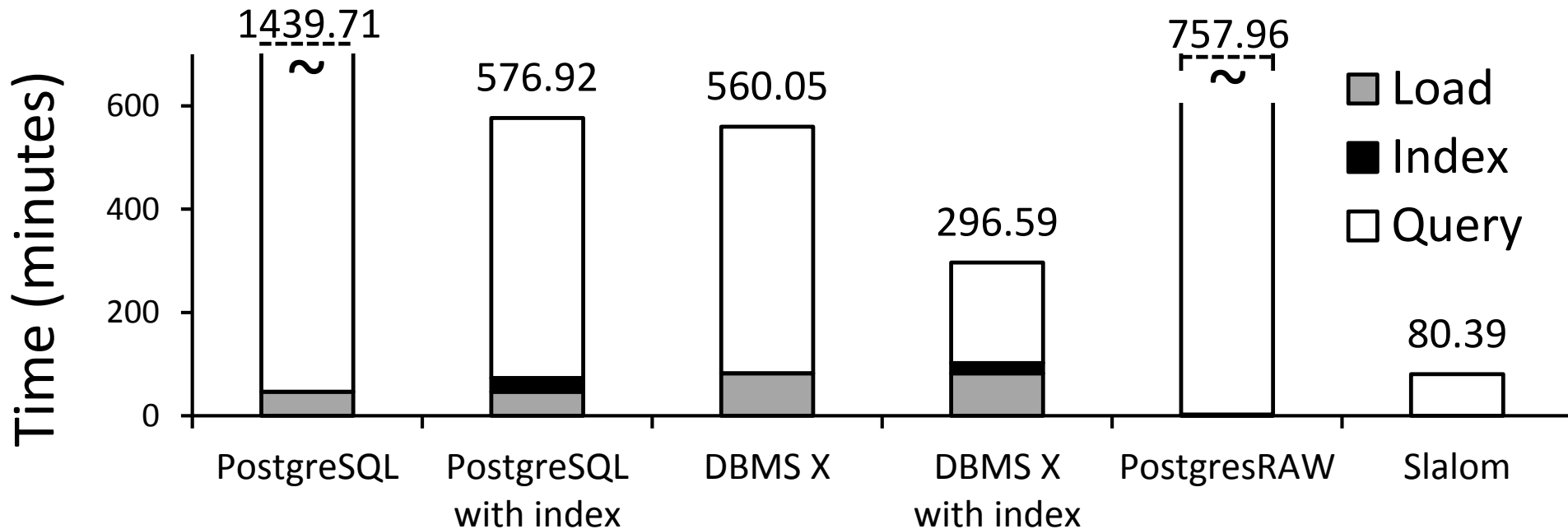
Combining Online Tuning with Adaptive Indexing

Adapt data access to queries and data at runtime

From raw data to results

59GB uniform dataset, 128GB RAM, cold caches

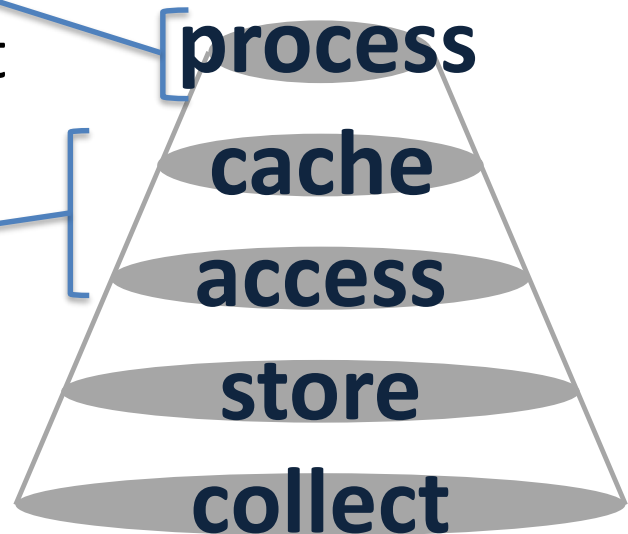
1000 point & range queries interchange on 2 attributes, sel: 0.5%-5%



