

OVERVIEW OF DATA EXPLORATION TECHNIQUES

Stratos Idreos, Olga Papaemmanouil, Surajit Chaudhuri
SIGMOD 2015, Melbourne

USER INTERACTION

express
interests

interests

collaborate

collaborate

visualize
results

results

query/results
recommendations

recommendations

annotate

annotate

assisted query
formulation

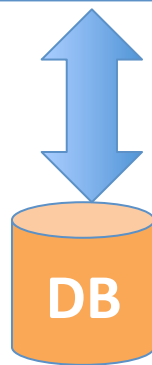
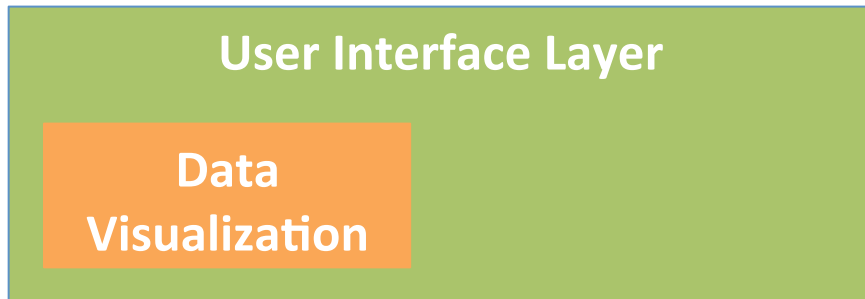
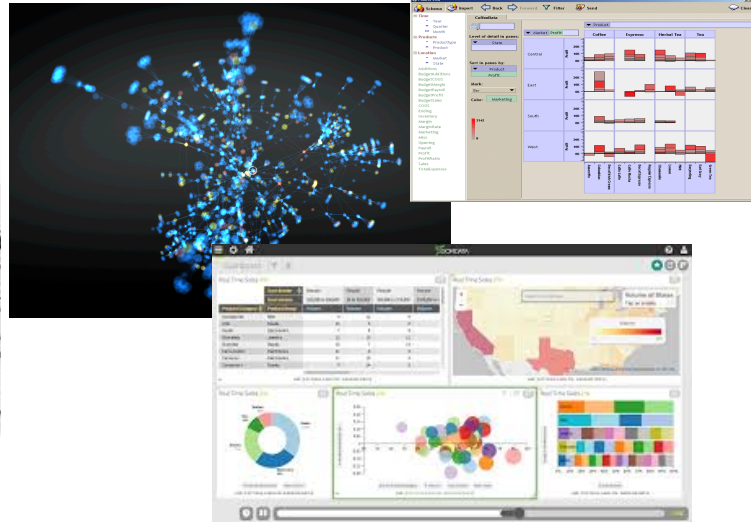
formulation

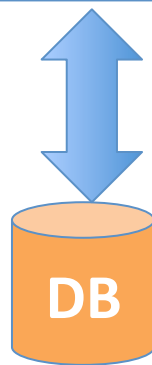
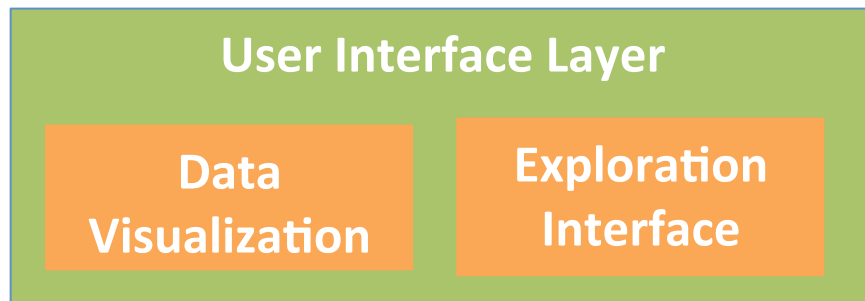
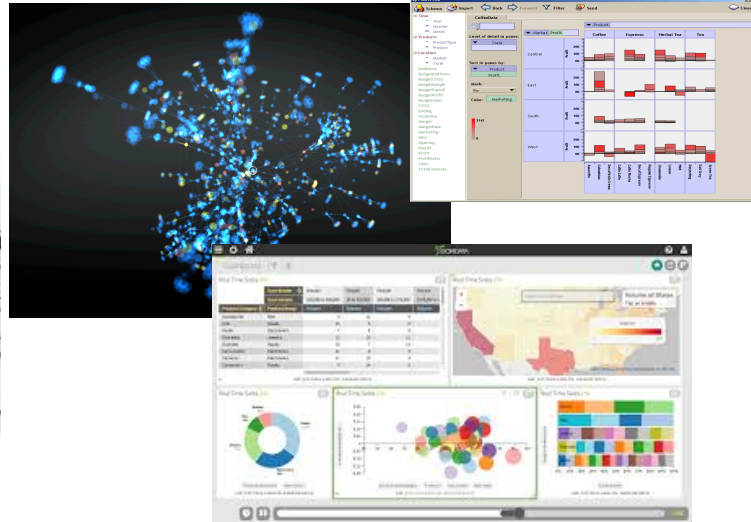


User Interface Layer



DB





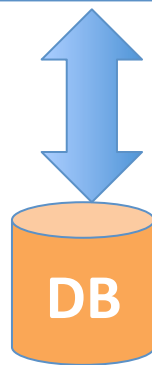
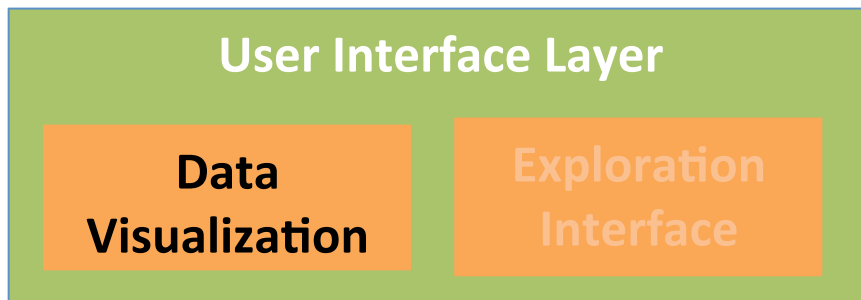


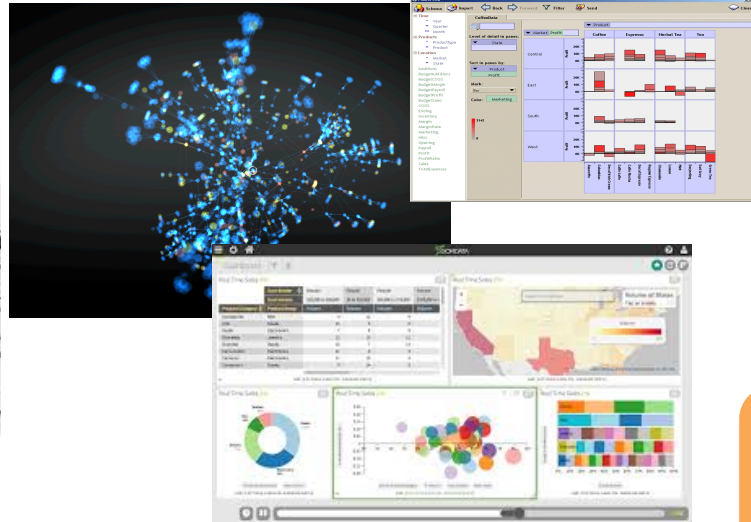
data visualization

visualization tools

visual optimizations

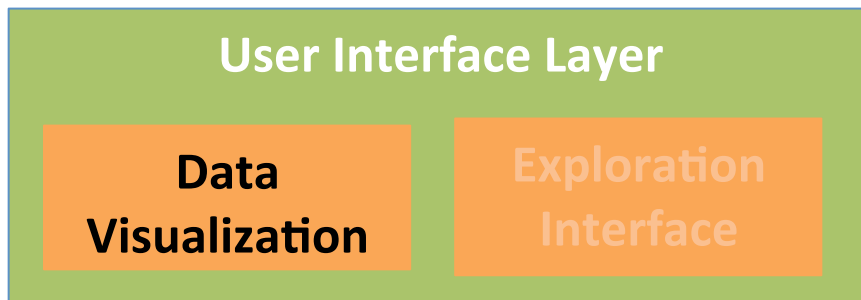
automatic visualization



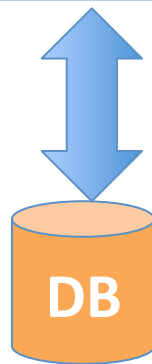


data visualization

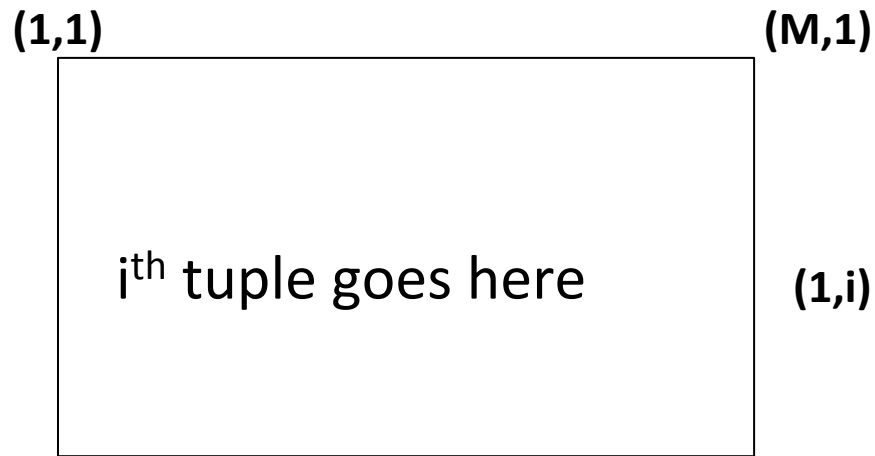
visualization tools



visual optimizations

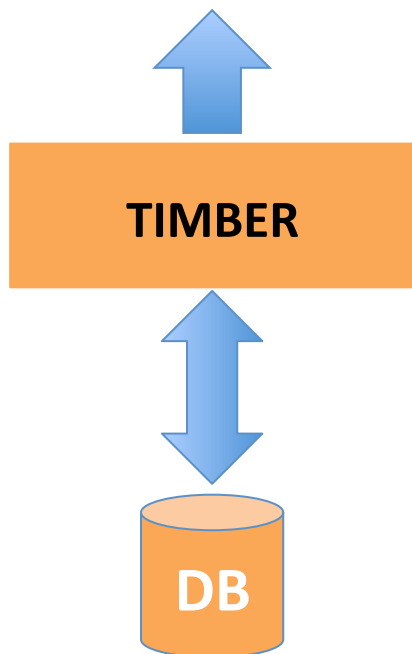


automatic visualization



Back in 1982...

window-based “sophisticated”
browser for relational DBs



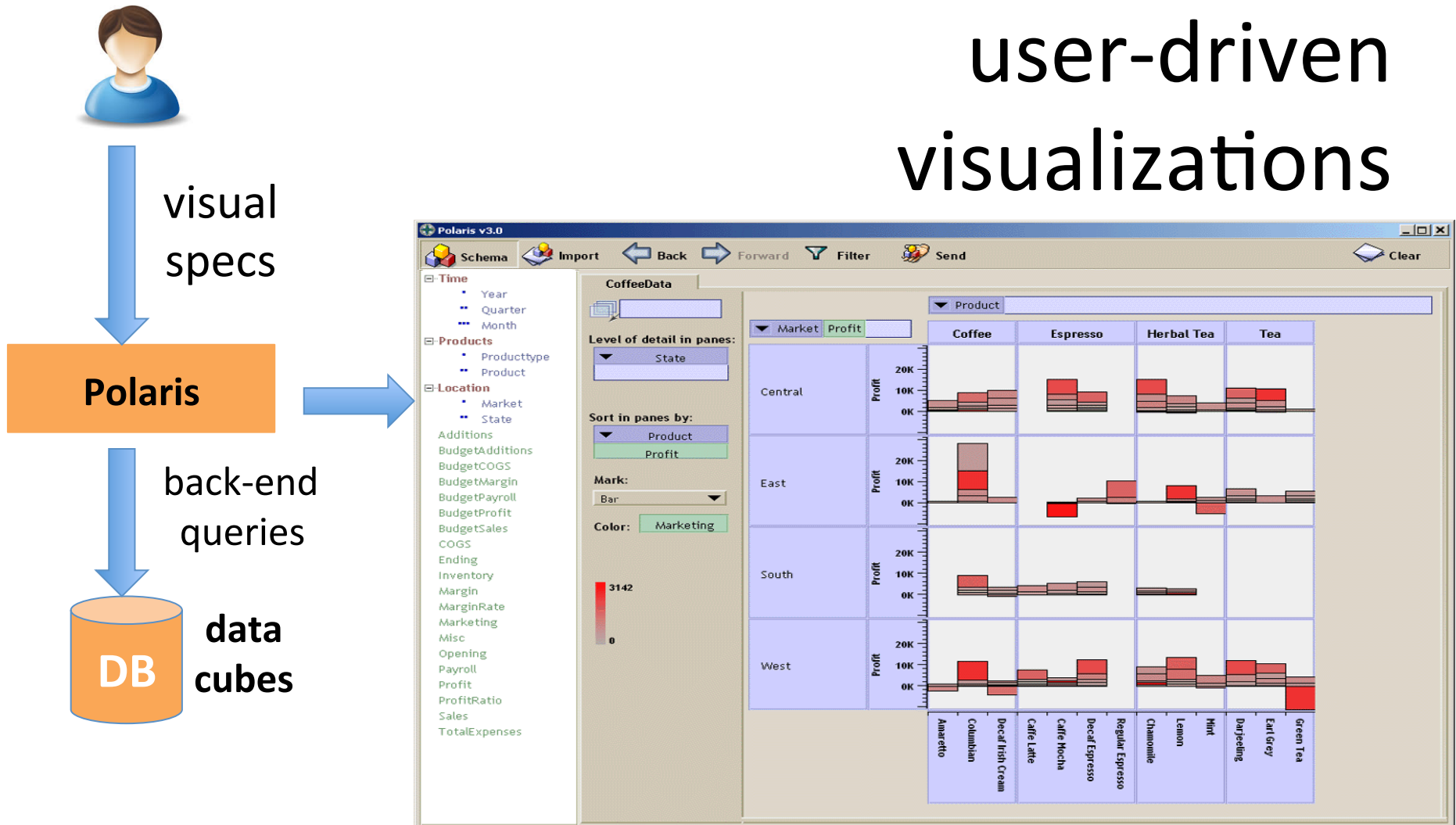
browser for multiple relations/tuples

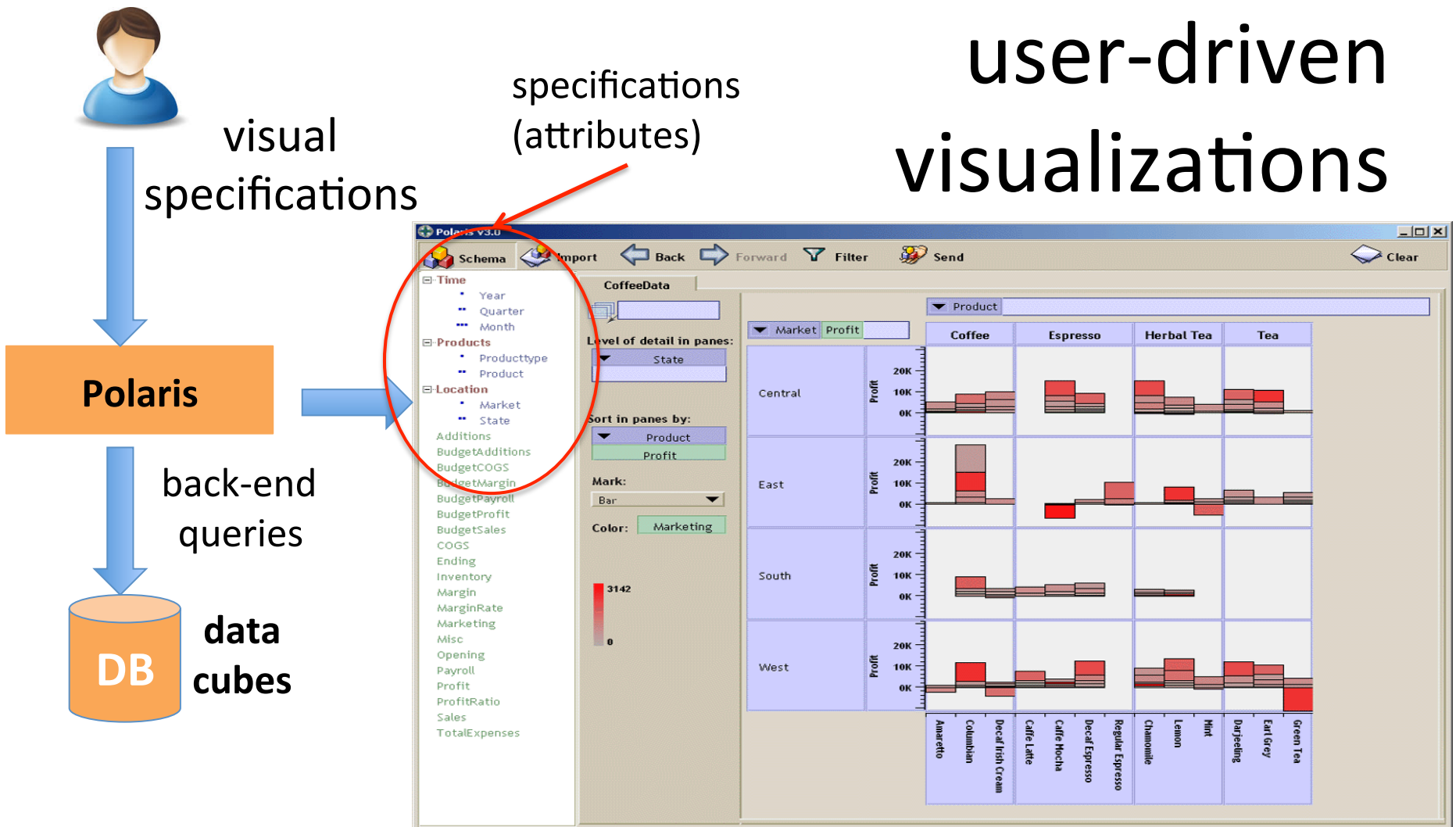
rich query language for icon-oriented DBs

visual editor of text objects

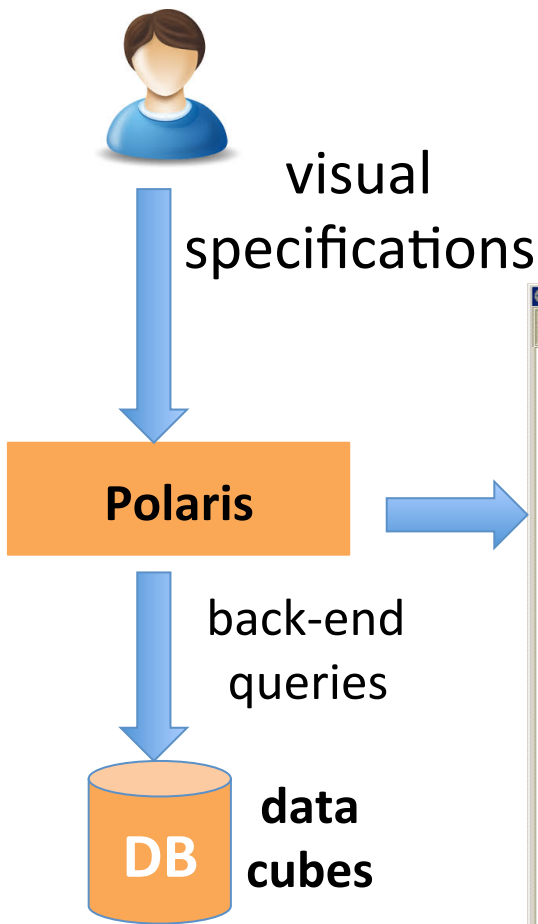
browser for geographical data

user-driven visualizations



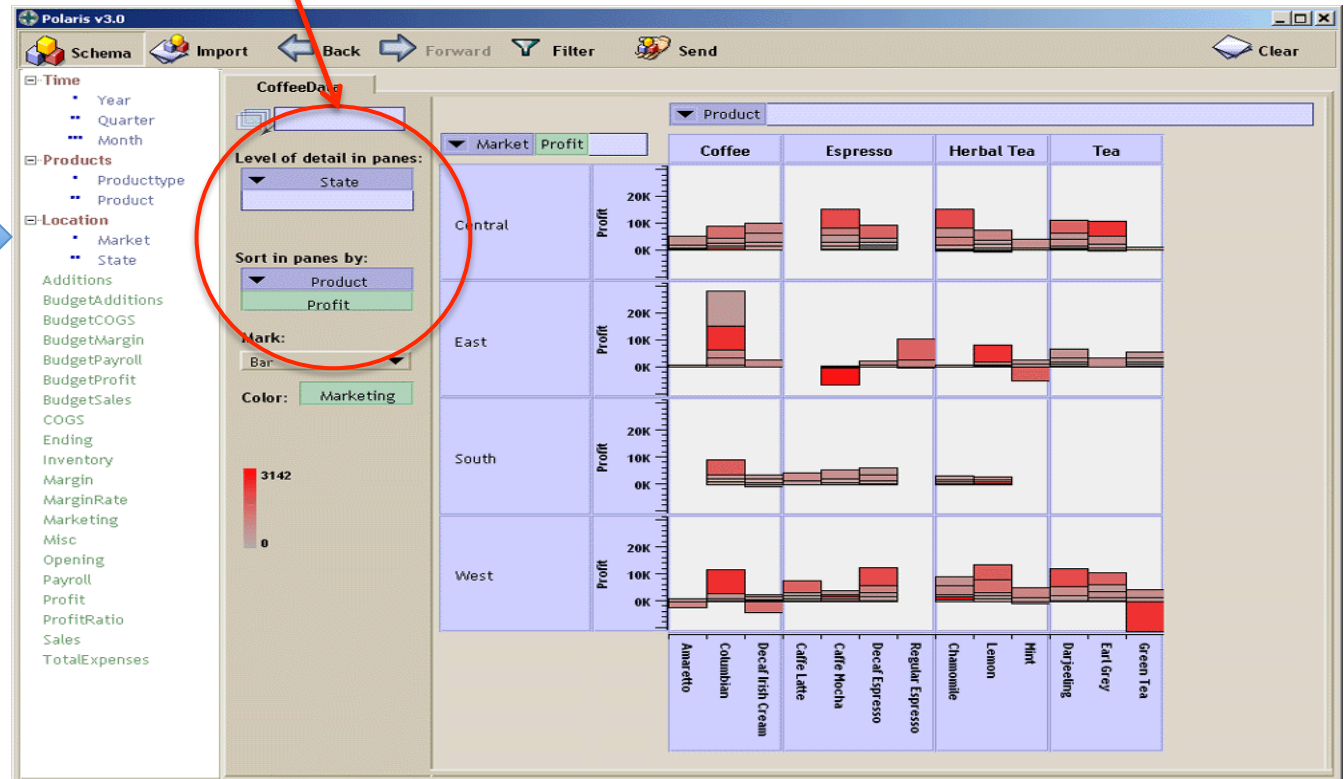


back-end queries: data selection, partition into panes

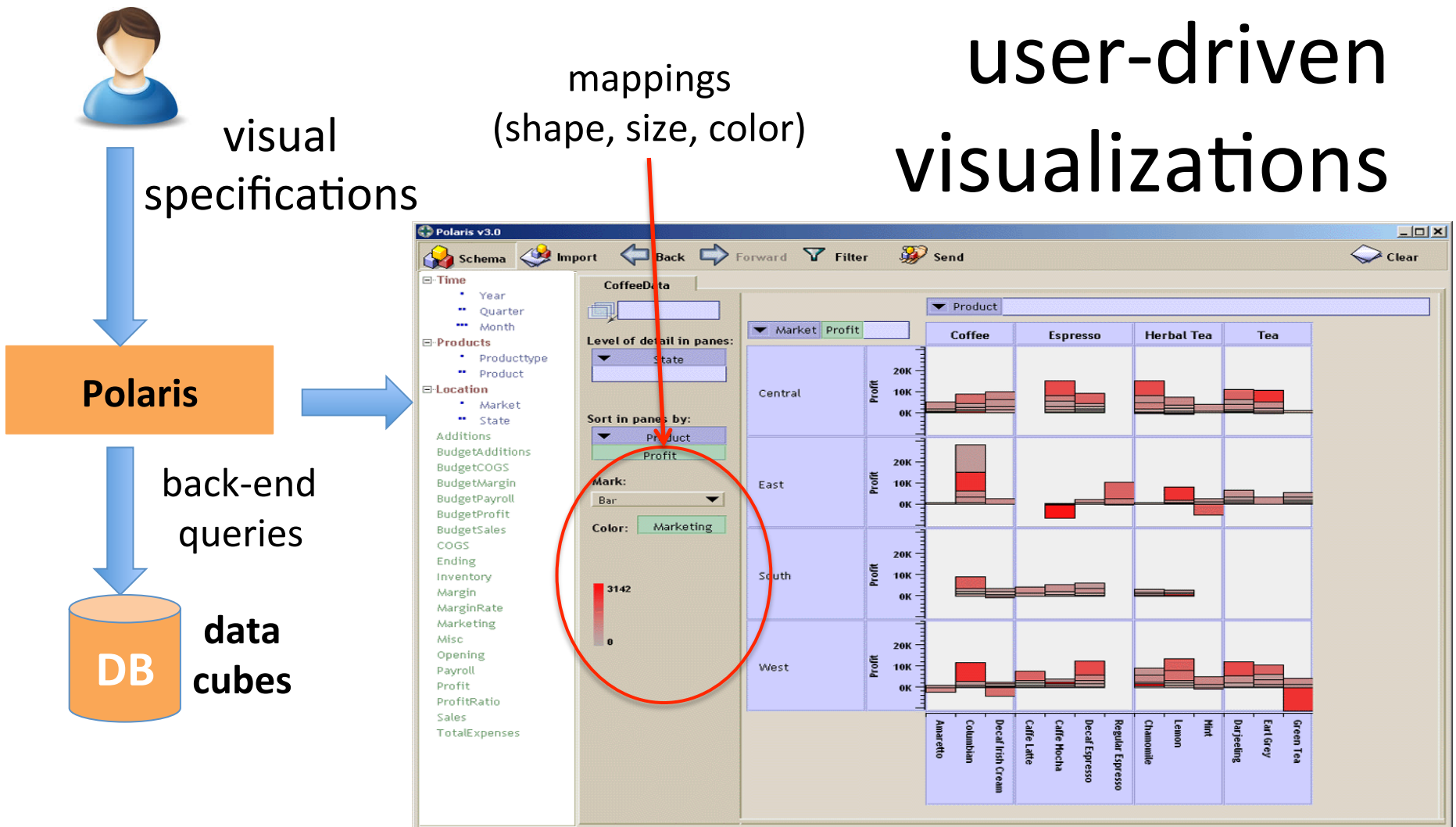


transformations
(group by, sort)

user-driven visualizations

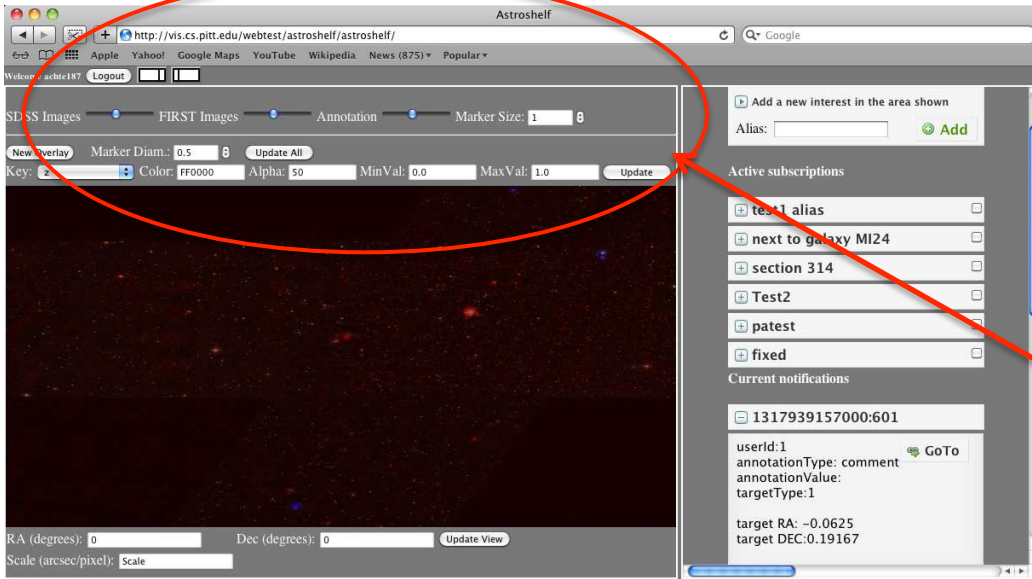


back-end queries: data transformations
(group, sort, aggregate within each pane)



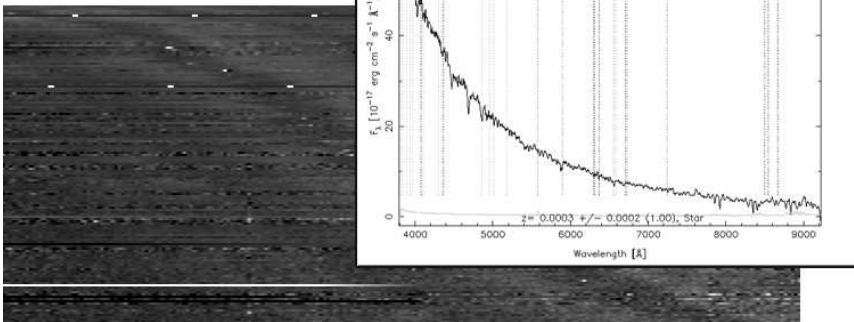
back-end queries: graphical transformations (rener and visualize)

collaborative exploration



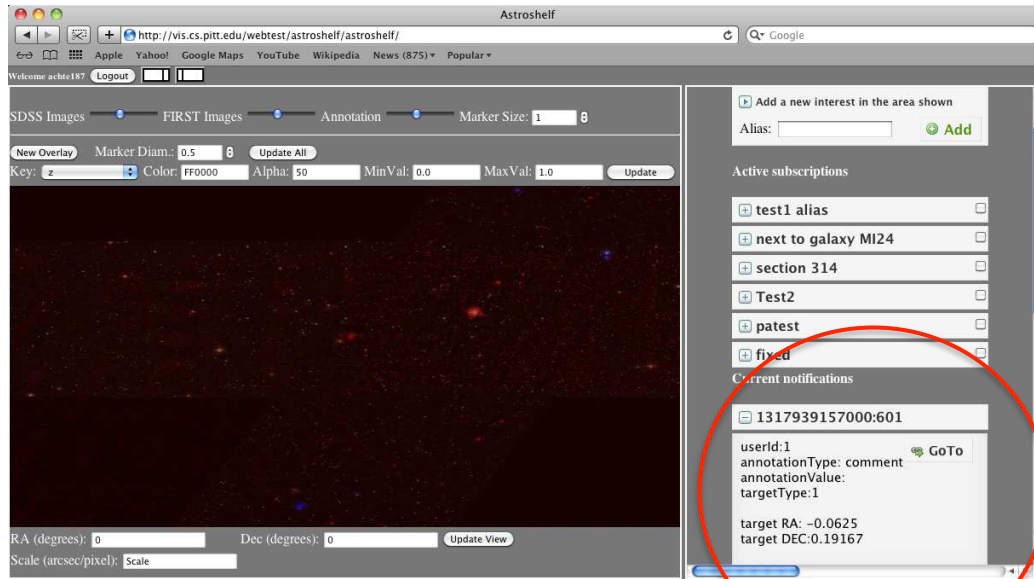
live annotations

Sky View



exploration for sky objects/patterns

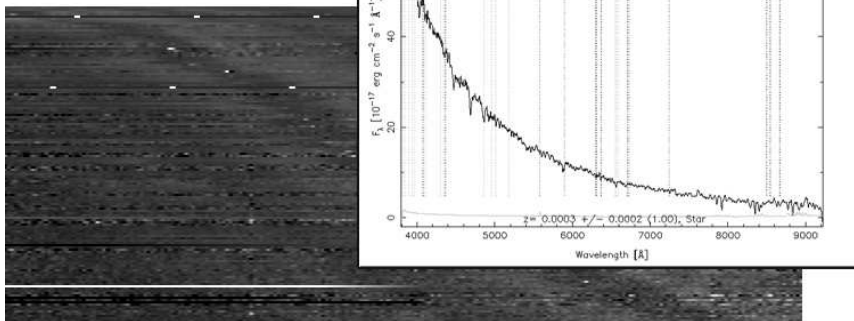
Live Annotations



collaborative
exploration

stream based
notifications

Sky View



exploration for
sky objects/patterns

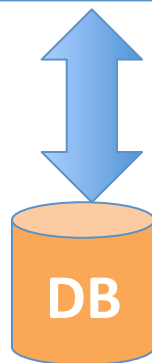
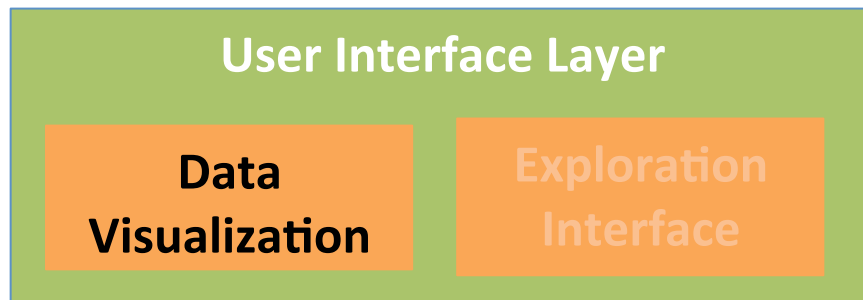


data visualization

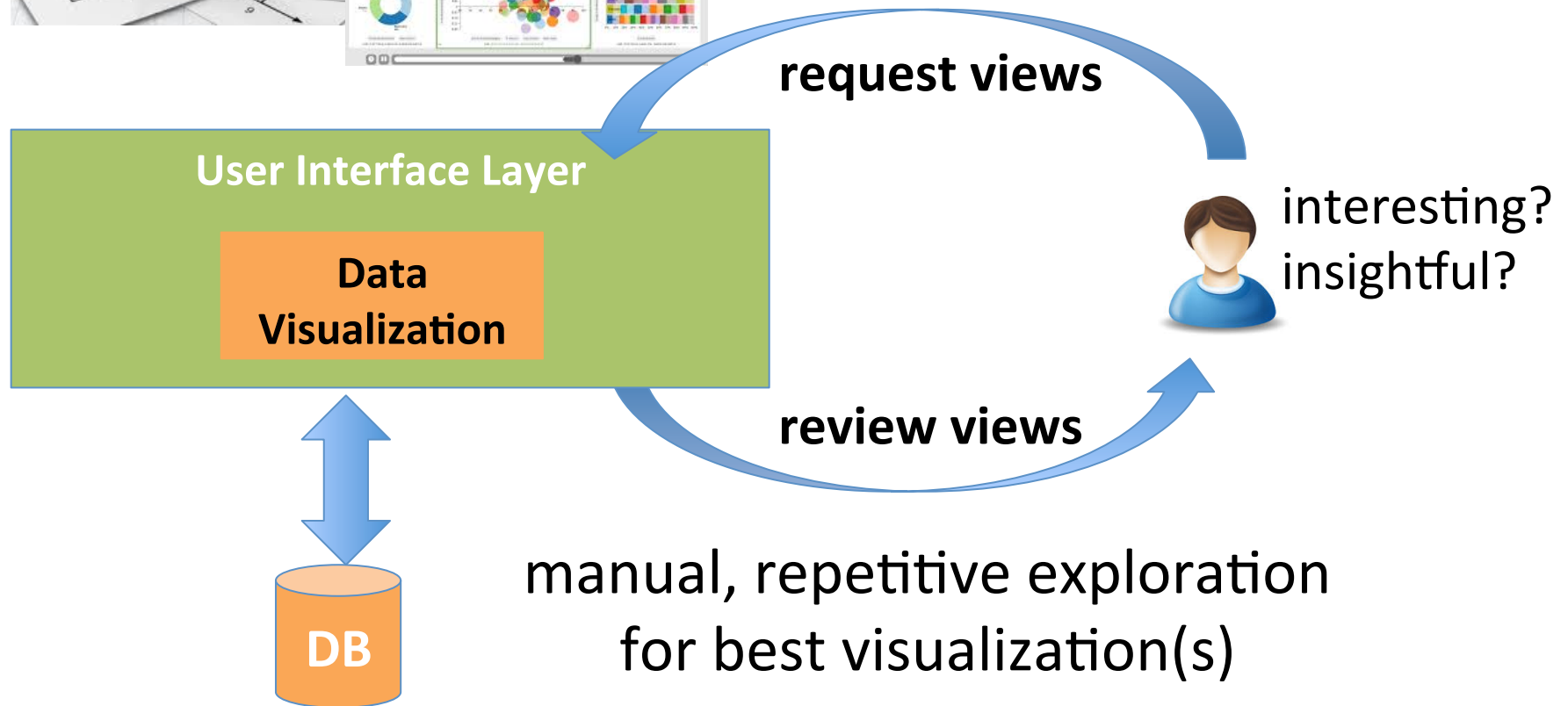
visualization tools

visual optimizations

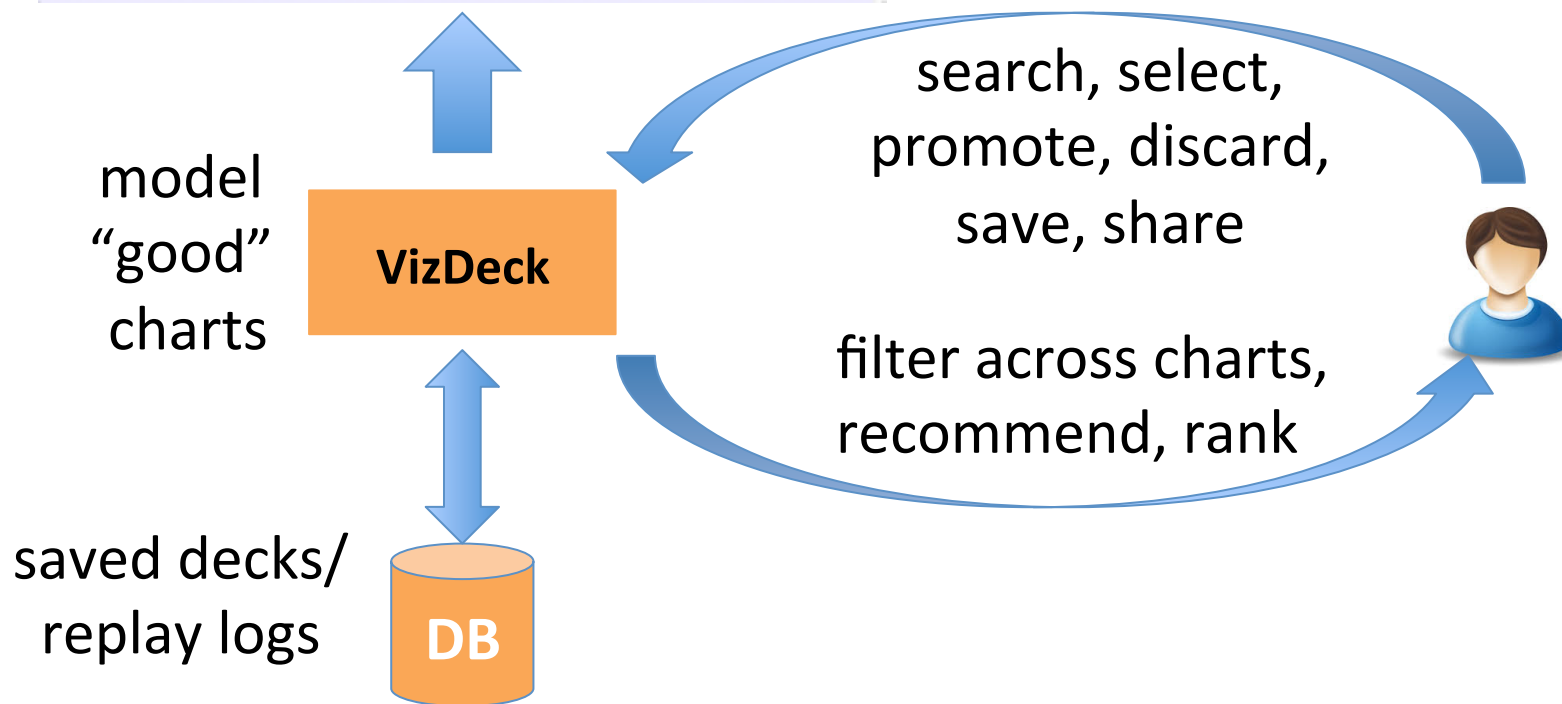
automatic visualization



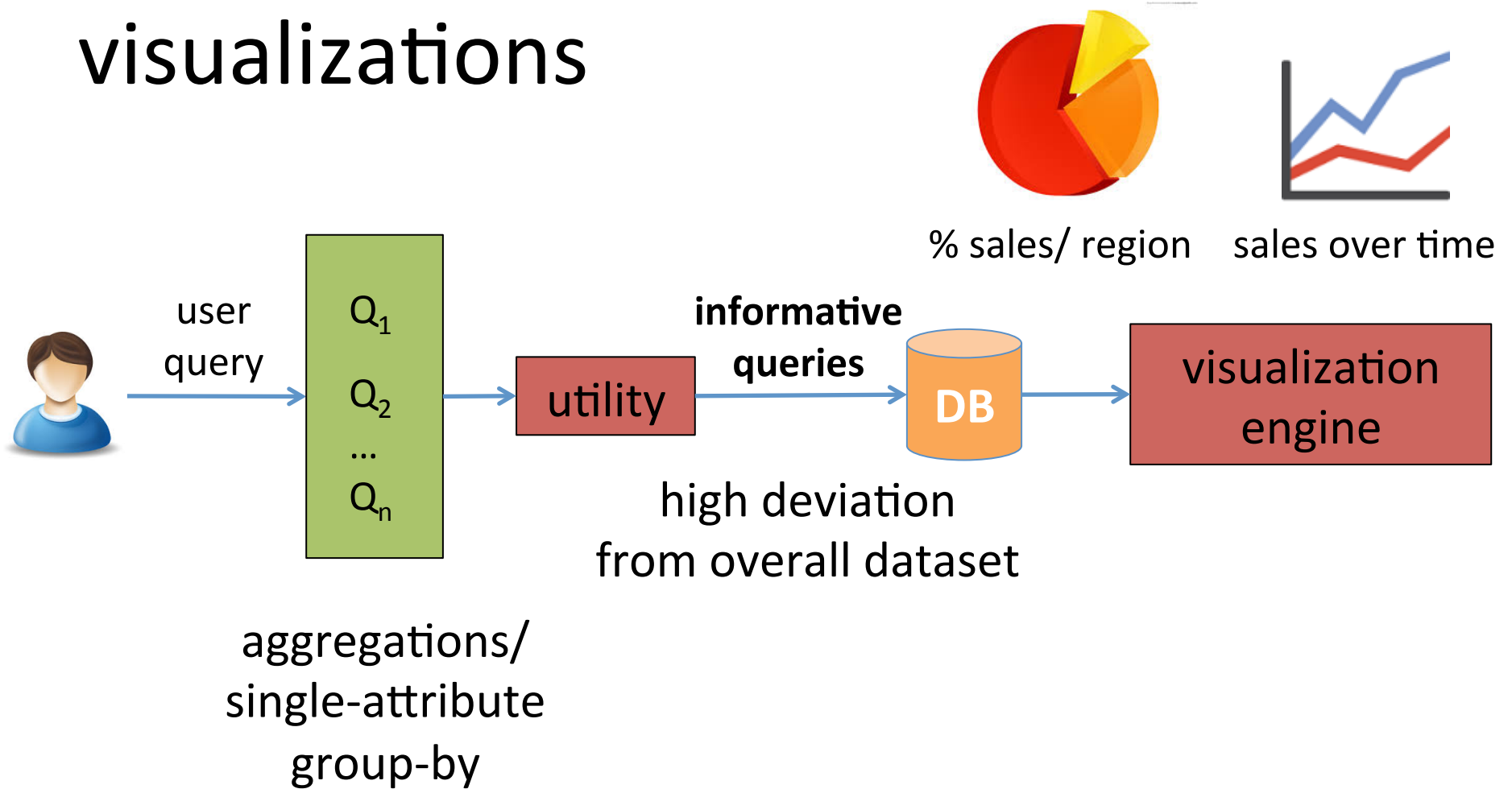
automatic visualization

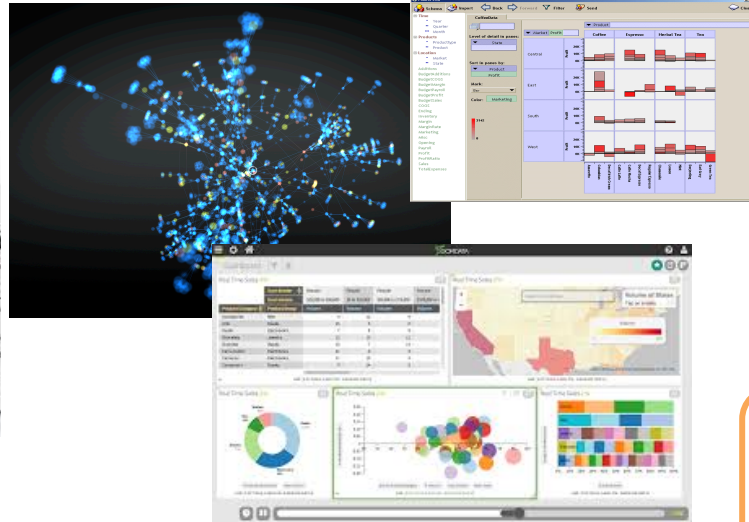


auto-ranked visualizations



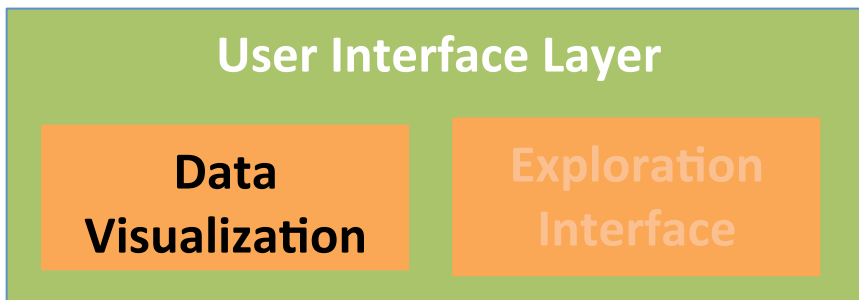
automatic visualizations



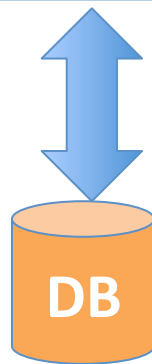


data visualization

visualization tools

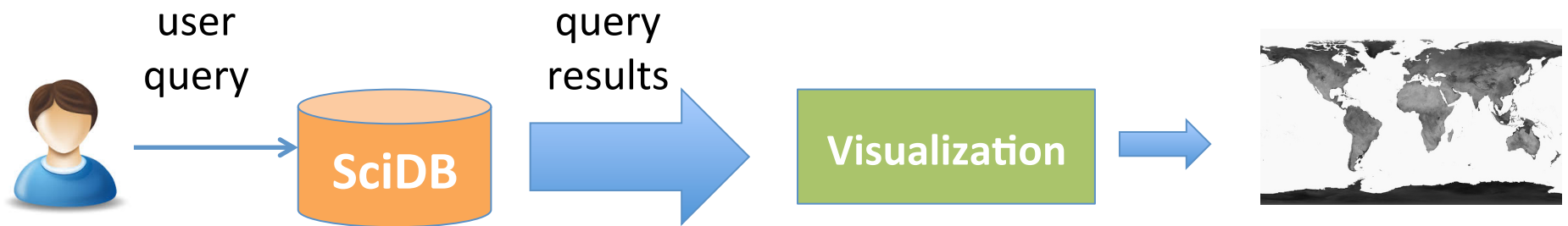


visual optimizations



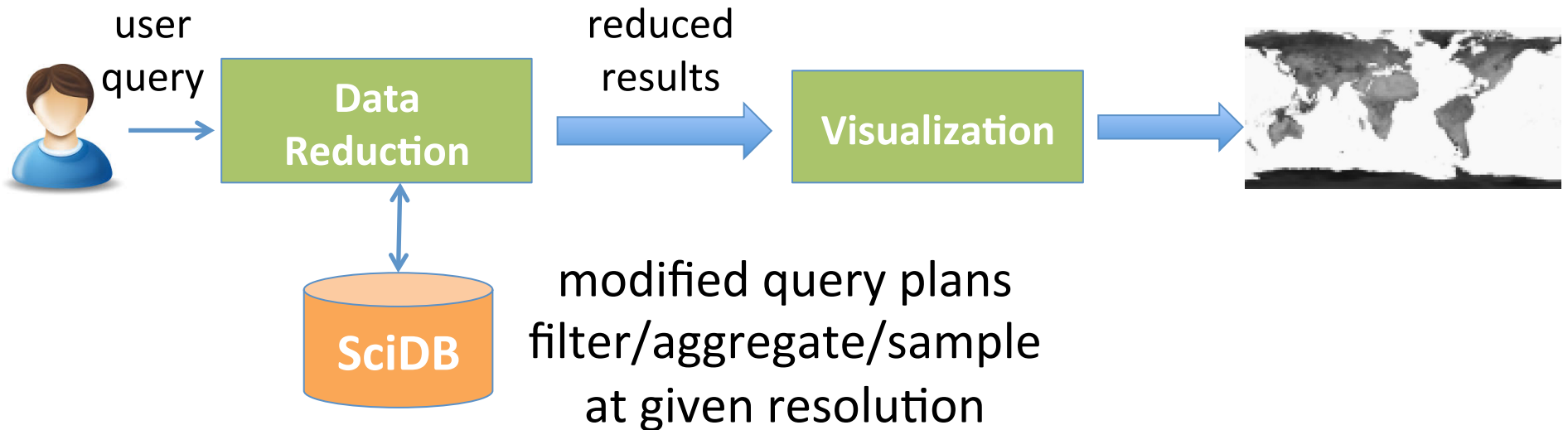
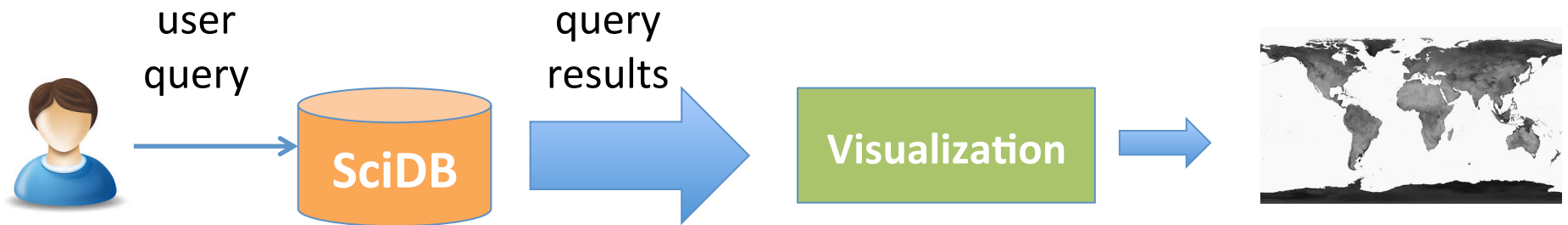
automatic visualization