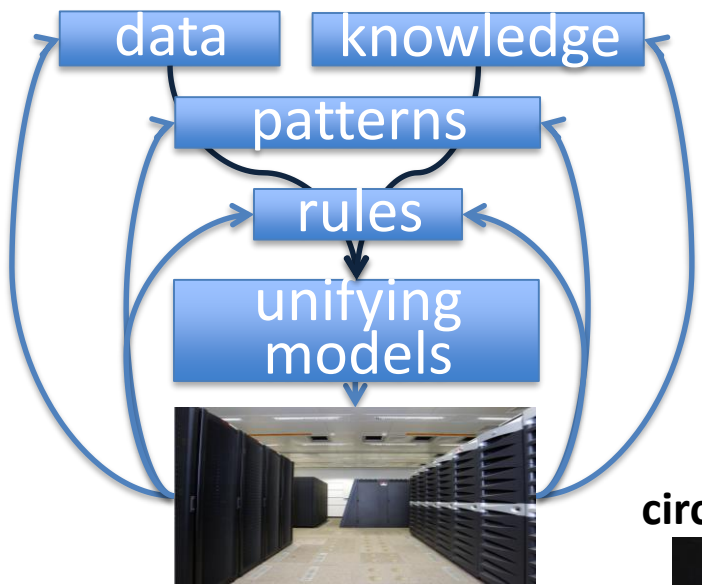
In-situ adaptive indexing achieves interactive access

RAW — A lean and agile engine

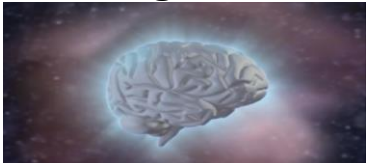
- Adaptive Query Processing
 - A database per query and dataset
 - SQL++ to query and clean all data
- Adaptive data access
 - Tune database dynamically
- **The Human Brain Project**
 - An inspiring use case



human brain project



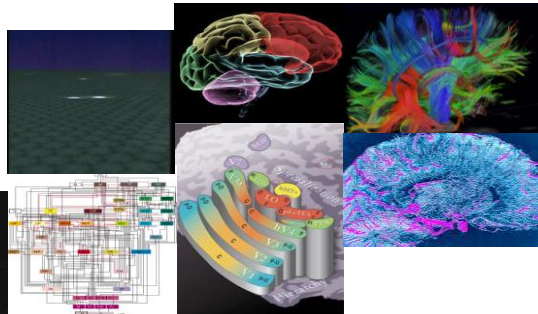
cognition



whole brain

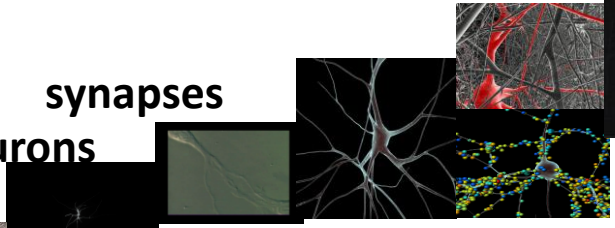


circuits



synapses

neurons



integrate clinical and simulation data

patient data: medical informatics

1 DATA FEDERATION



health services

epidemiology data visualization

hypothesis testing

clinical trials

5 BIOLOGICAL SIGNATURE OF DISEASE

→ DISEASE DEFINITION

→ PHARMACOLOGY

→ CLINICAL TRIAL

2 DATA INTEGRATION

3 DATA MINING

biological disease signatures

the coupling of
clinical measurements with
validated biomarkers

Example: Alzheimer's disease

Clinical - Phenotype	Proteomic Biomarkers	Genomic Biomarkers	Imaging Biomarkers
Cognition: memory	CSF protein: beta amyloid	APP, PSEN1, PSEN2	Volumetric change: hippocampus, inferior temporal cortex...
Functional capacity		Common genetic variant: APOE e4e4 ...	Beta amyloid imaging
General physical health			