resolution reduction

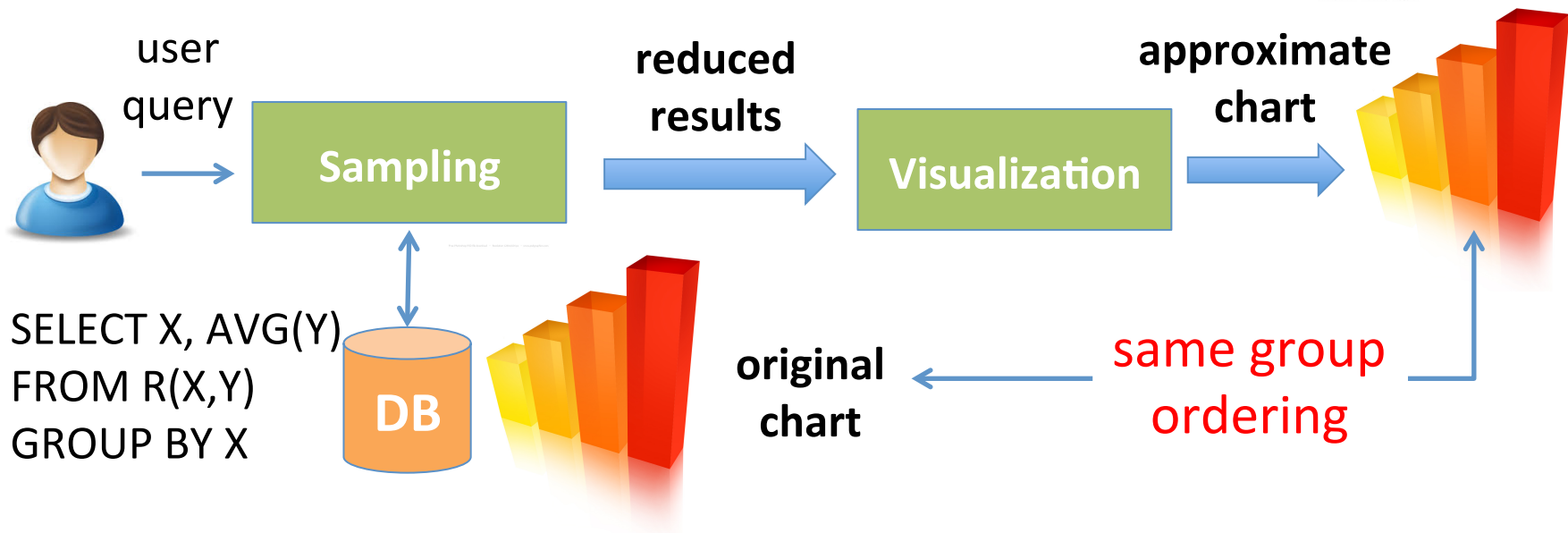


expensive, ineffective on big data sets

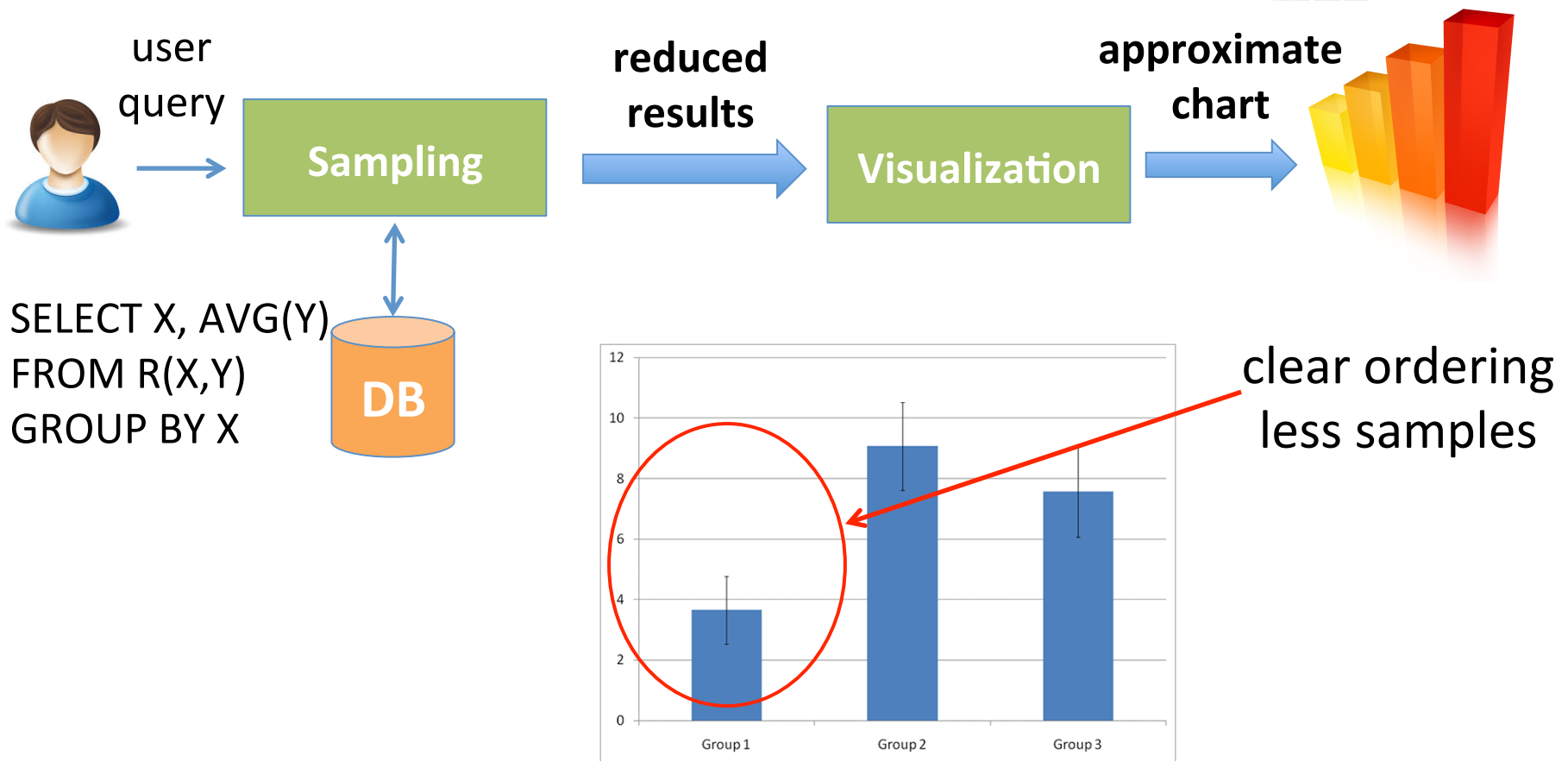
resolution reduction



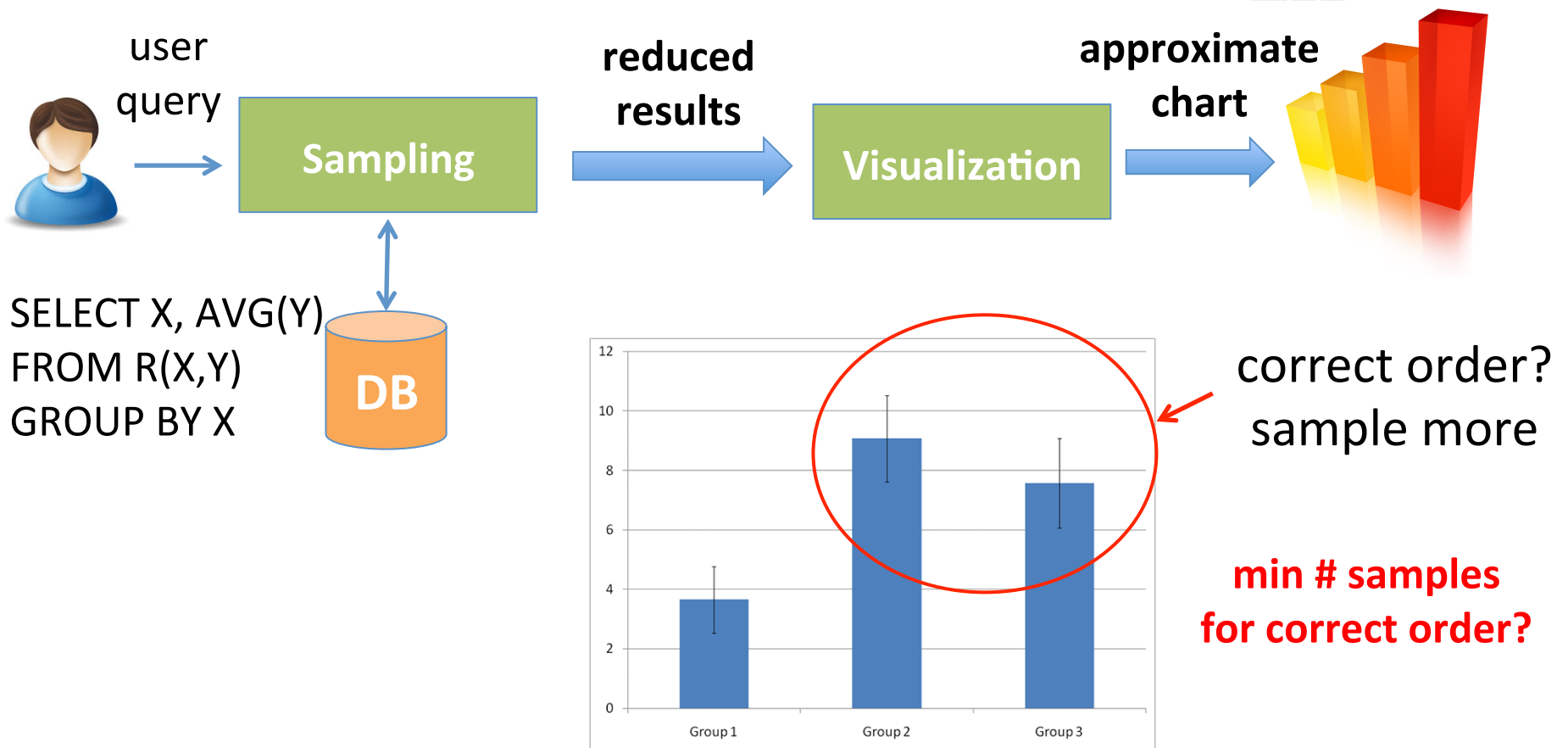
approximate visualizations



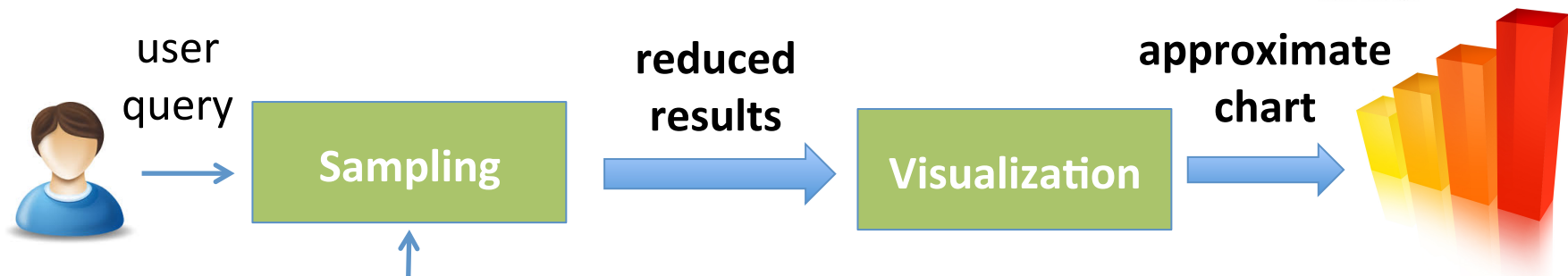
approximate visualizations



approximate visualizations



approximate visualizations



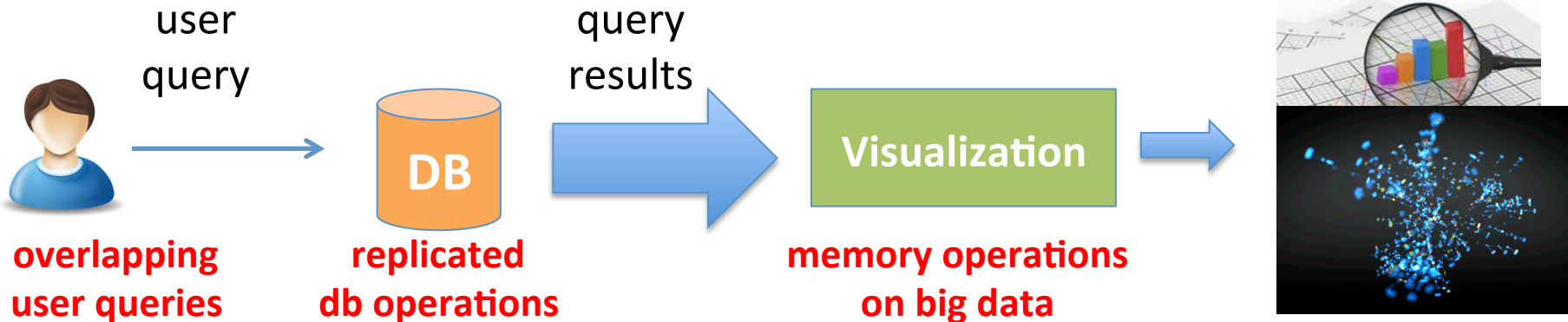
SELECT X, AVG(Y)
FROM R(X,Y)
GROUP BY X



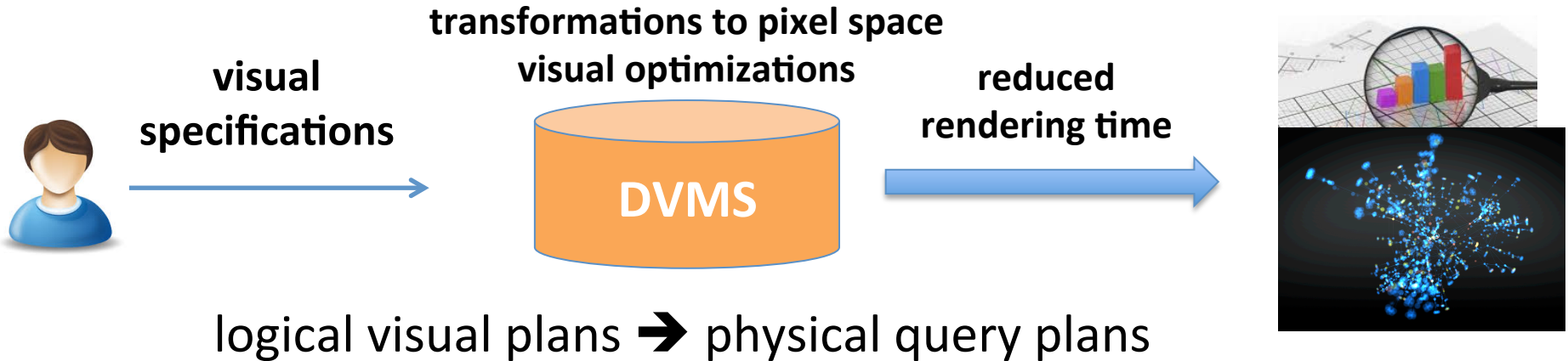
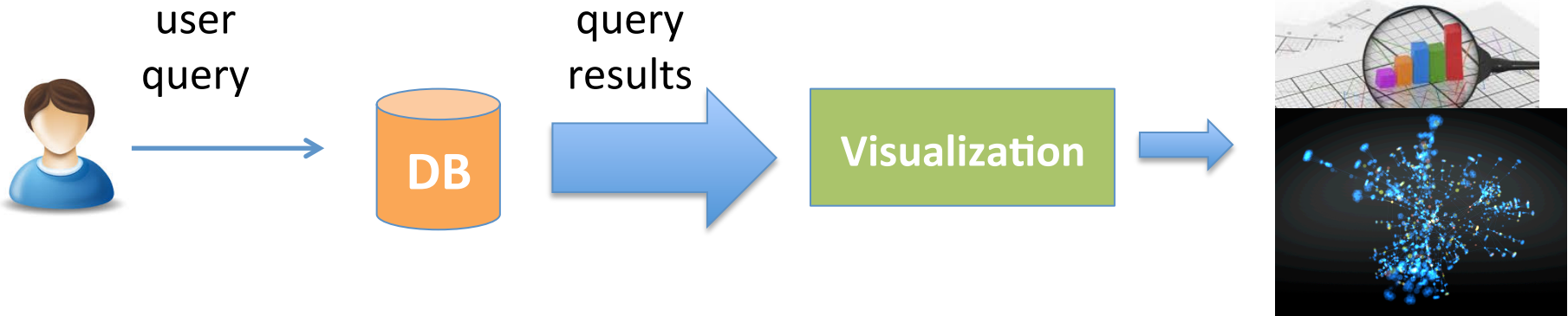
#samples	Group 1	Group 2	Group 3	Group 4
1	[60,90]	[20,50]	[10,40]	[40,70]
20	[64,84]	[30,48]	[15,35]	[45,65]
21	[66,84], I	[30,48]	[17,35]	[46,64]
70	[66,84], I	[40,47]	[17,32], I	[46,53]

sampling phases/ confidence intervals

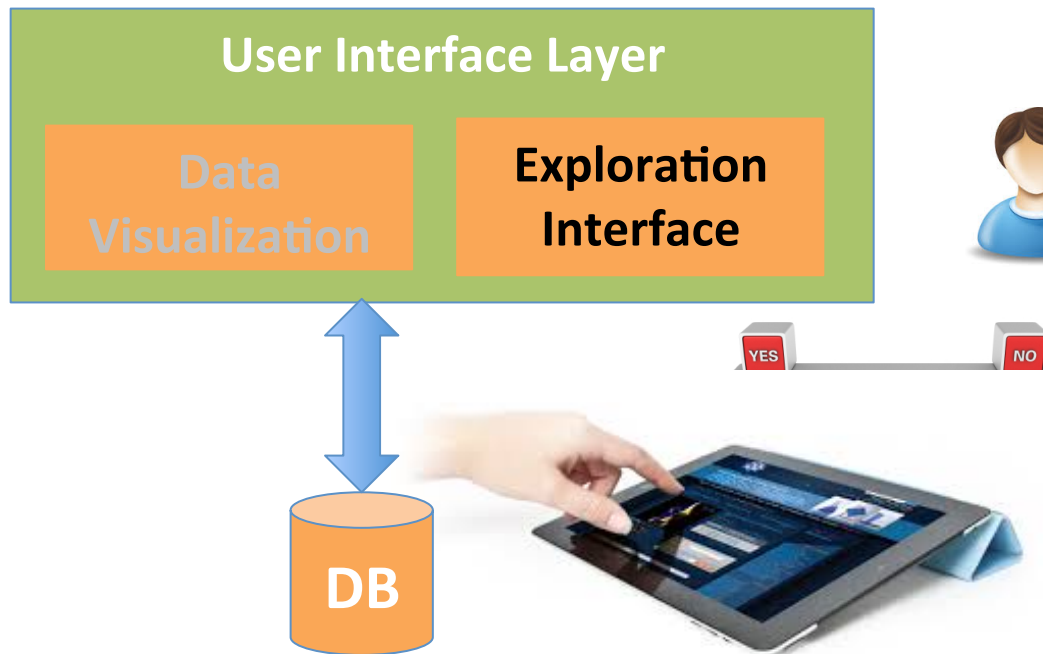
visualization management



visualization management



exploration interfaces

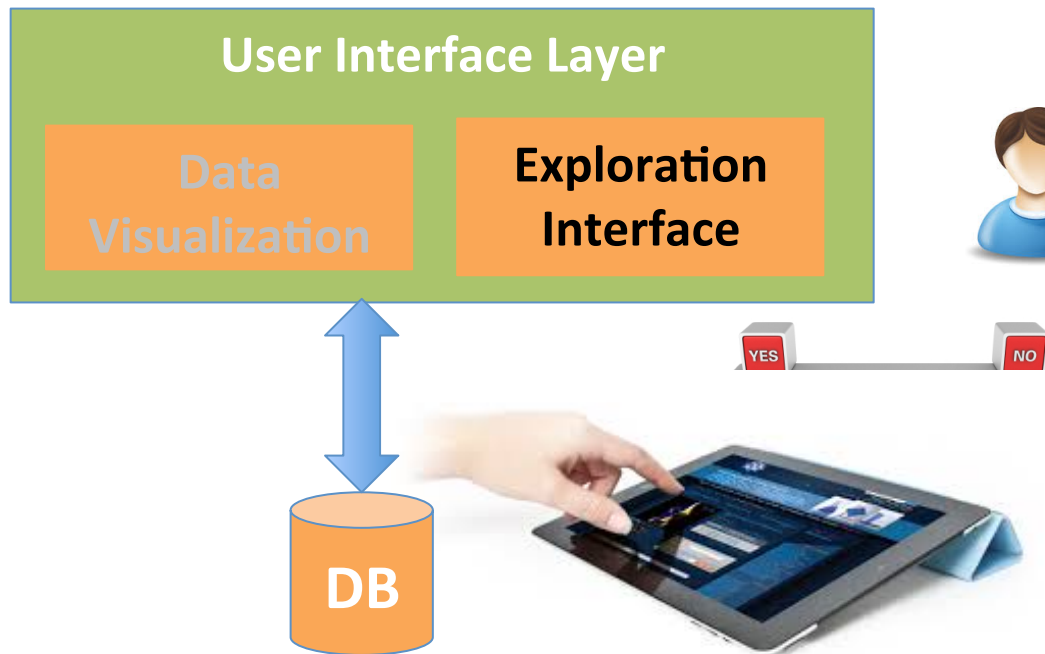


**automatic
exploration**

**assisted query
formulation**

novel query interfaces

exploration interfaces

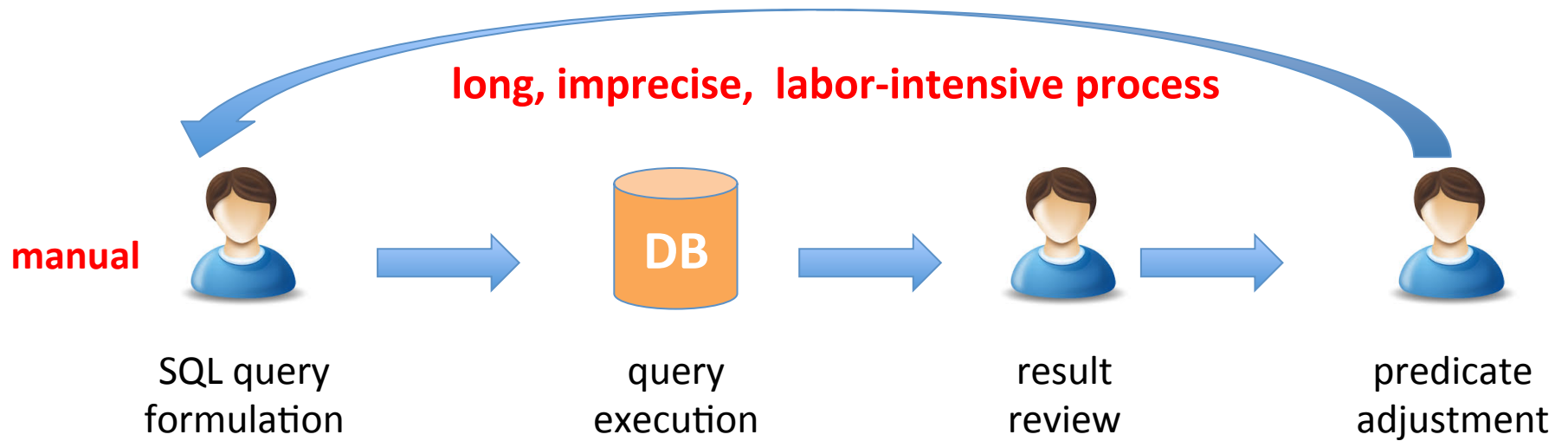


**automatic
exploration**

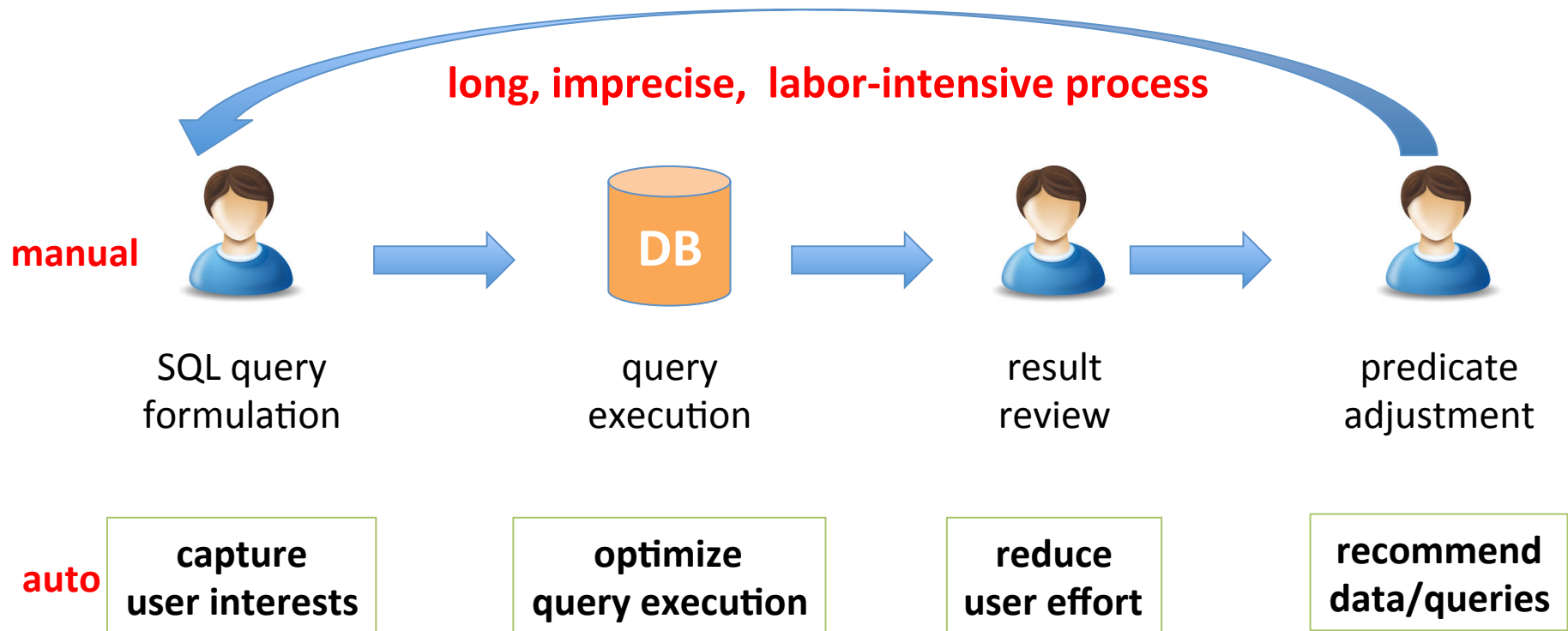
**assisted query
formulation**

novel query interfaces

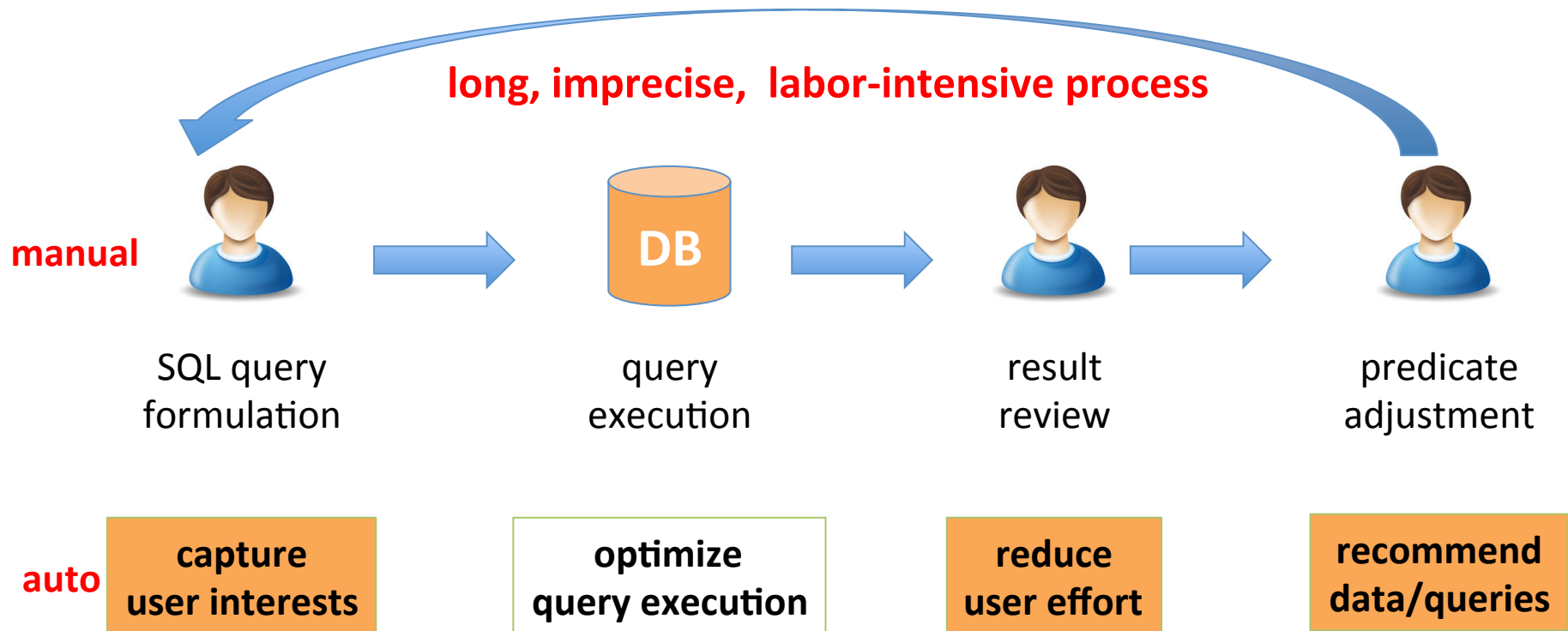
manual vs automatic data exploration



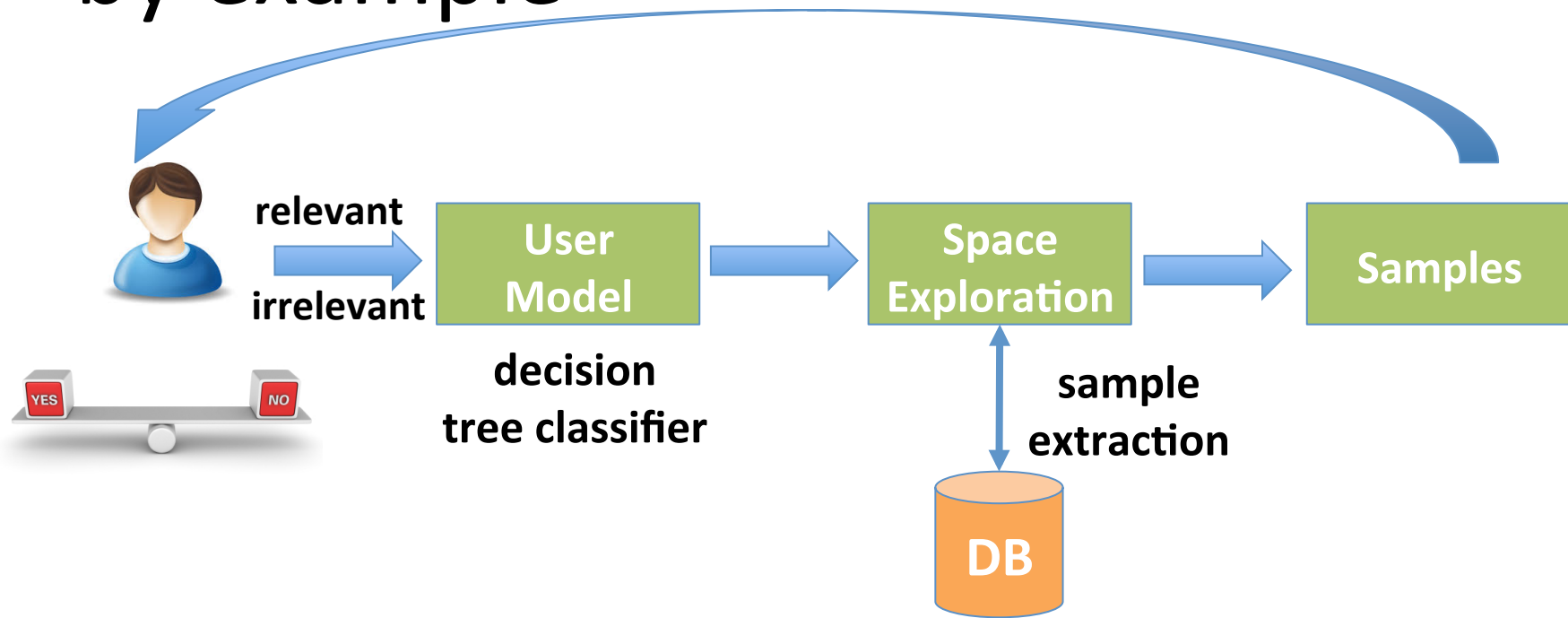
manual vs automatic data exploration



manual vs automatic data exploration

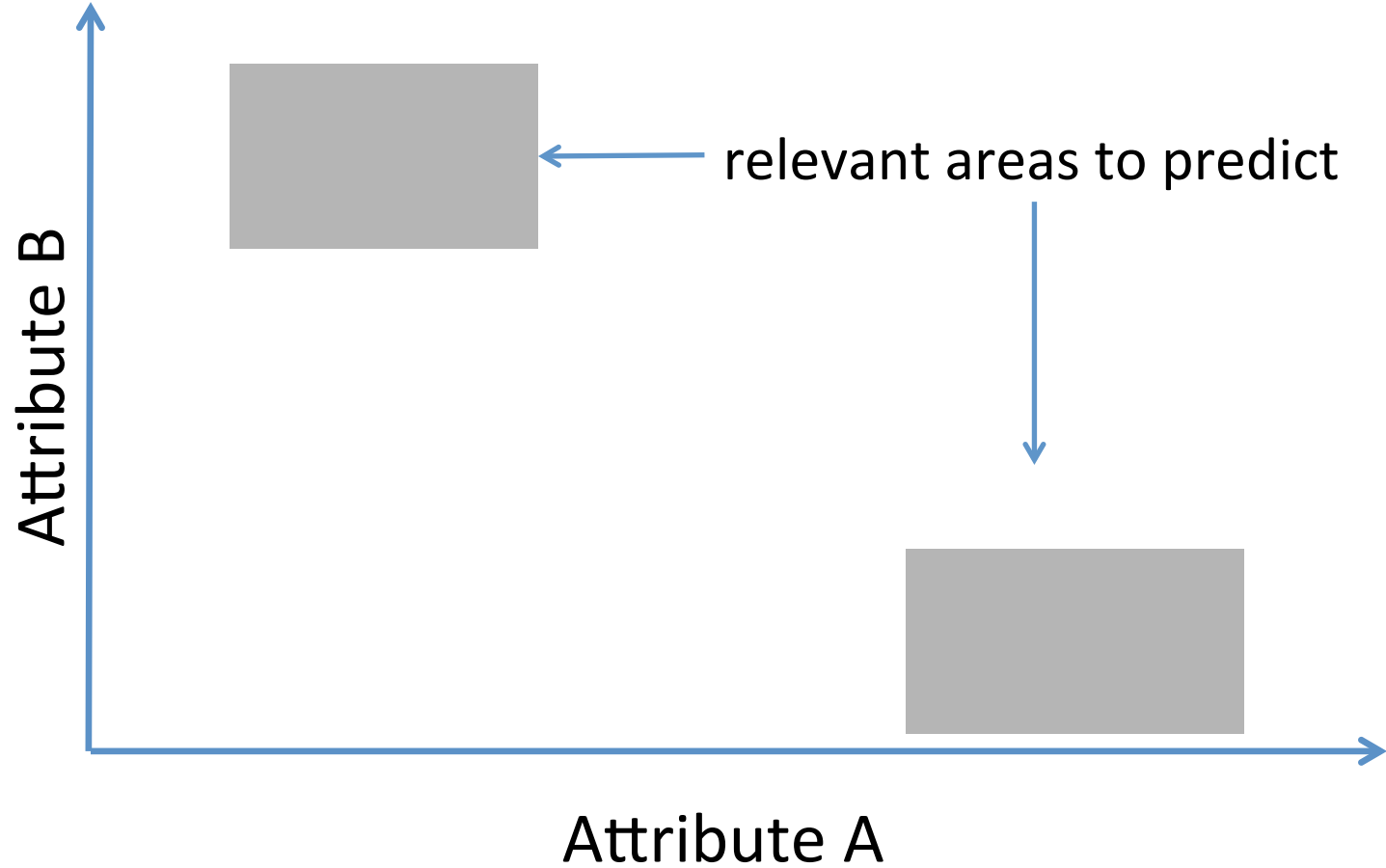


explore by example



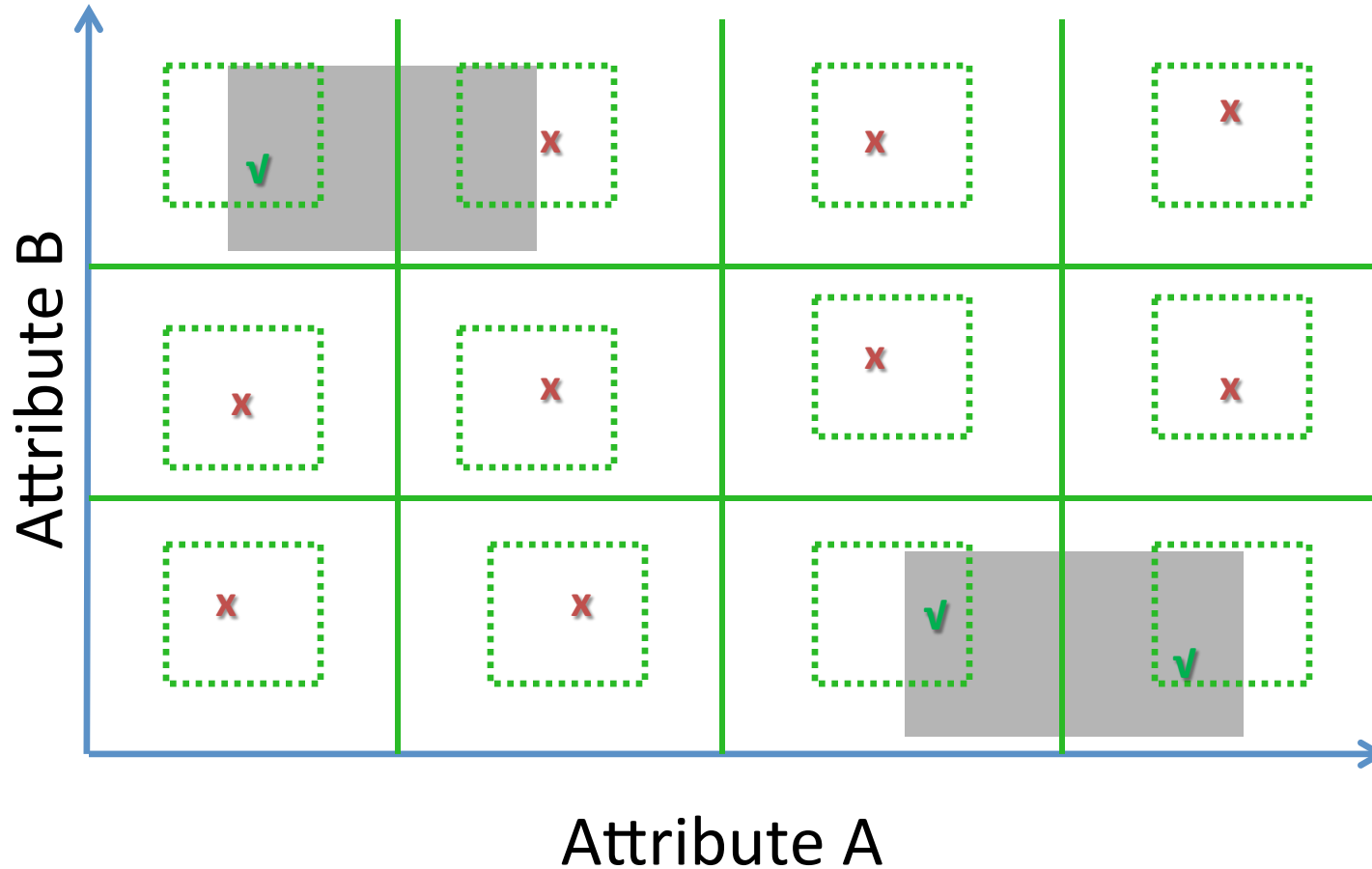
effectiveness vs efficiency
sampling areas? sampling size?