memory loss/tau level/APOE e4e4/hippocampus atrophy/amyloid

medical informatics platform

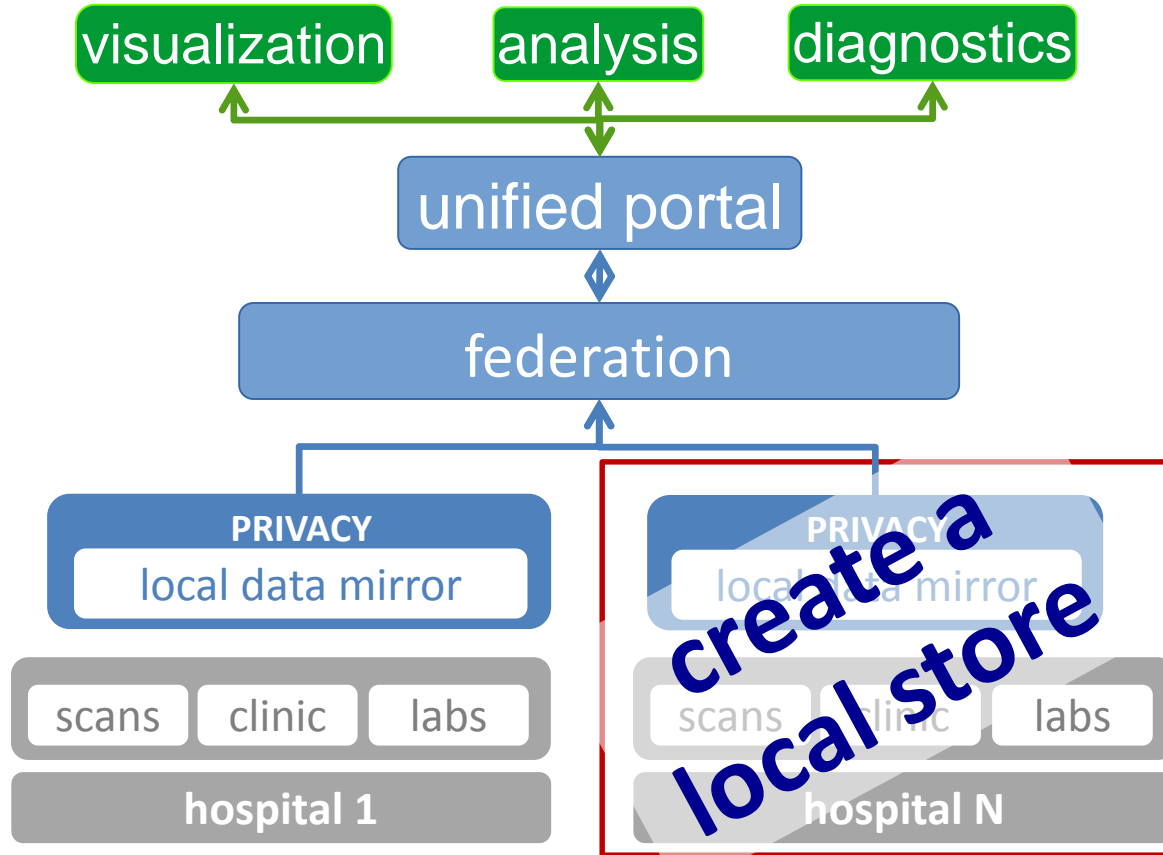
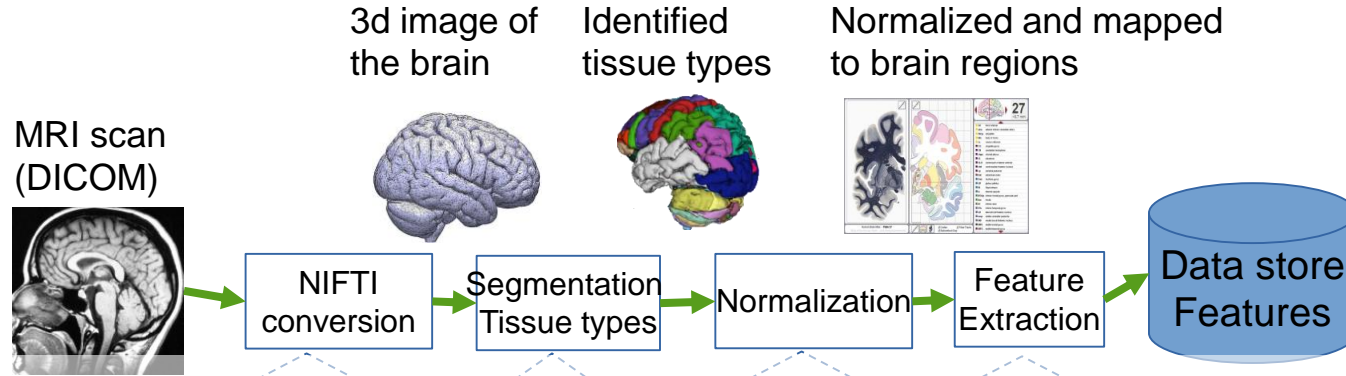
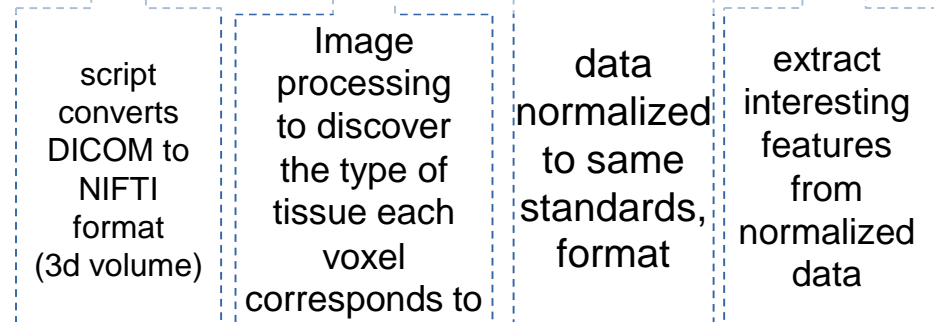