explore by example

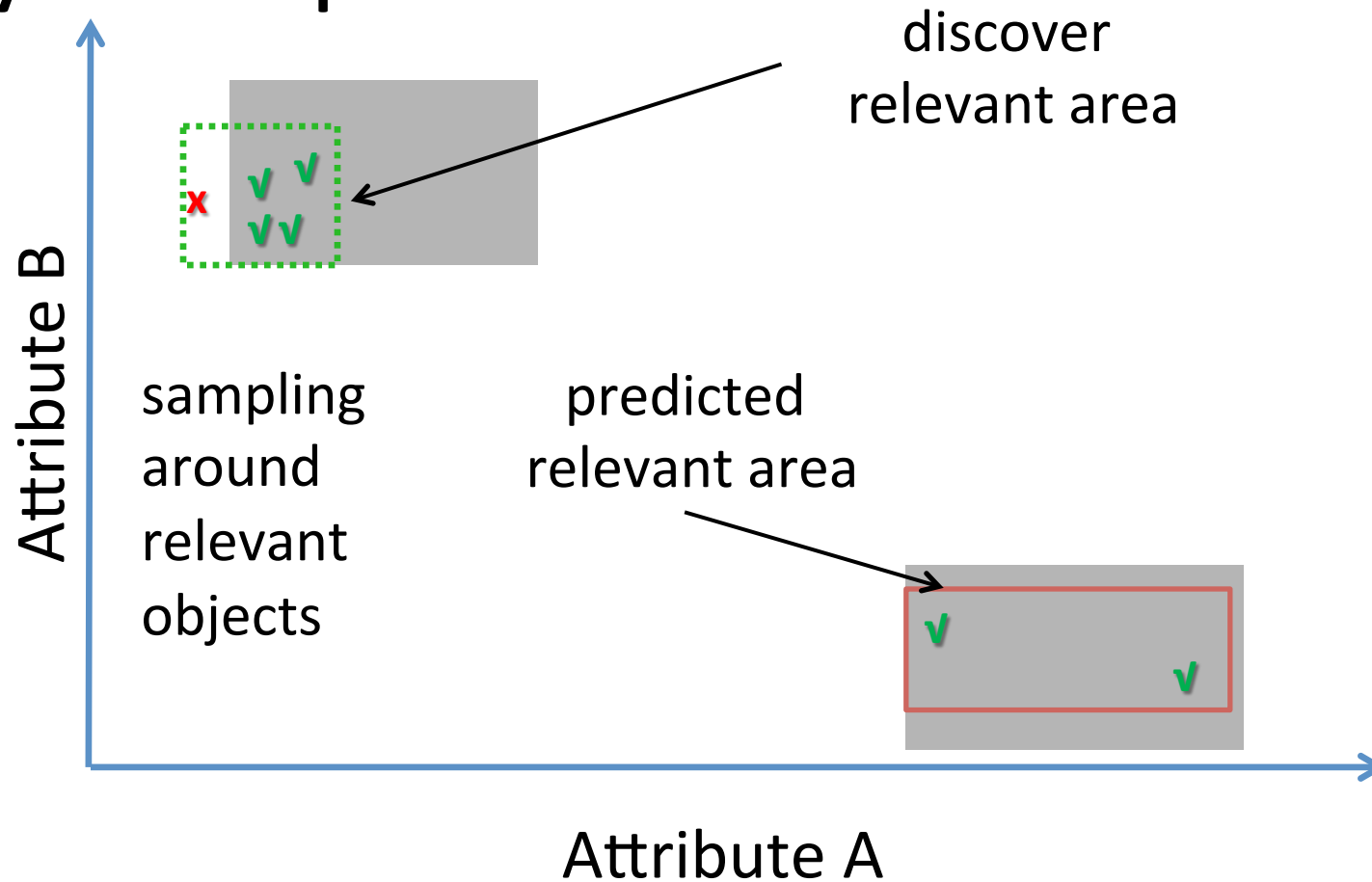


explore by example

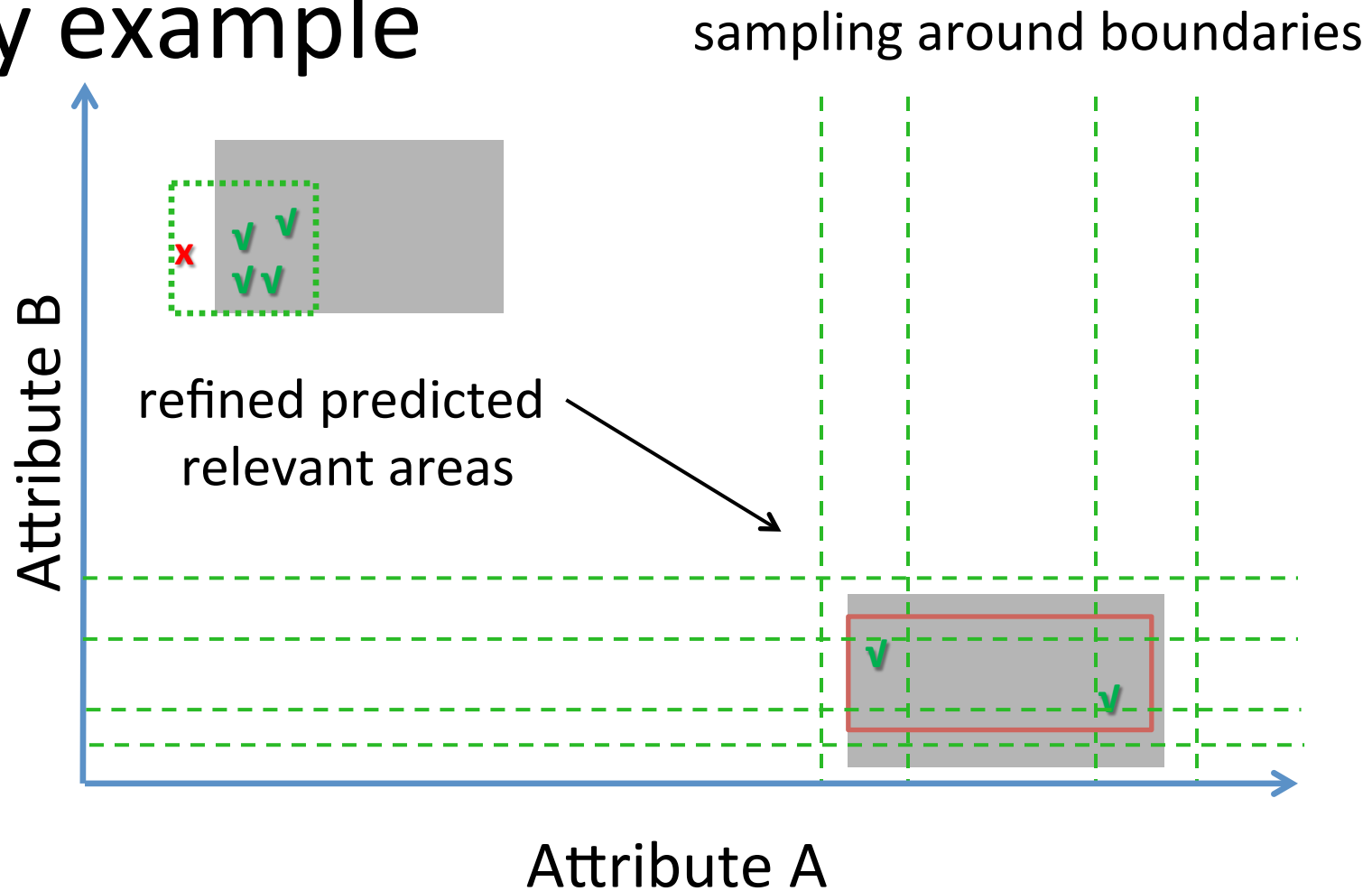
uniform sampling across domain



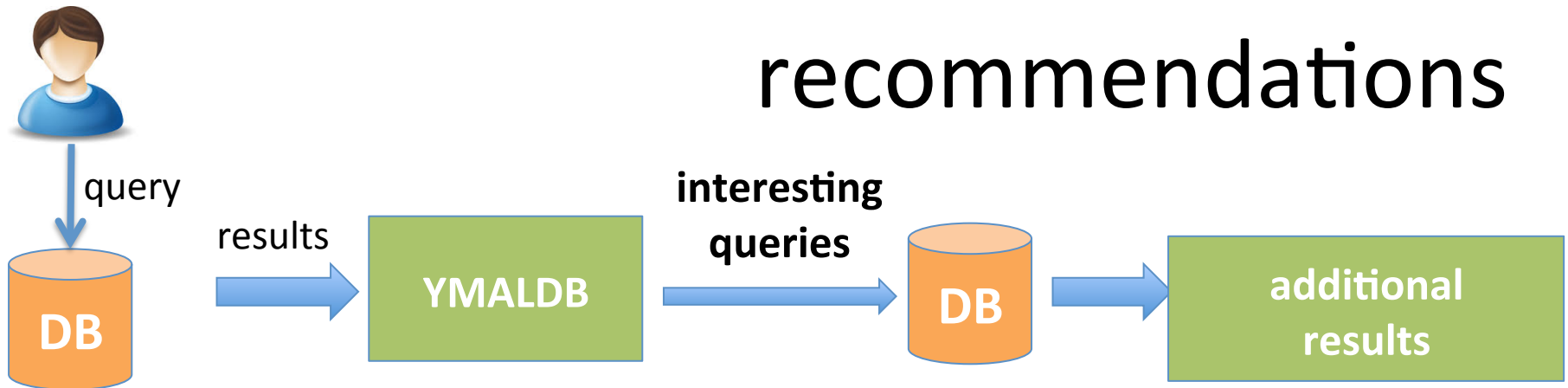
explore by example



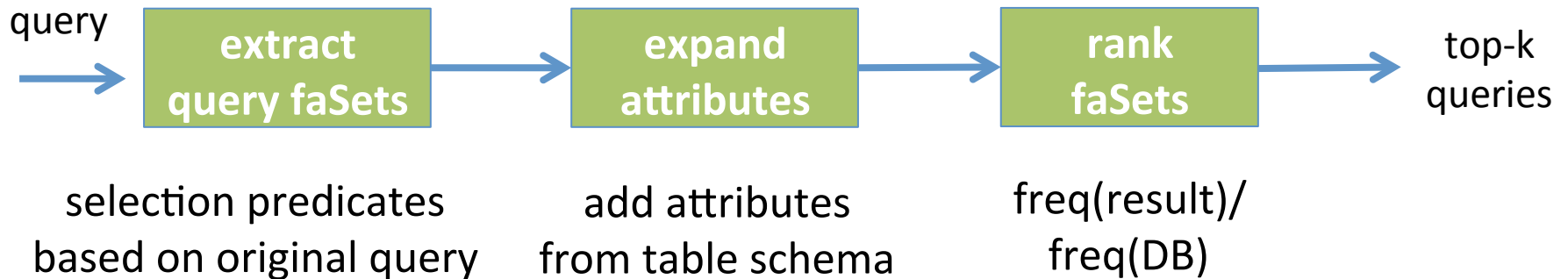
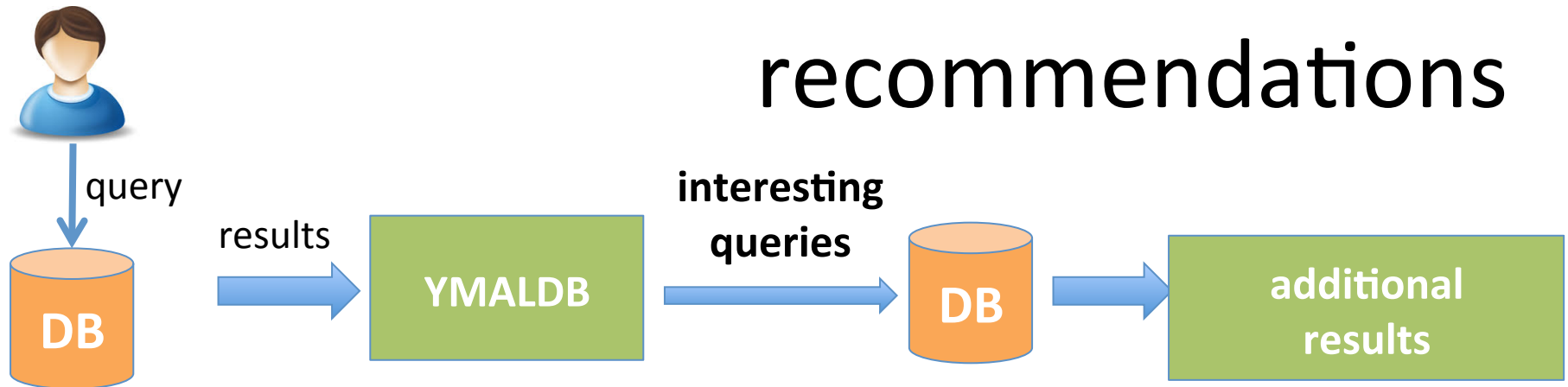
explore by example



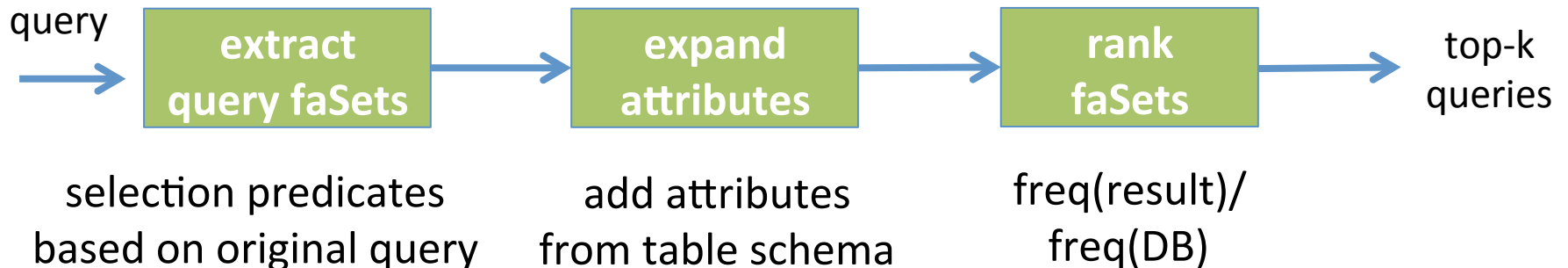
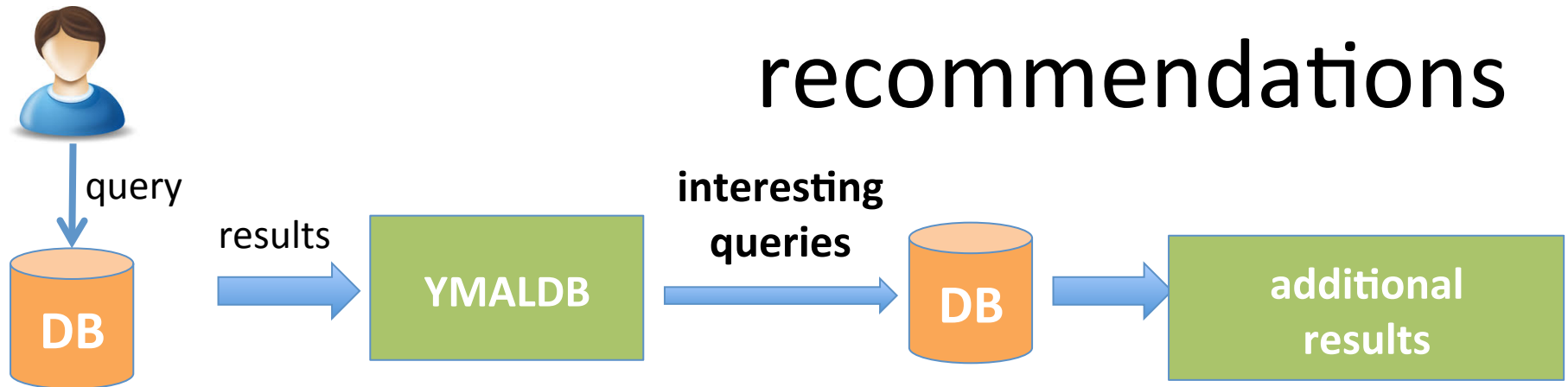
result recommendations



result recommendations

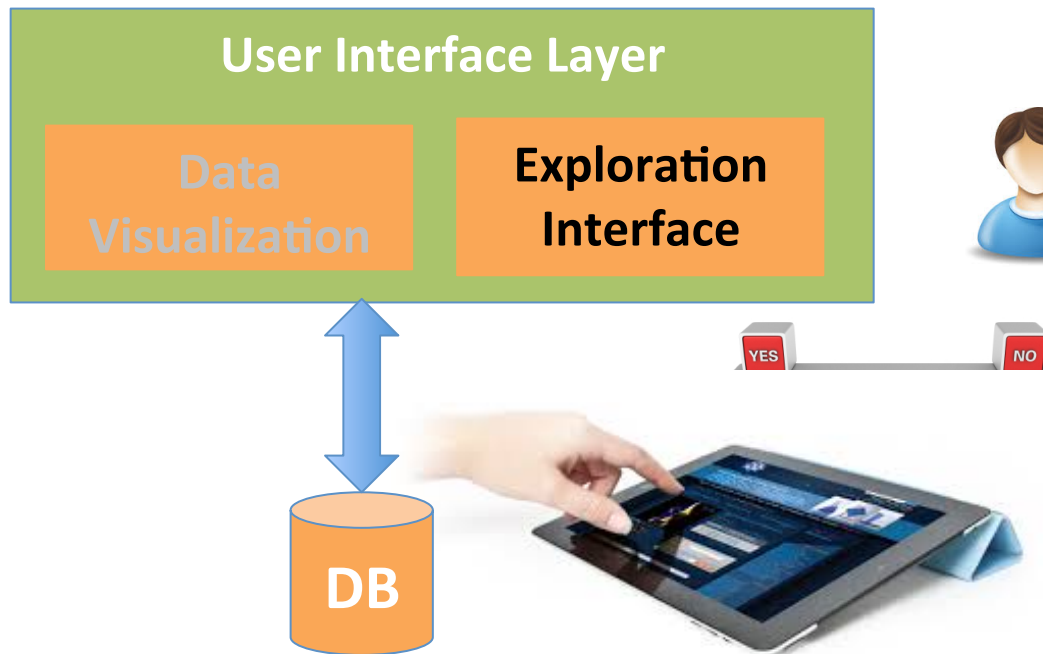


result recommendations



title, year, genre of Scorsese movies + *title, year, genre, **country** of Scorsese movies* = *many Scorsese movies are related to Italy*

exploration interfaces

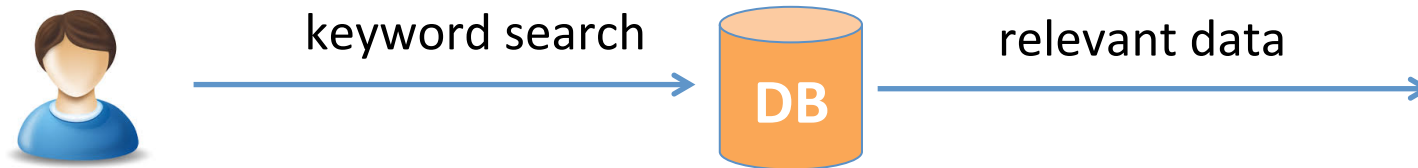
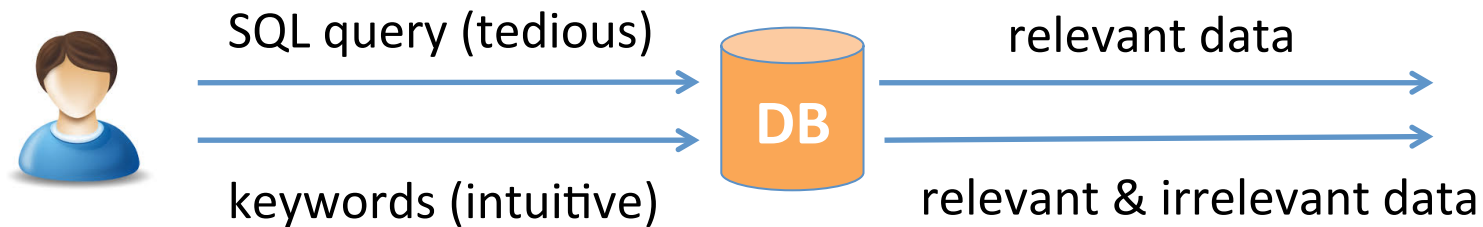


**automatic
exploration**

**assisted query
formulation**

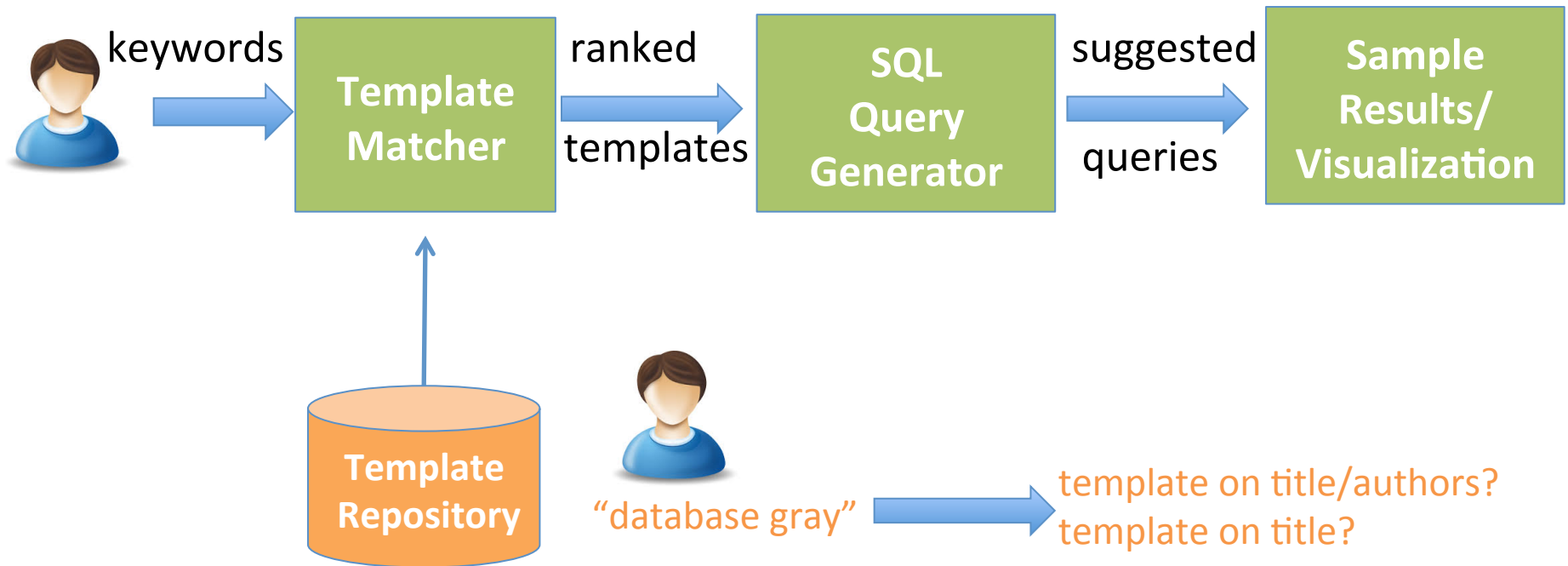
novel query interfaces

keyword-based query suggestions

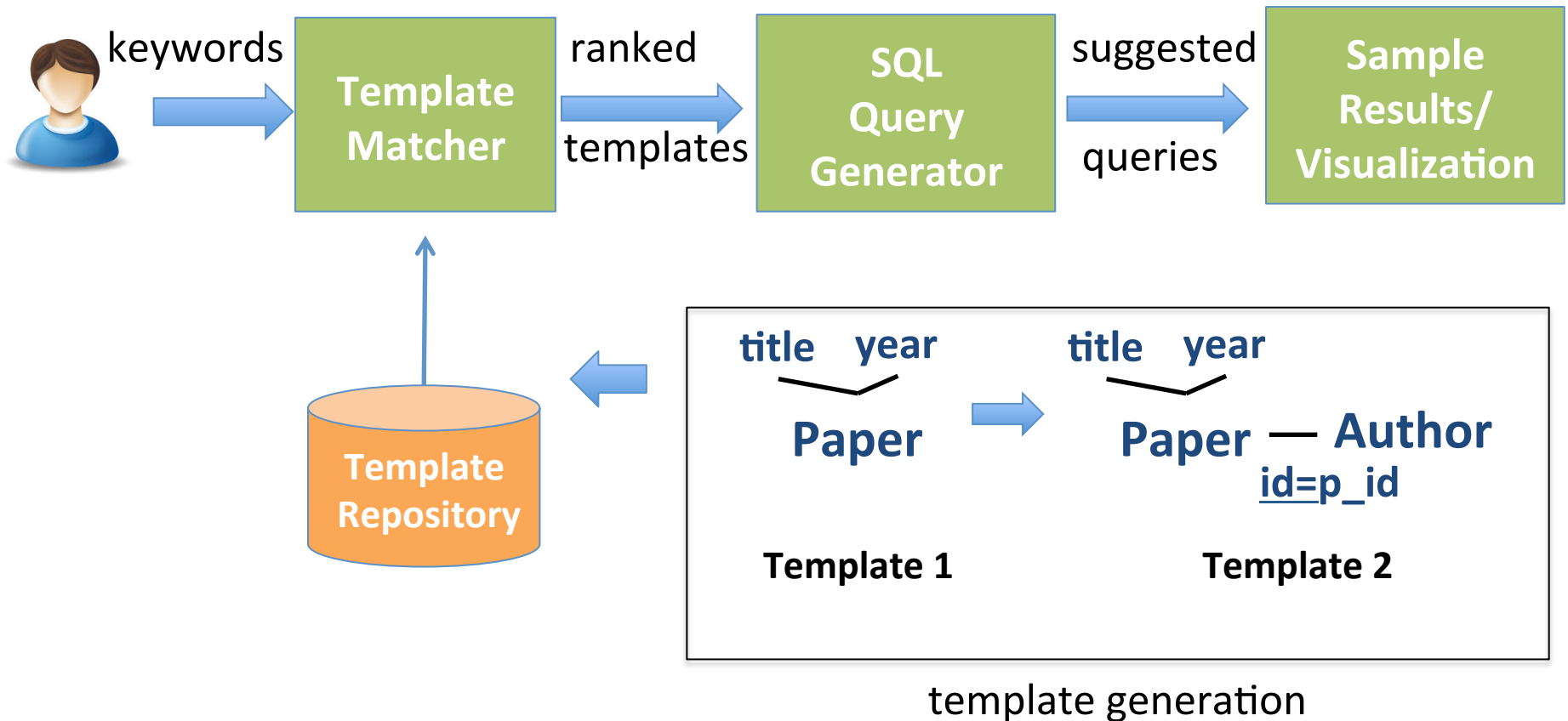


how we can discover **relevant** queries?