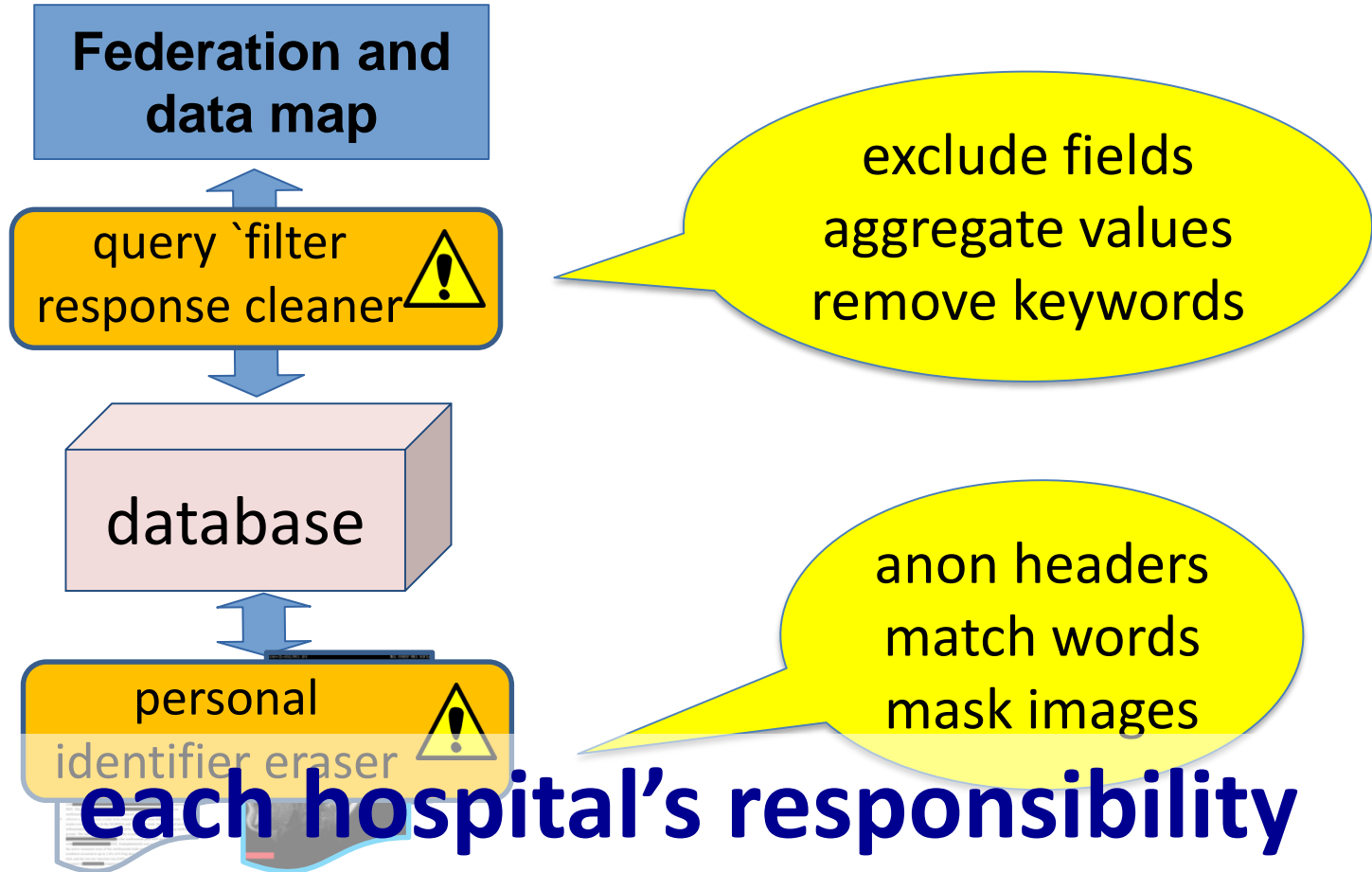
image preprocessing (Lausanne)