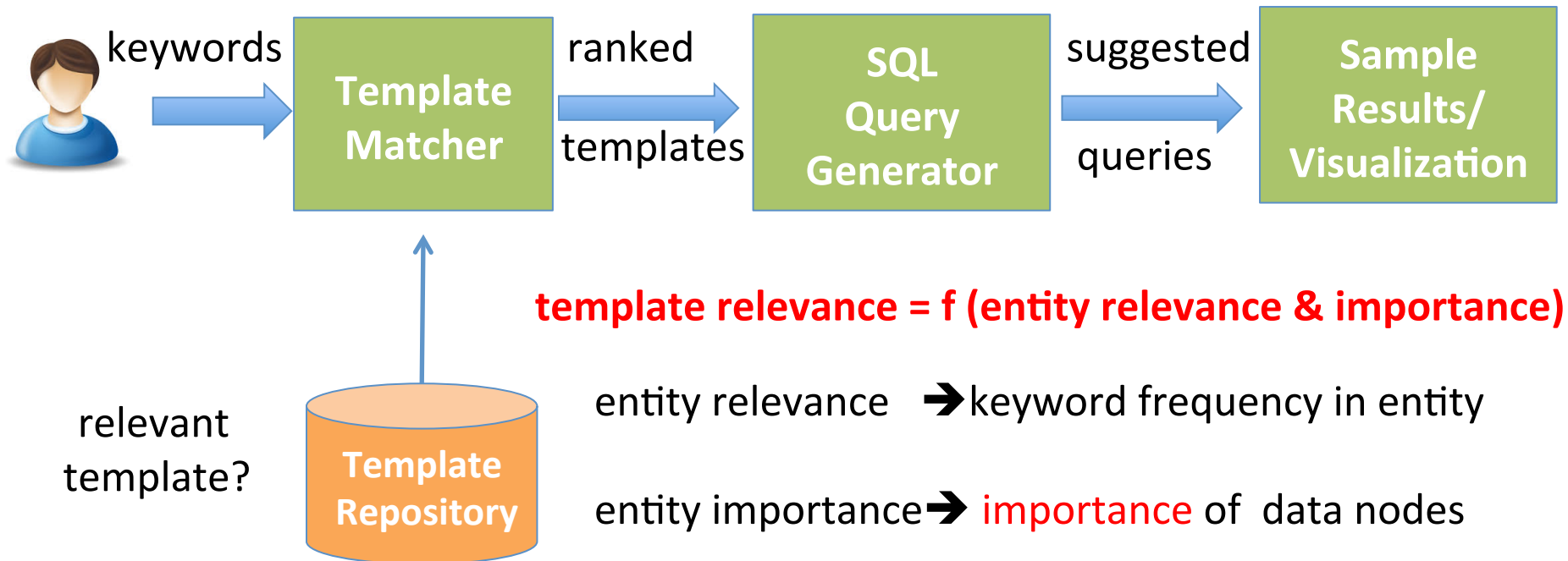
keyword-based query suggestions



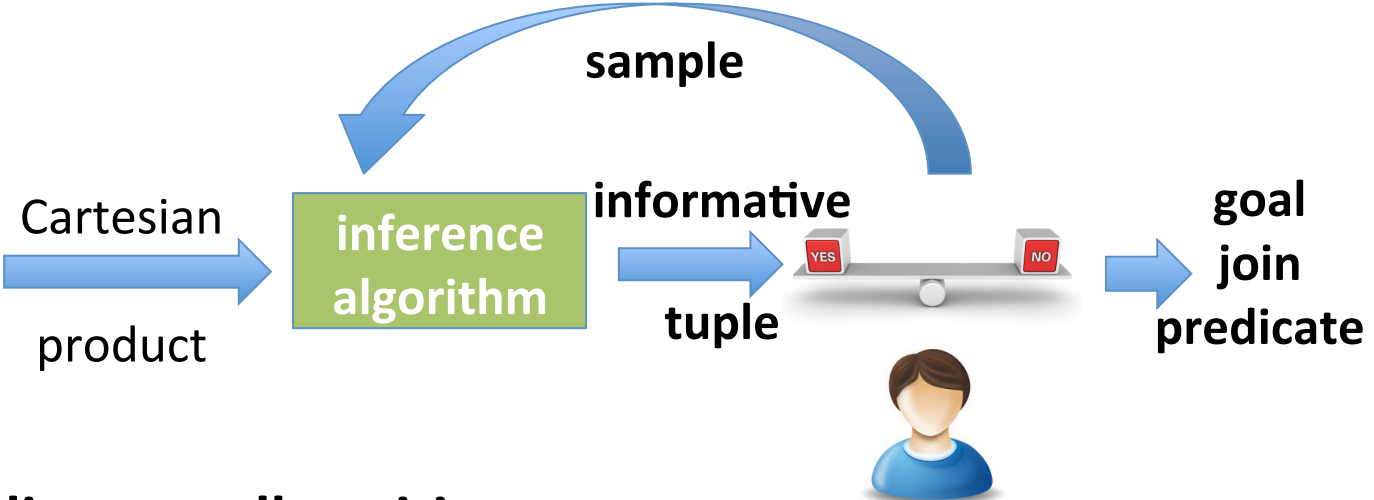
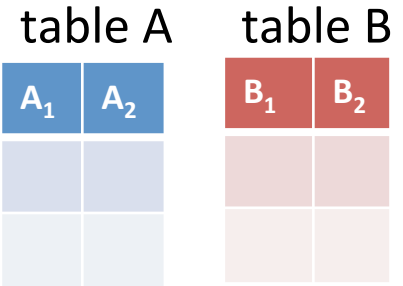
keyword-based query suggestions



keyword-based query suggestions



equi-join inference

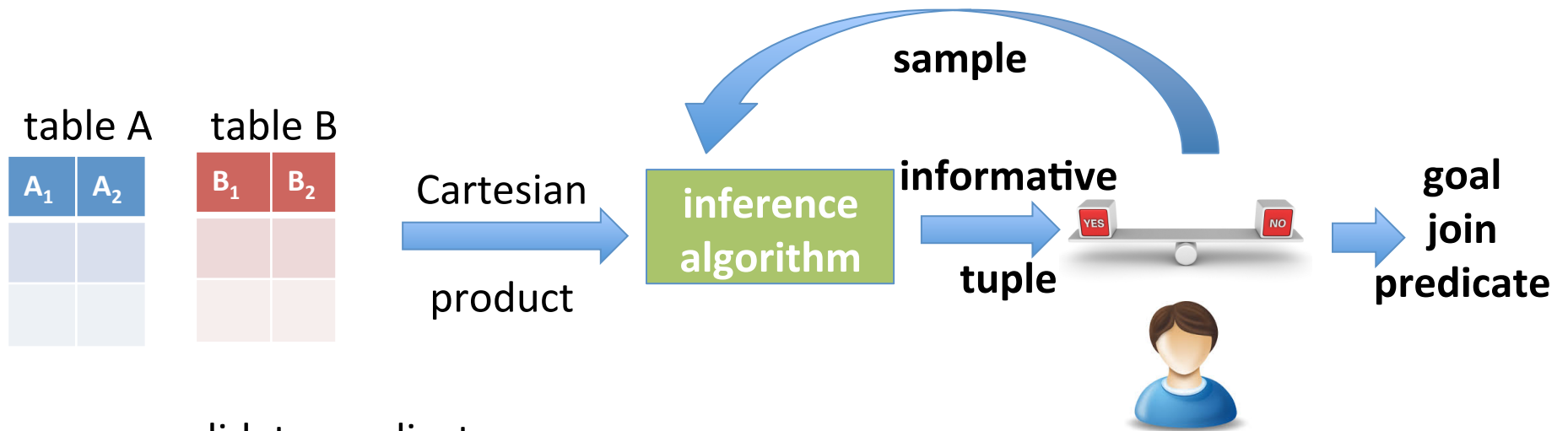


goal predicate:

**discover all positives
eliminate all negatives**

minimize user effort

equi-join inference



candidate predicates

(A1, B1)	(A1, B1)	(A1, B1)
(A1, B2)	(A2, B1)	(A2, B2)
(A1, B1)	(A1, B2)	(A2, B1)

prune predicates with uninformative tuples

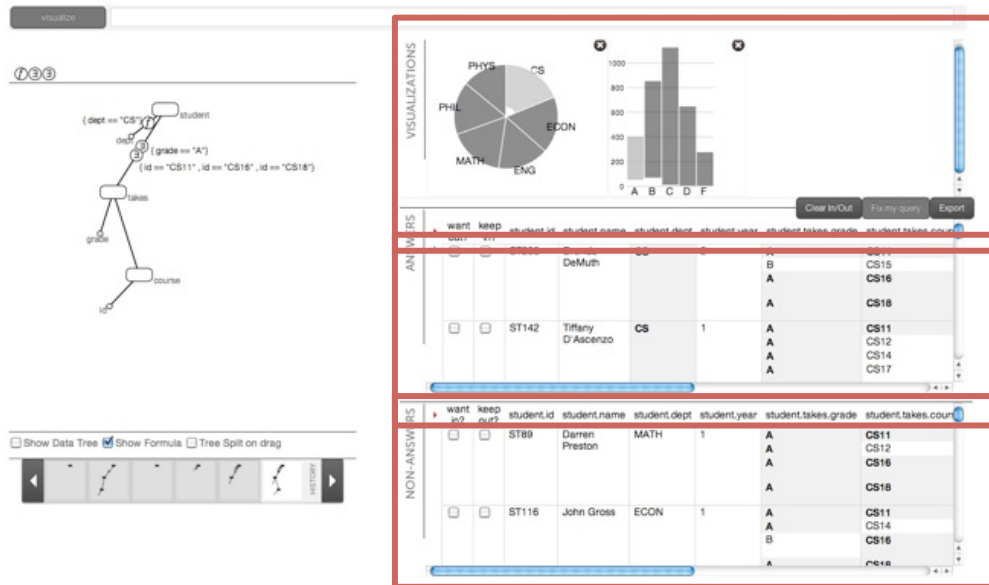
label tuple that prunes as many predicates as possible