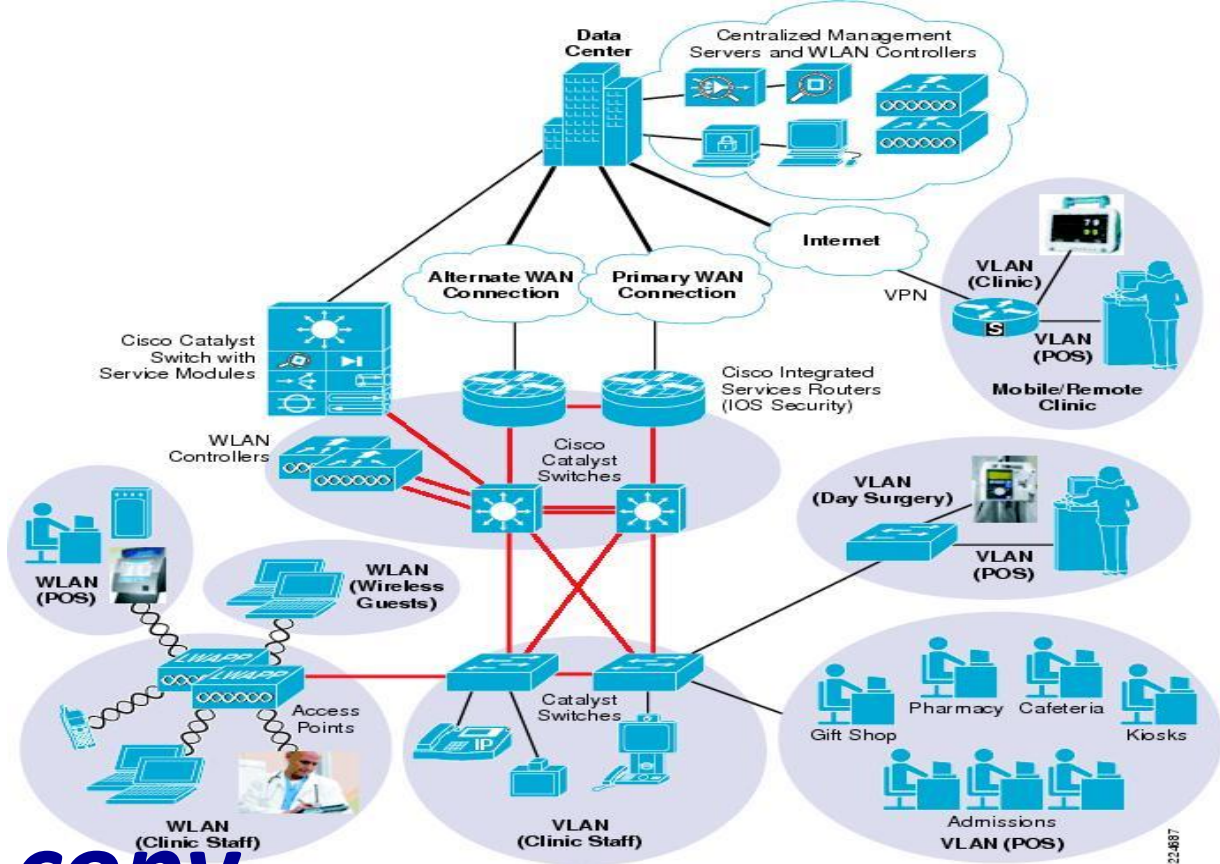
intermediate data also available



two-level anonymization



access to private hospital data



no move, no copy

clinical+genetic+imaging data → signature

Patients (CSV)

id	Protein: AACT	Age	Phenotype	...
1	1.4	45	Trauma	...
2	2	55	Chronic Symptoms	...
3	0.2	56

Brain_GrayMatter (Binary)

	0	1	...	n
0	0.45	0.75	...	0.1
1	0.33	0.3	...	0.38
...
m	0.12	0	...	0.47

signature:

age > 50

AND

amygdala.Vol > 0.3

AND

AACT < 1

BrainRegions (JSON)

```
[{"id": 1,
  "amygdala": {"X":15,"Y":20, "Vol": 0.5},
  "hippocampus": {"X":17, "Y":10, "Vol":0.2}},
{"id": 2, ...},
{"id": 3, ...}]
```