graphical query specification

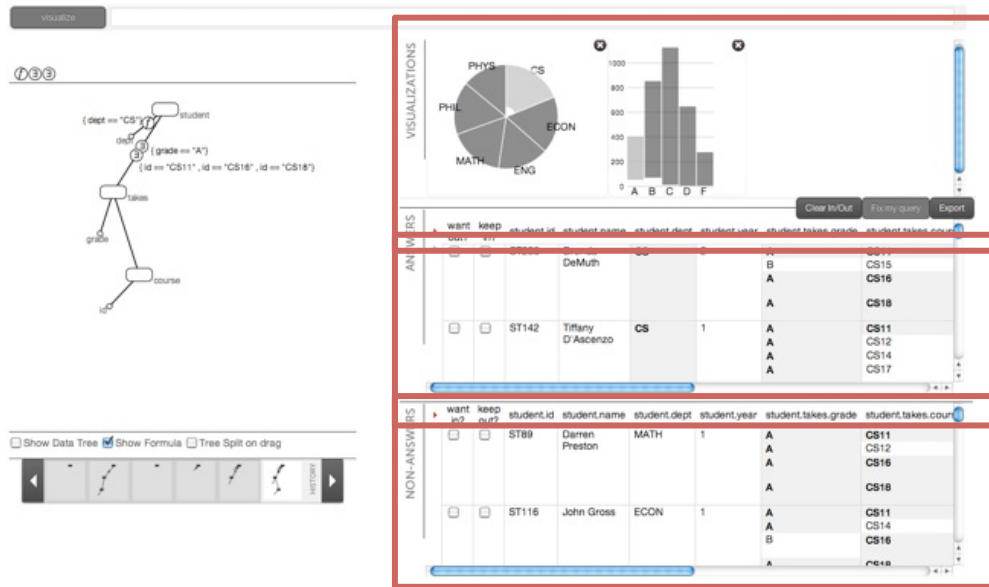
result visualization

answers

non-answers



graphical query specification



result visualization

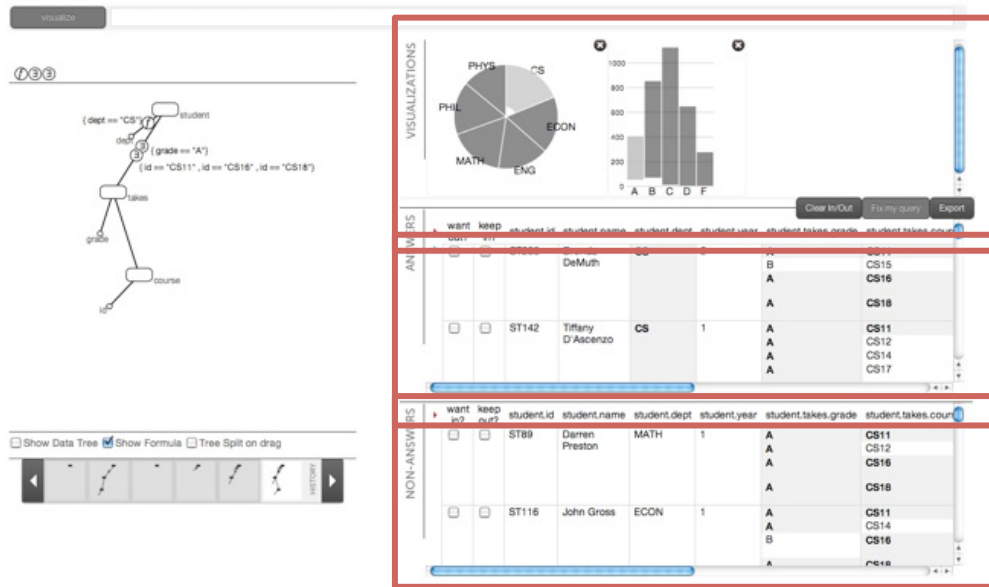
answers

non-answers



semantic query tuning by local syntactic modifications

graphical query specification



result visualization

answers

non-answers

pivot relation



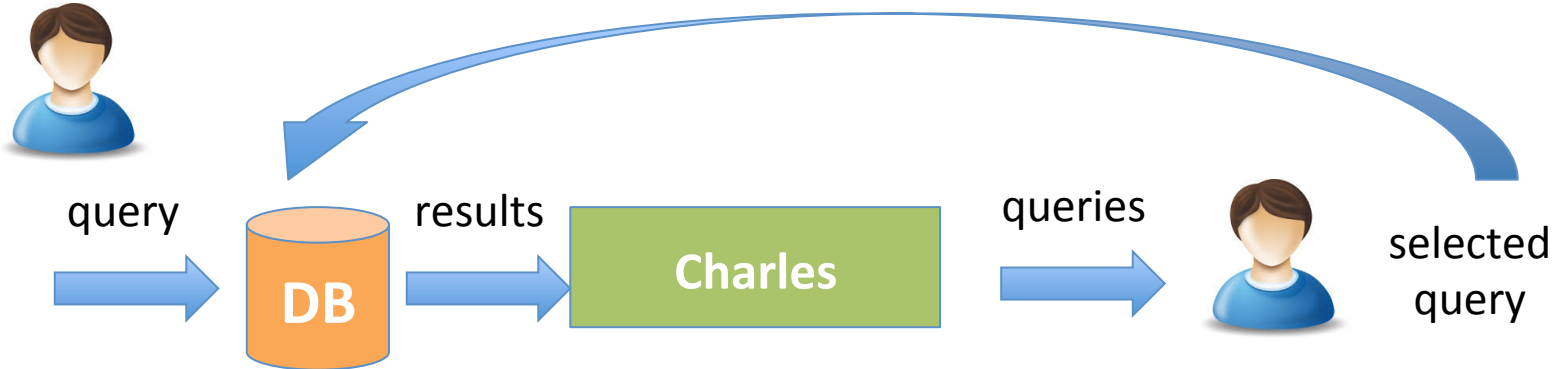
add, remove results



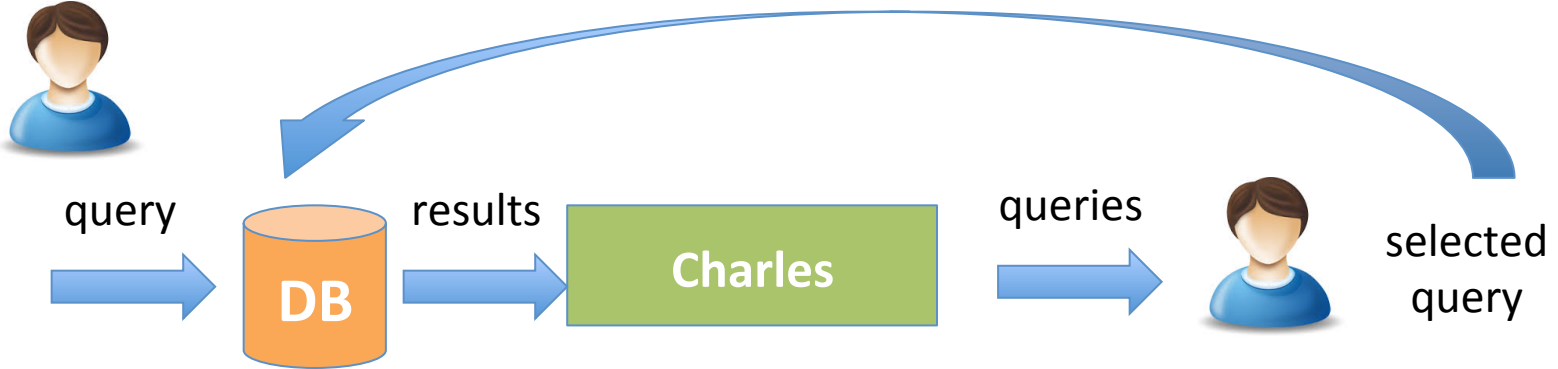
query corrections

search limited to local modifications

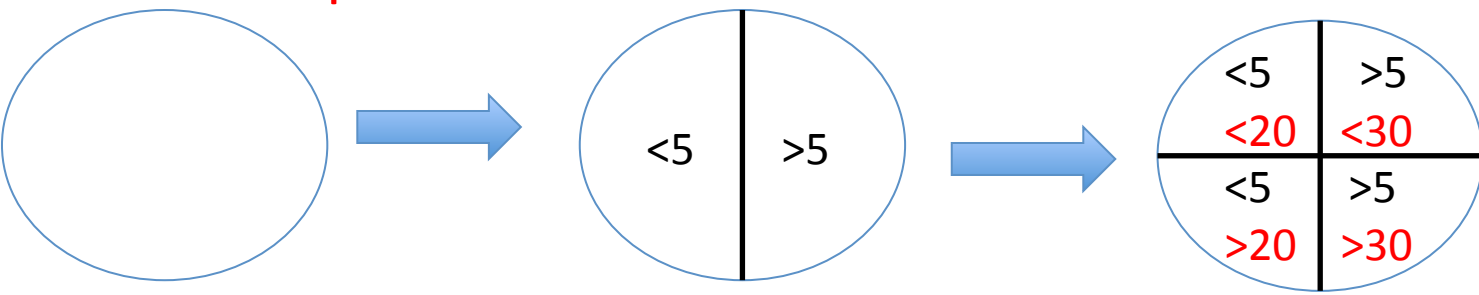
query recommendations



query recommendations



different data partitions



quality: simplicity, breadth, balance

query refinement



conditional query



ranked results by match probability

sensitivity of user predicates

query refinements w/ quality improvement

select species from birds
where color= {red: 80%, blue: 20%}

rank	species
1	Bluebird
2	Blue Jay

attr	sensitivity
color	18.6

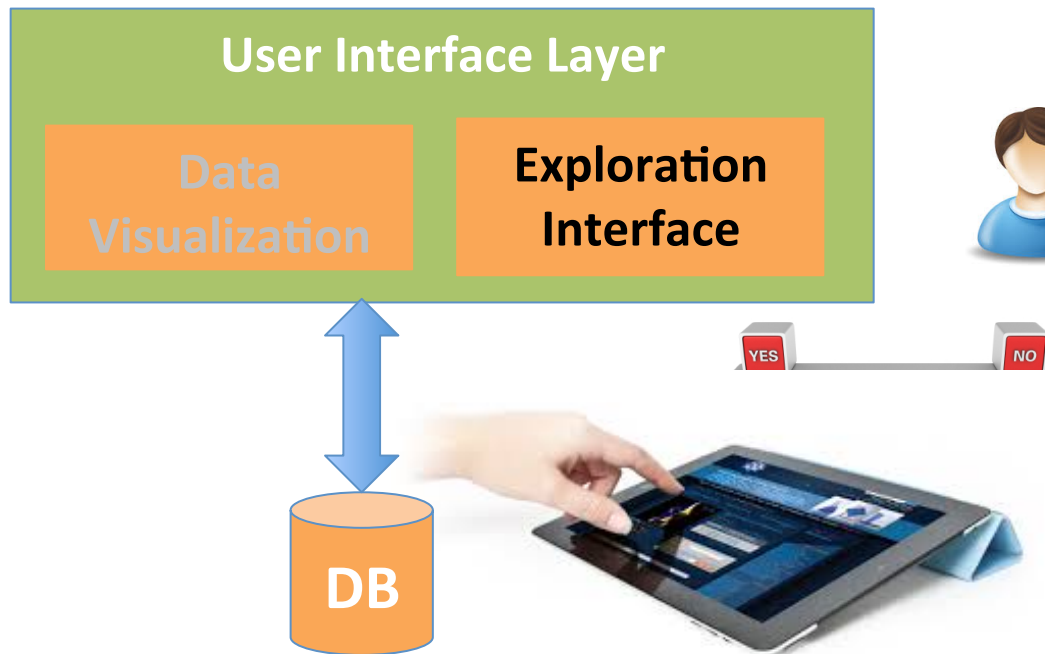
impact on ranking

attr	quality score
size	83.3
legcolor	57.1

remaining attributes

result quality if added in the query

exploration interfaces

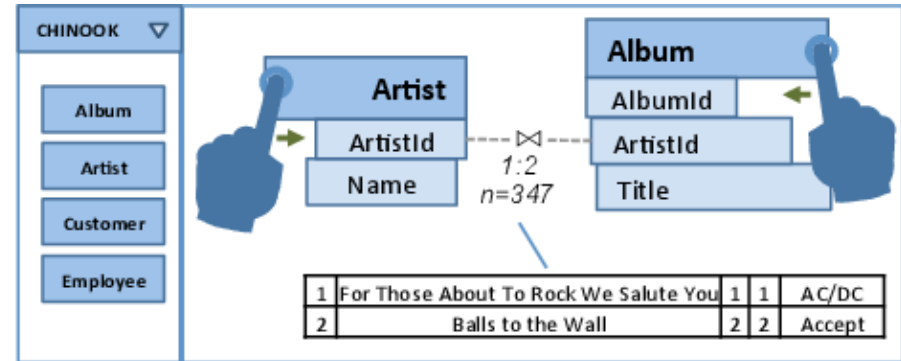


**automatic
exploration**

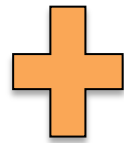
**assisted query
formulation**

novel query interfaces

no-keyboard interfaces



query
context

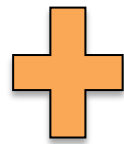


gesture
recognition



query
intend

query
space

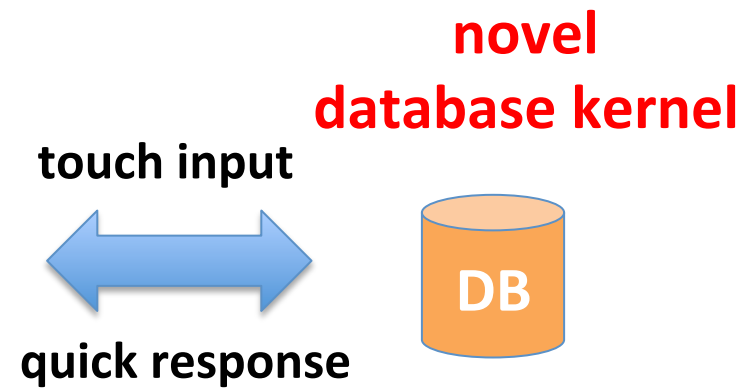


search
pattern



query
template

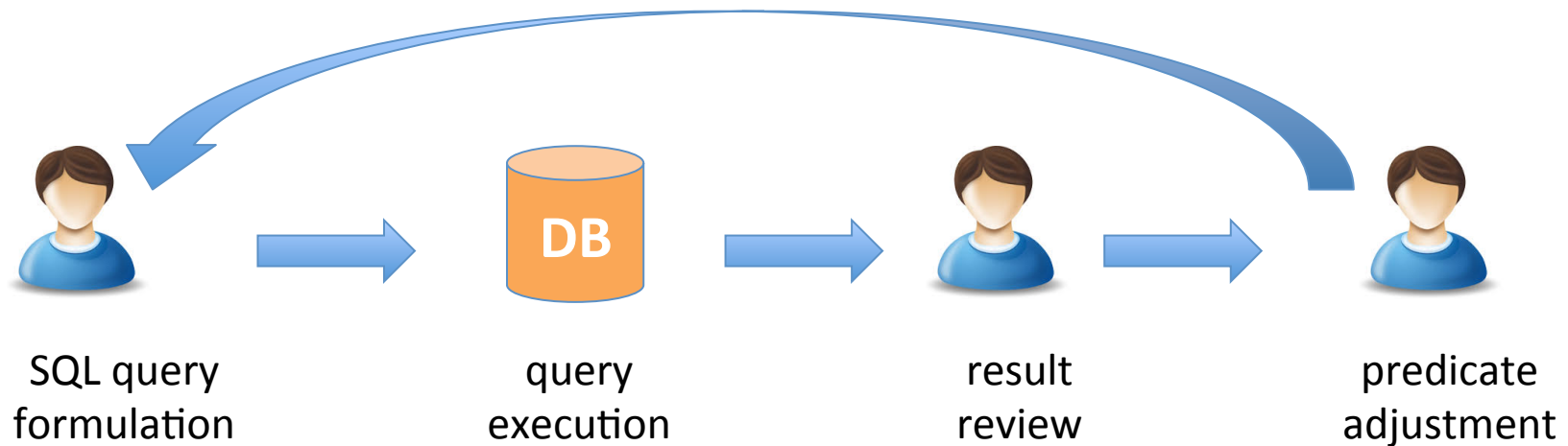
no-keyboard interfaces



Interactive Exploration through
Data Prefetching & Query Approximation

MIDDLEWARE TECHNIQUES

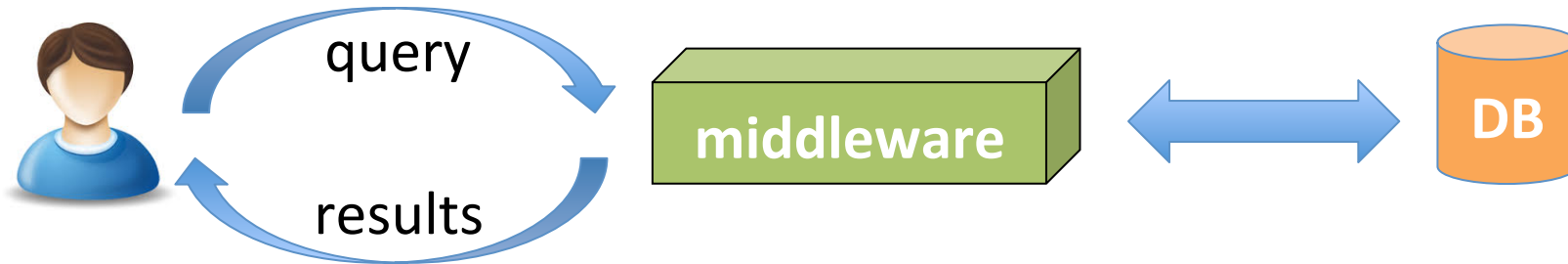
interactive data exploration



ad-hoc, non-optimized, labor-intensive process

interactive: small latency bounds on user wait time

middleware optimizations



query approximation

online processing

sample-based
processing

prefetching

speculative
query execution

result
reuse

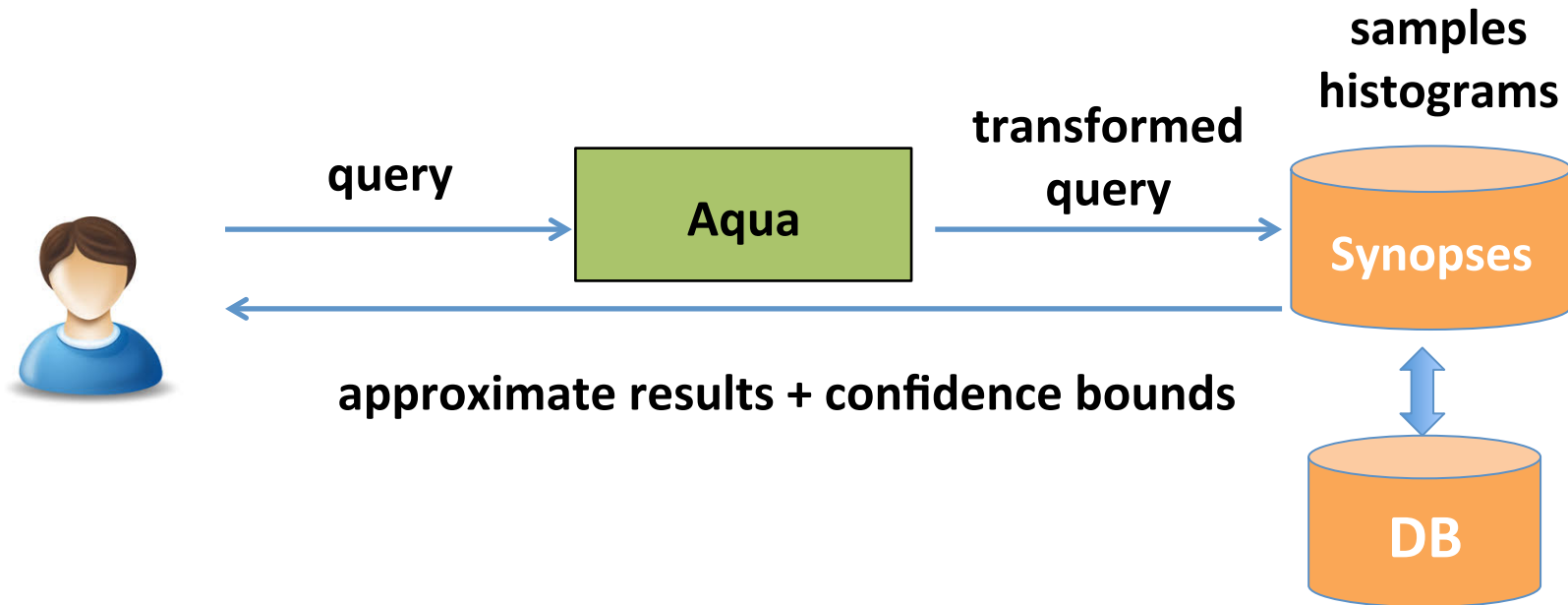
structure-aware
prefetching

sample-based processing



- accuracy vs response times
- sample construction & selection
- error approximation

off-line data synopses

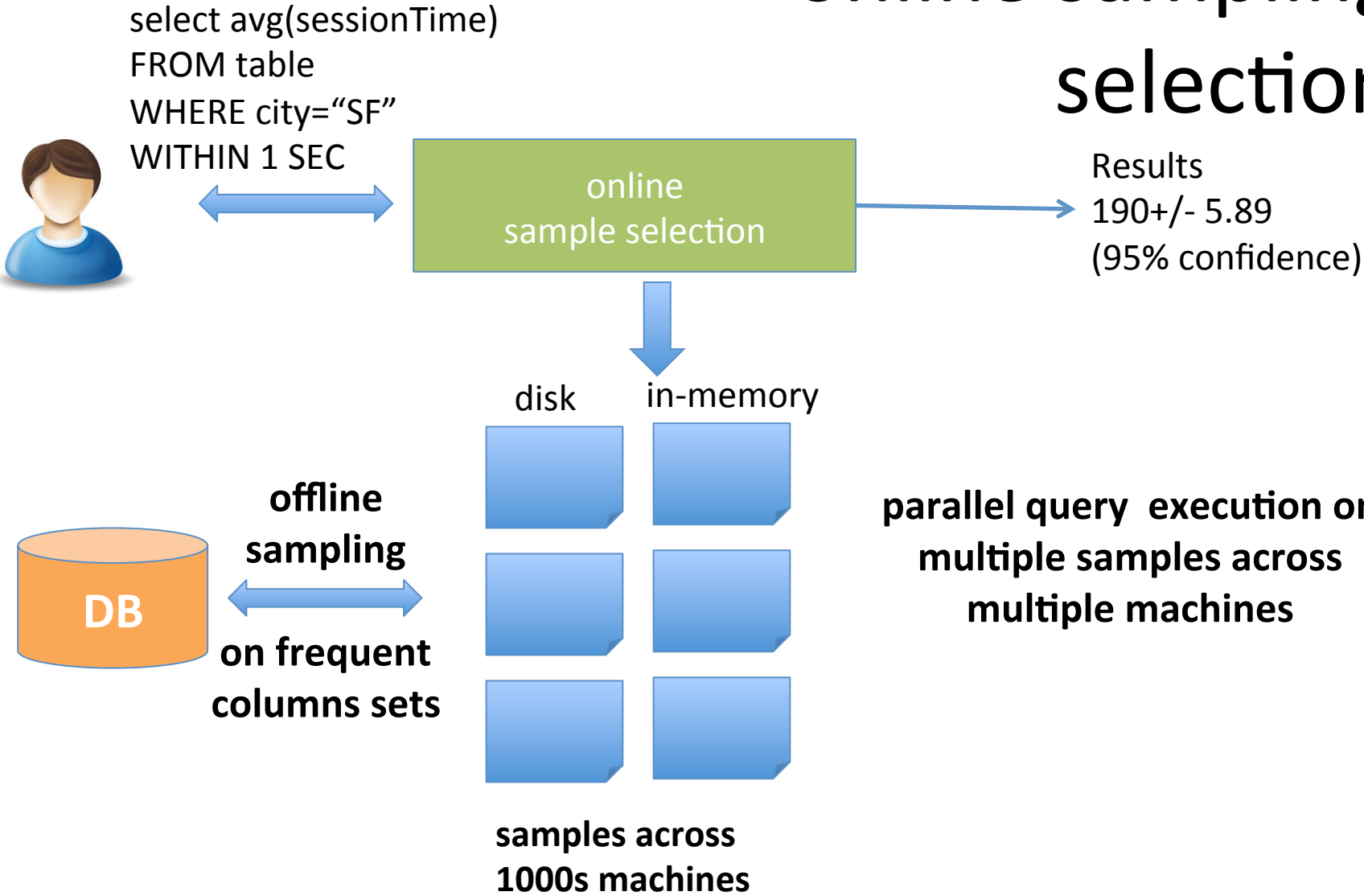


join synopses: sample distinguished joins

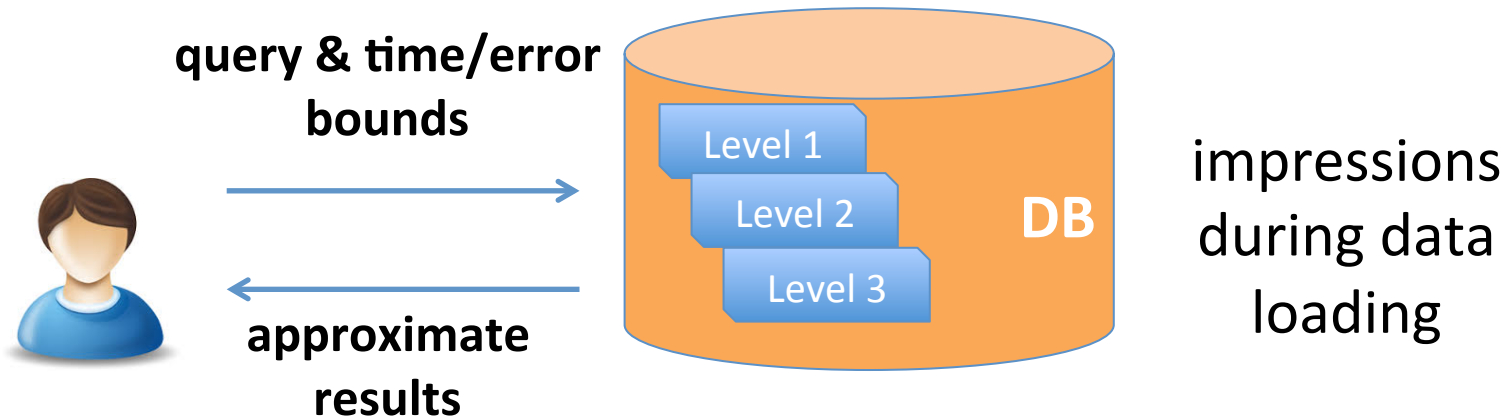
congressional samples: biased sampling for group-by queries

incremental maintenance: equi-depth & compressed histograms

online sampling selection



data impressions

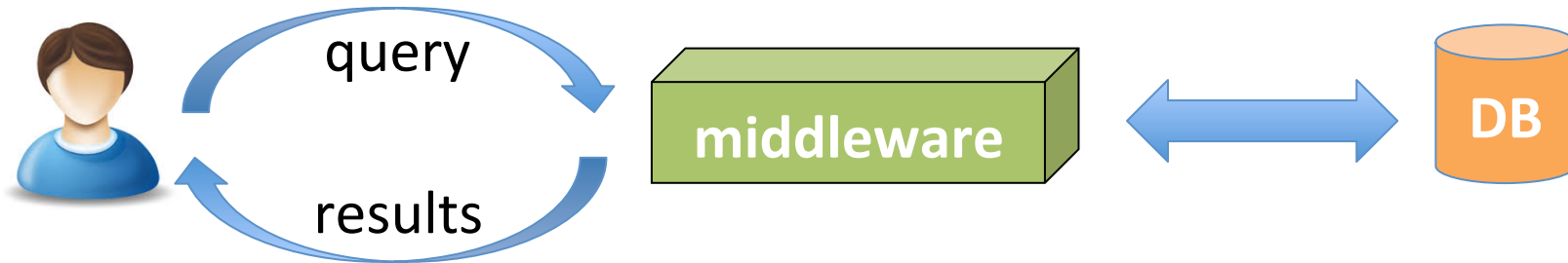


adaptive sampling to exploration focus



multi layer sampling and processing to meet user bounds

middleware optimizations



query approximation

online processing

sample-based processing

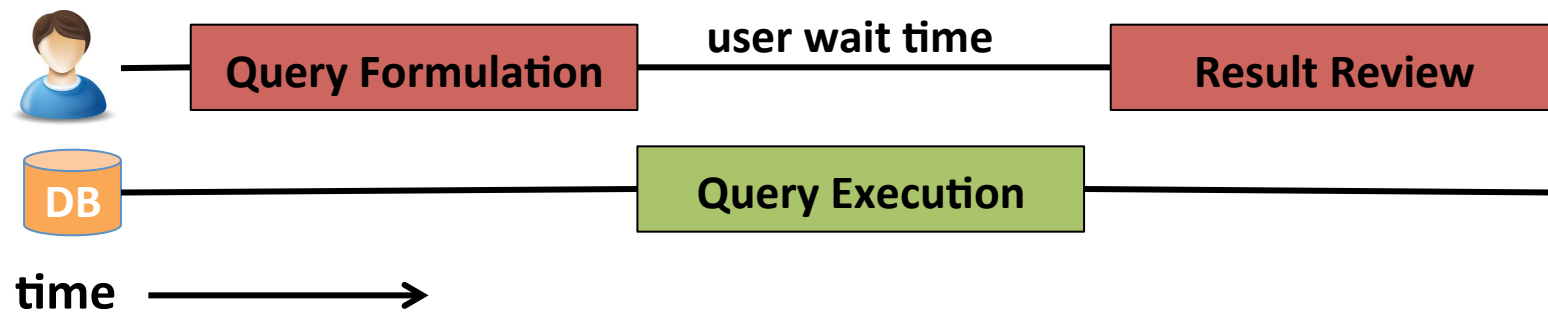
prefetching

speculative query execution

result reuse

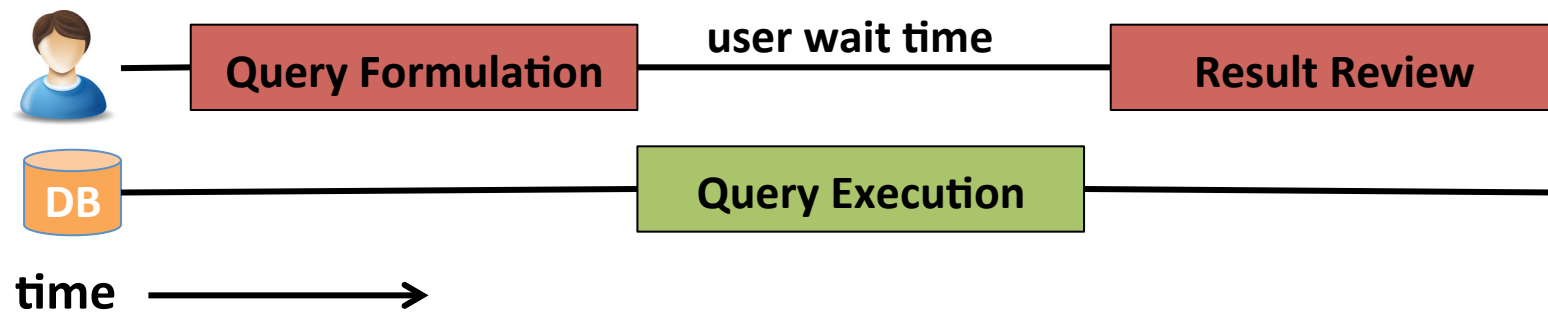
structure-aware prefetching

speculative query execution



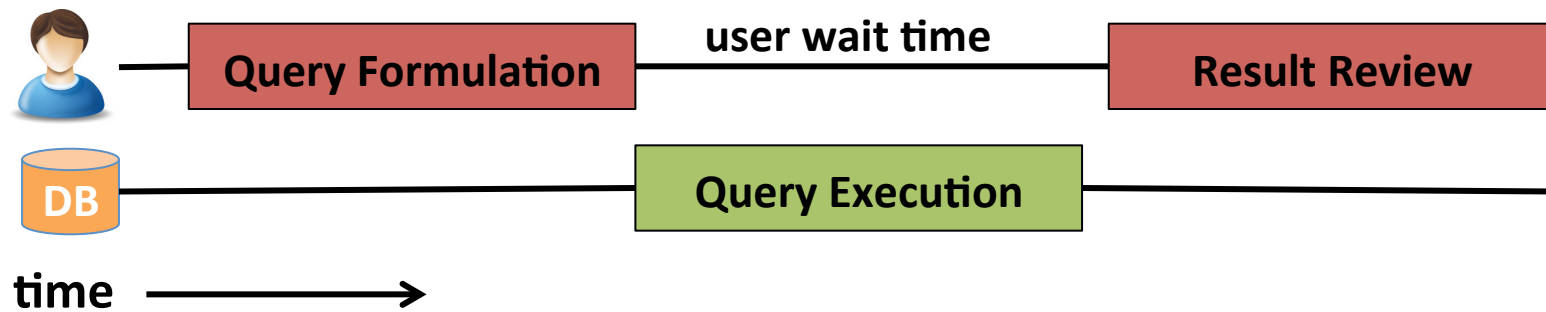
1. predict follow-up queries
2. execute queries
3. cache results

speculative query execution

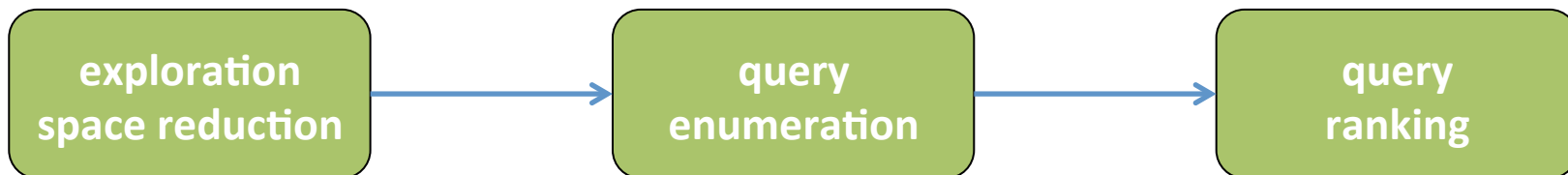


1. predict follow-up queries
2. execute queries
3. cache results

speculative query execution



1. predict follow-up queries
2. execute queries
3. cache results

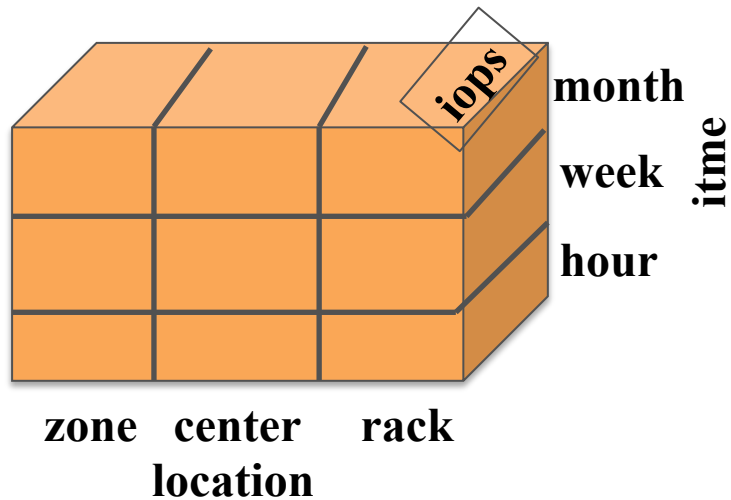


cube exploration

exploration
space reduction

user query

```
SELECT AVG (iops) FROM events  
WHERE month="m1" AND week="w1"  
GROUP BY zone
```

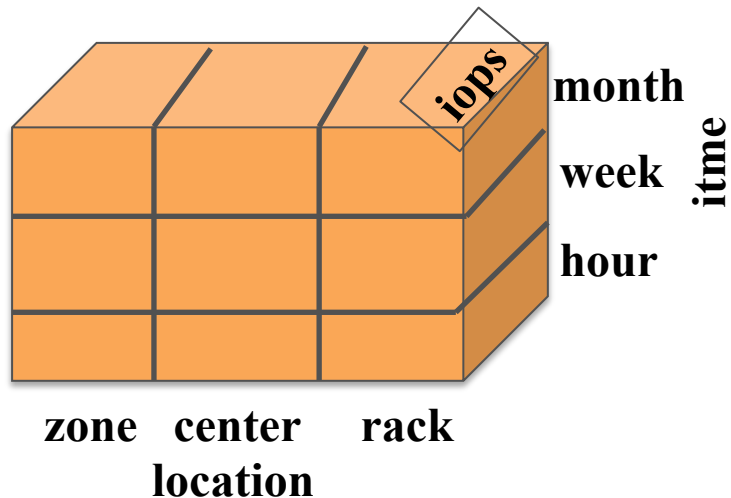


cube exploration

exploration
space reduction

user query

```
SELECT AVG (iops) FROM events  
WHERE month="m1" AND week="w1"  
GROUP BY zone
```



cube exploration operators

WHERE month="m1"

parent

WHERE month="m1"
AND week="w1"
AND hour="h1"

child

WHERE month="m1"
AND week="w2"

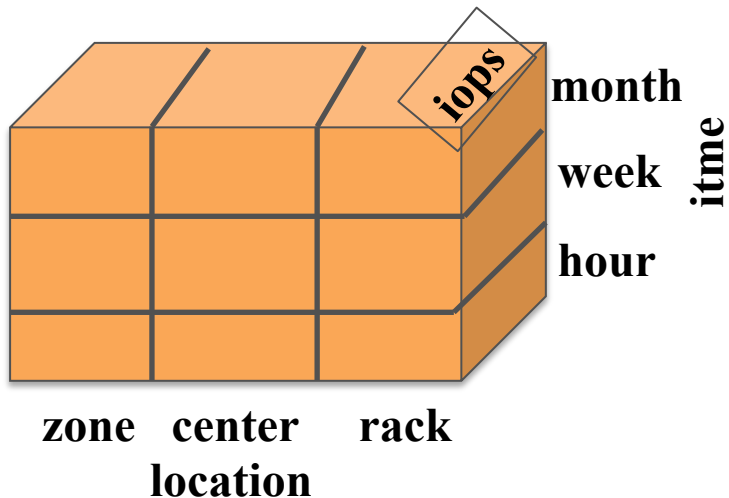
sibling

cube exploration



user query

```
SELECT AVG (iops) FROM events  
WHERE month="m1" AND week="w1"  
GROUP BY zone
```



speculative queries

Q(month="m1")

...

Q(month = "m12")

Q(hour = "h1")

...

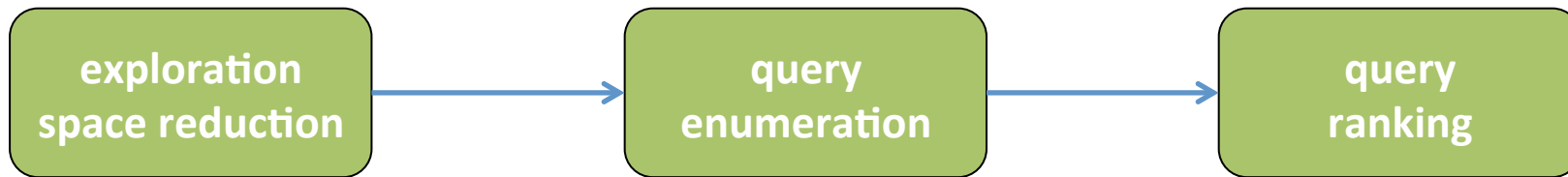
Q(hour = "h24")

Q(week="w2")

...

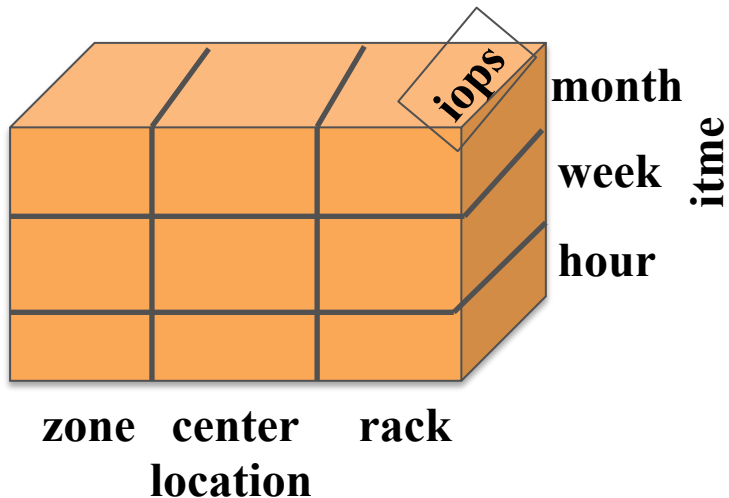
Q(week="w3")

cube exploration



user query

```
SELECT AVG (iops) FROM events  
WHERE month="m1" AND week="w1"  
GROUP BY zone
```



speculative queries

Q(month="m1")

...

Q(month = "m12")

Q(hour = "h1")

...

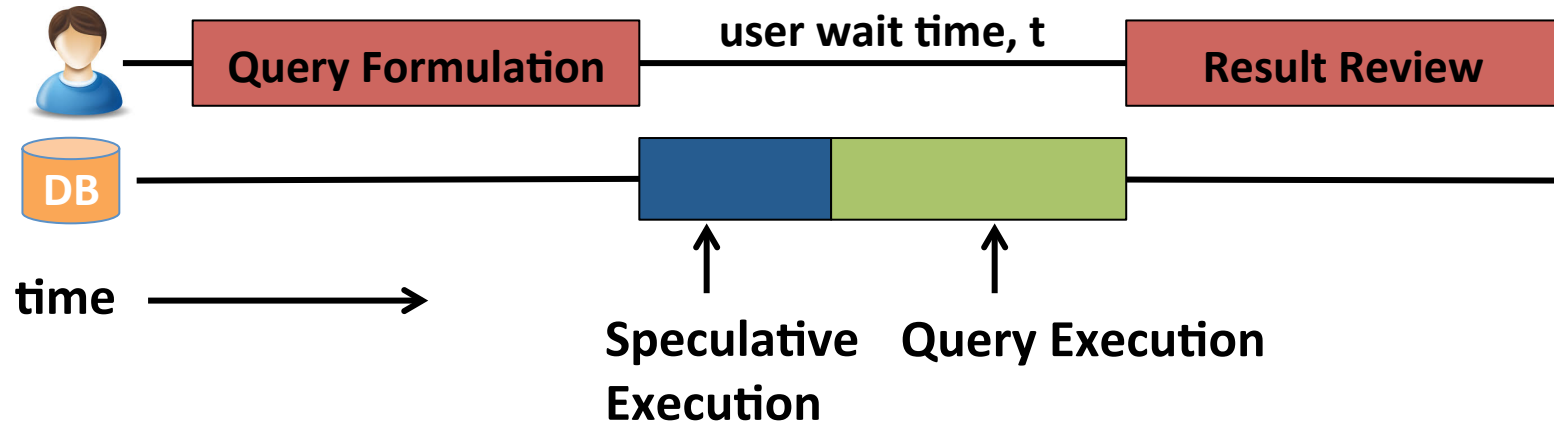
Q(hour = "h24")

Q(week="w2")

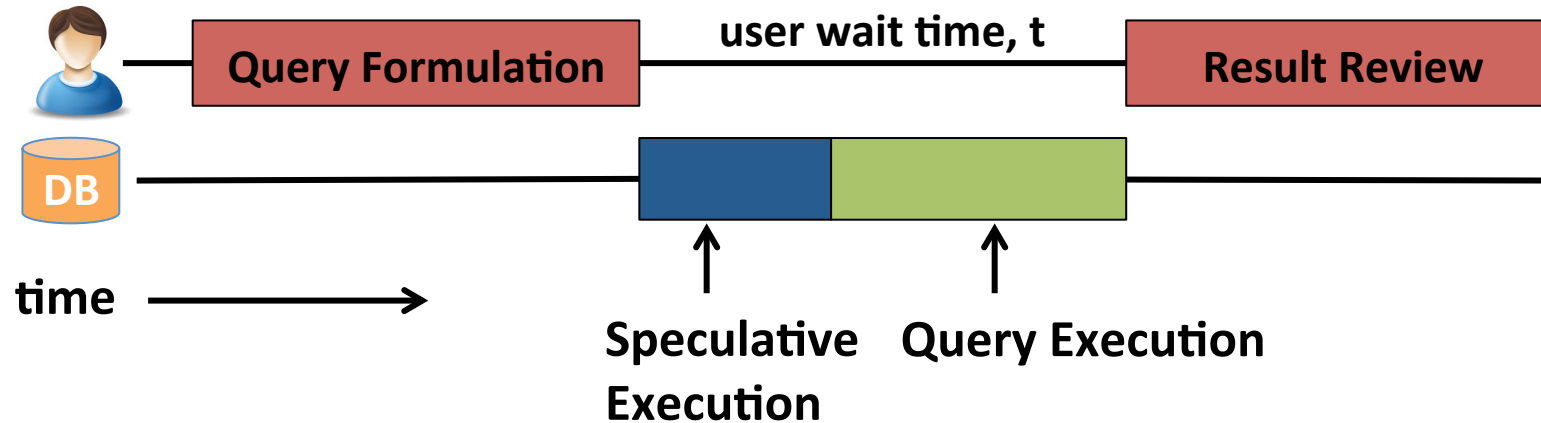
...

Q(week="w3")

cube exploration



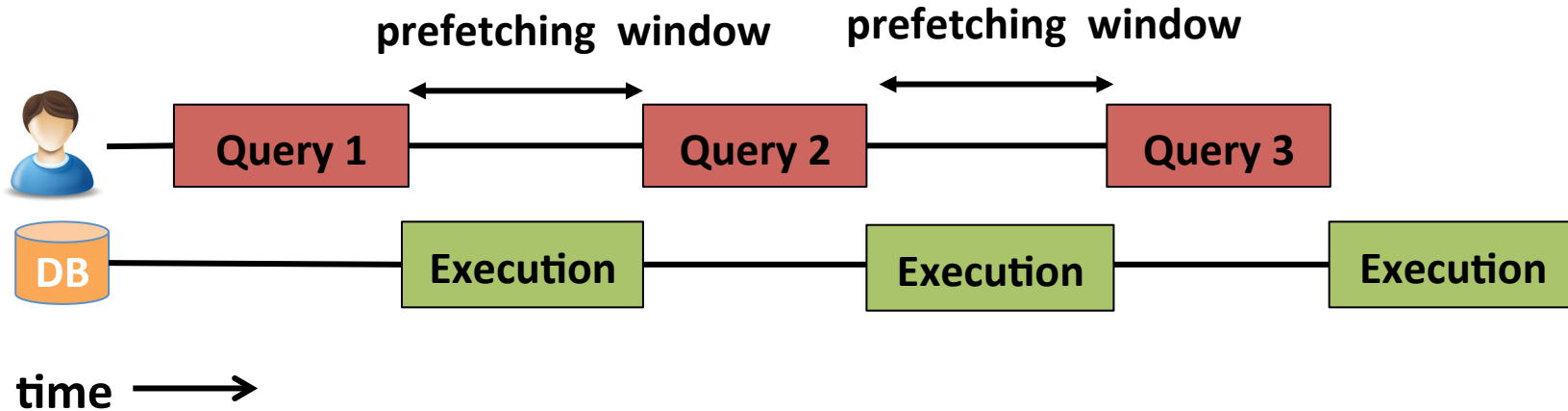
cube exploration



QUERY	Probability	Exec Time
Q ₁	0.3	22
Q ₂	0.25	20
Q ₃	0.25	35
Q ₄	0.15	70
Q ₅	0.05	35

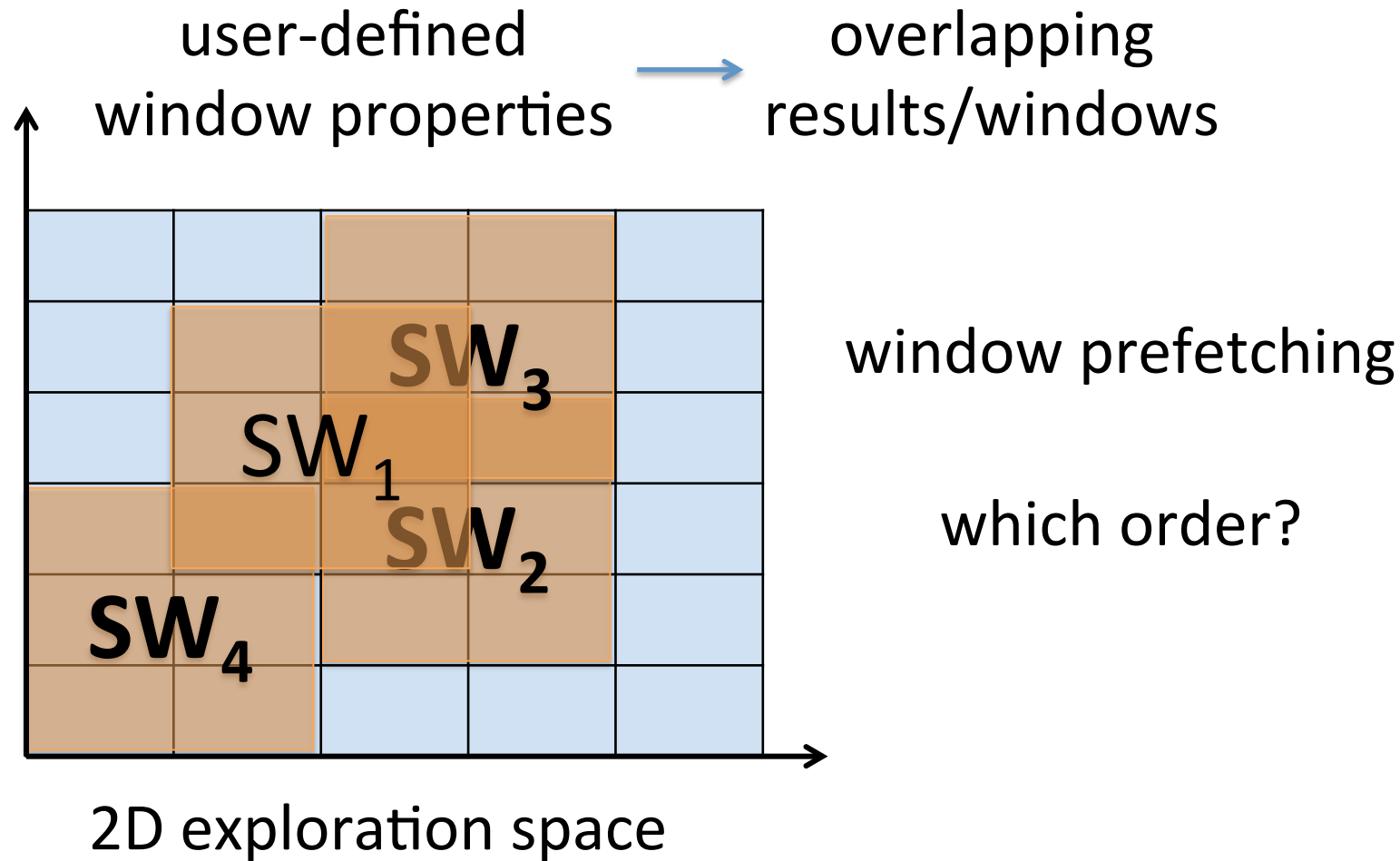
maximize query probability
total speculation time < t

result reuse

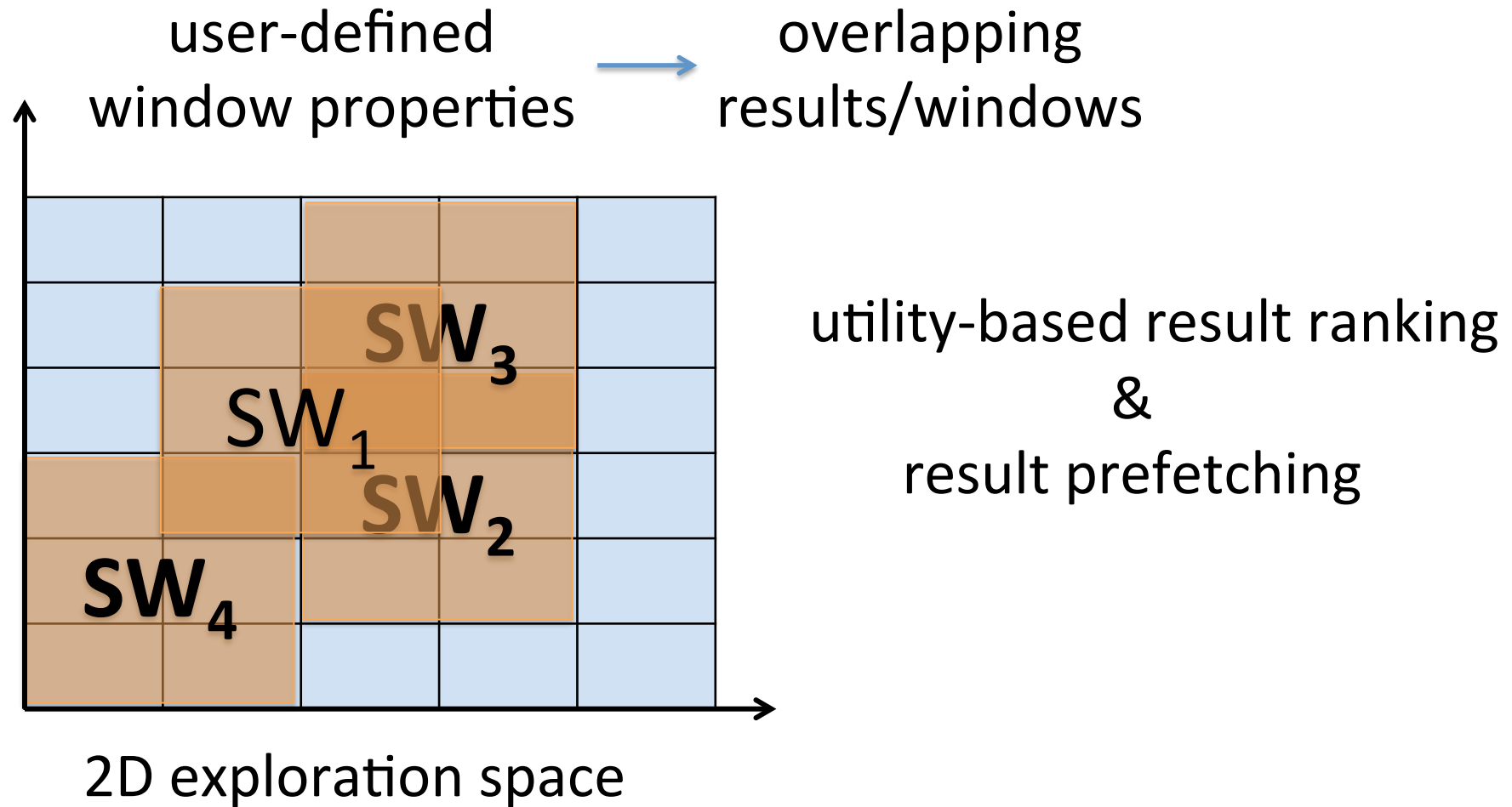


- identify (likely) **overlapping** results
- **cache** them
- reduce query execution time (user wait time)

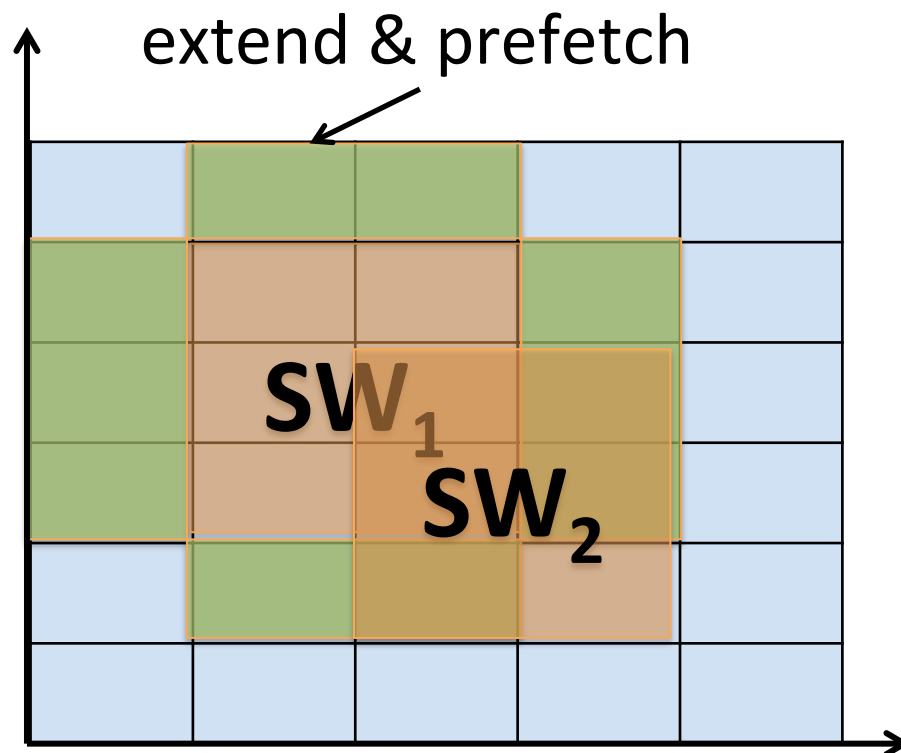
semantic windows



semantic windows



semantic windows



online performance
vs
query completion time



adjust prefetching size
to output progress



query



diversified
results



k representative tuples with
max total pairwise distance

data diversification

data diversification



T_1	$d(T_1, T_3)$
T_2	$d(T_2, T_3)$
T_3	← random tuple
T_4	$d(T_4, T_3)$
T_5	$d(T_5, T_3)$

data diversification



query



diversified
results



k representative tuples with
max total pairwise distance

Query
Output

Max Diversified Set Search



Diversified
Output k= 3

T_1
T_2
T_3
T_4
T_5

$d(T_1, T_3)$

$d(T_2, T_3)$

← random
tuple

$d(T_4, T_3)$

$d(T_5, T_3)$

T_1
T_2
T_3
T_4
T_5

$d(T_2, T_1) + d(T_2, T_3)$

$d(T_4, T_1) + d(T_4, T_3)$

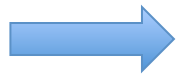
$d(T_5, T_1) + d(T_5, T_3)$

data

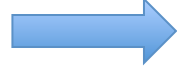
diversification



query



diversified results



k representative tuples with max total pairwise distance

Query Output

Max Diversified Set Search



Diversified Output k= 3

T_1
T_2
T_3
T_4
T_5

$d(T_1, T_3)$

$d(T_2, T_3)$

← random tuple

$d(T_4, T_3)$

$d(T_5, T_3)$

T_1
T_2
T_3
T_4
T_5

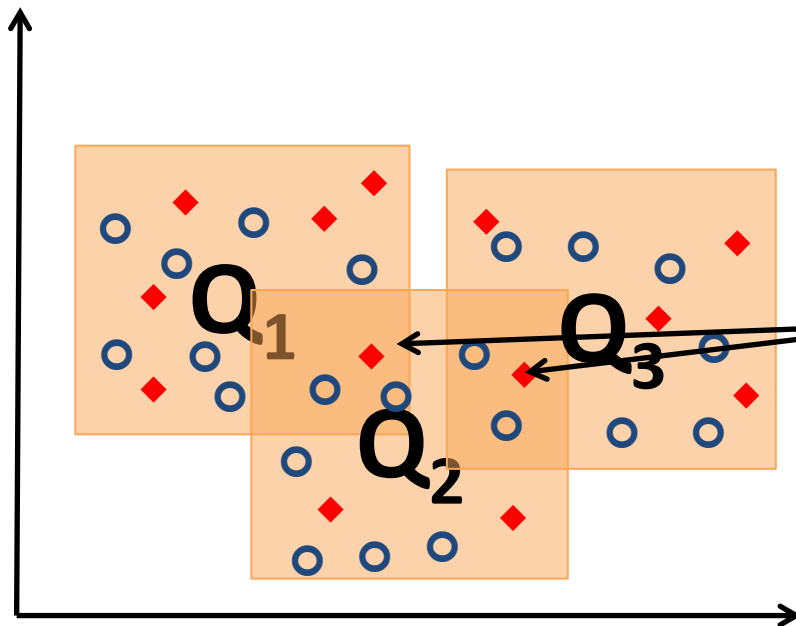
$d(T_2, T_1) + d(T_2, T_1)$

$d(T_4, T_1) + d(T_4, T_3)$

$d(T_5, T_1) + d(T_5, T_3)$

T_1
T_2
T_3
T_4
T_5

interactive data diversification

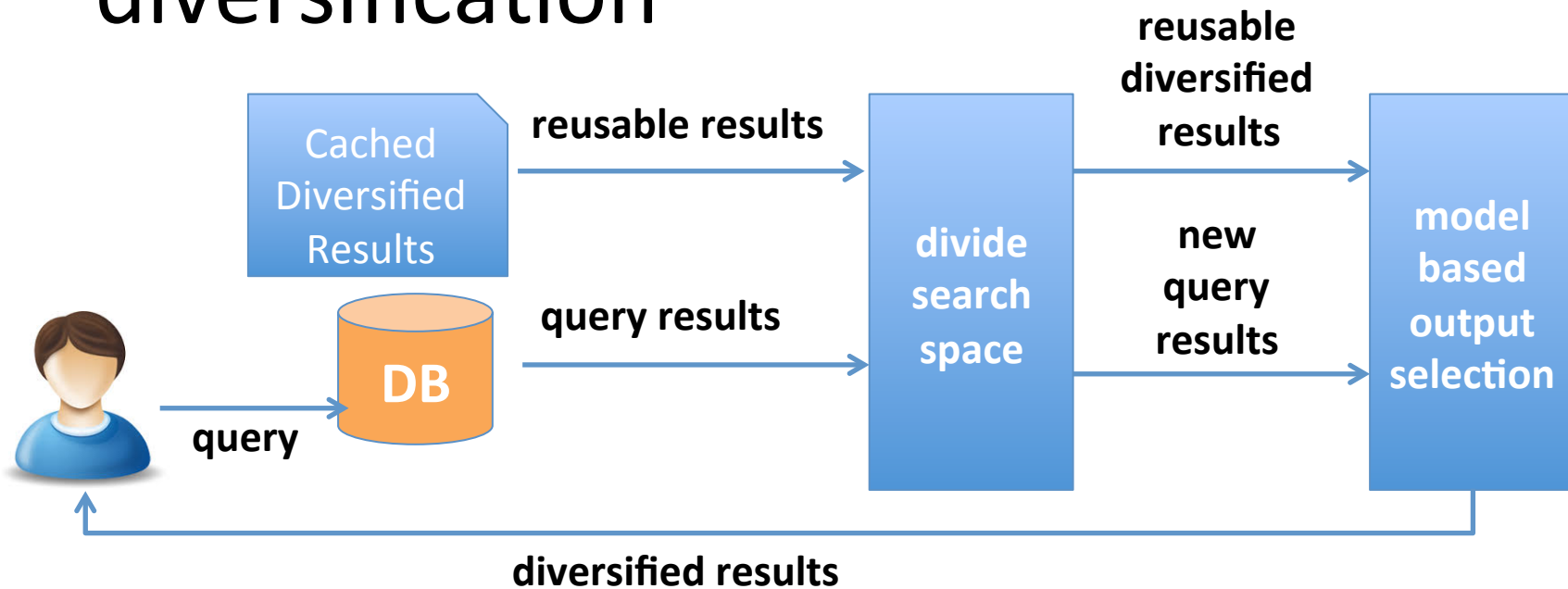


overlapping
diversified results

long **Time-To-Insight**

- cache **diversified** results and use most promising
- regression model **predicts** max diversification of a set

interactive data diversification



- search space **pruning** through regression model
- best/first fit search for max total diversification among cached and new results

structure-aware prefetching

- prefetching for interactive spatial query sequences
- model structures of past spatial queries in **graph**
- identify guiding structure in past two queries : **iterative pruning**
- cache the predicted next location