How **RAW** works

1. Ask a question



2. Generate the needed software tools



3. Discover interesting data



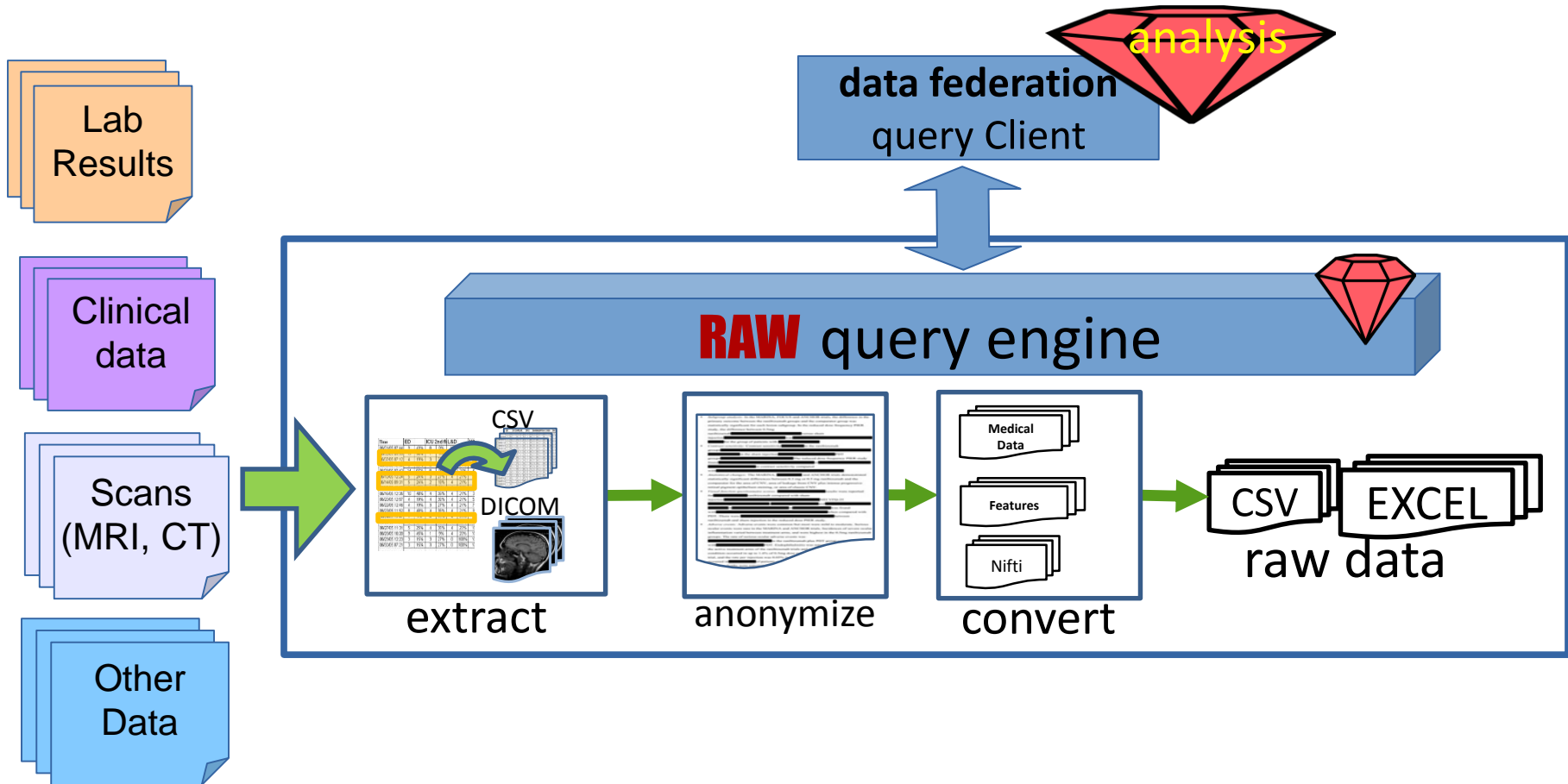
Data is accessed and integrated in real time

Why **RAW** is fast



As queries run, RAW remembers information on data accessed and generated code. Its “database” is only the useful data.

deployment: hospital data mirror



Medical Informatics Platform

visualization

analysis

diagnostics



unified portal

Niguarda (IT)

User Interface

Federation and
Data Map

RAW Query Engine

Hospital Data

UniClinic (DE)

User Interface

Federation and
Data Map

RAW Query Engine

Hospital Data

CHUV (CH)

User Interface

Federation and
Data Map

RAW Query Engine

Hospital Data

What we learned

- currently data management cost **grows with data owned**
- **impossible** to pre-cook a database system suitable for all data
- from *manual ingestion* to **automatic adaptation**: rethinking DB stack with just-in-time queries and storage

RAW

Just ask.

dias.epfl.ch

raw-labs.com