

 POLITECNICO DI MILANO

Dipartimento di
Elettronica e Informazione

Search Computing: Query Execution Framework

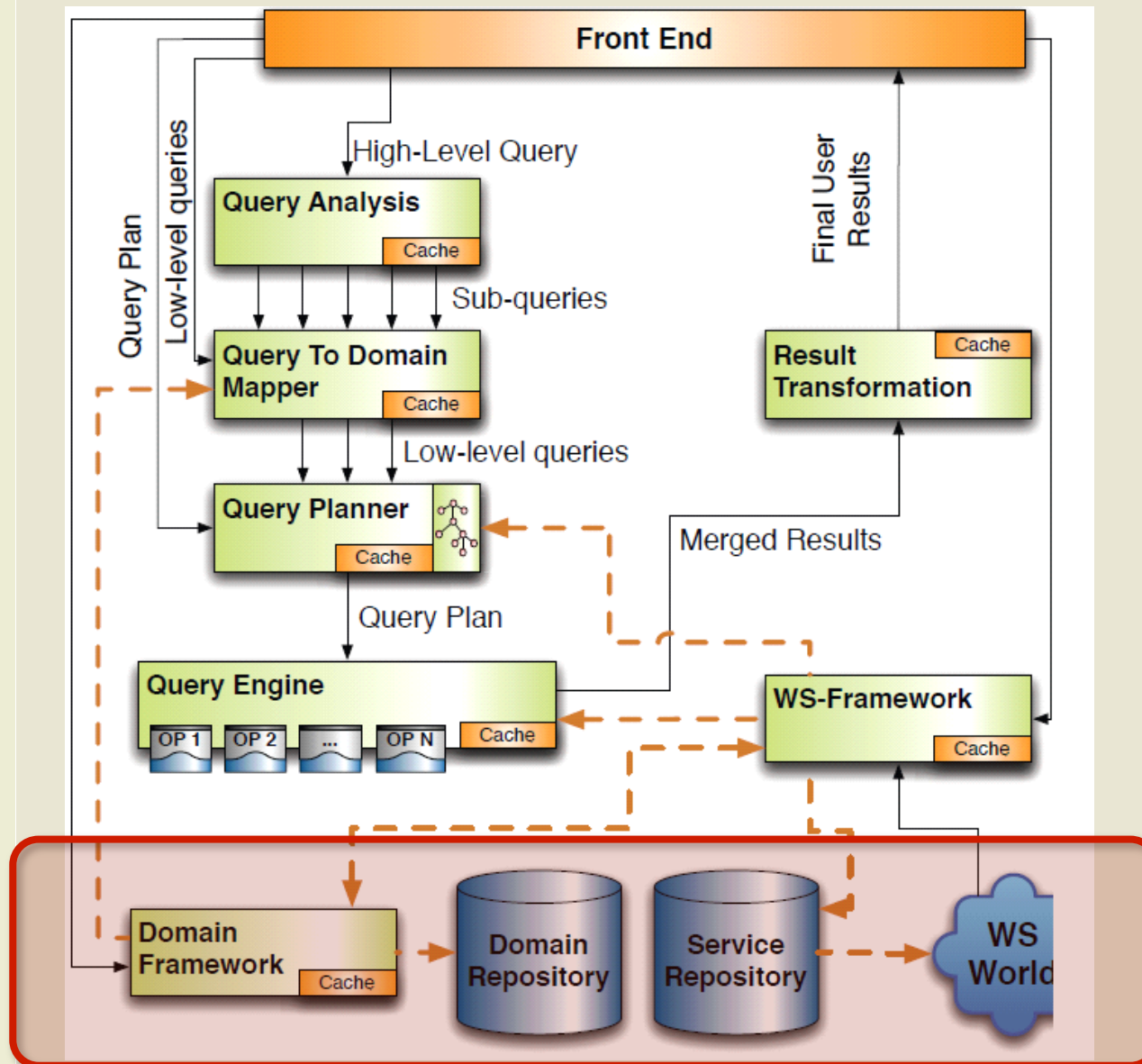
Stefano Ceri

... and the SeCo Project Team

Adnan Abid, Mamoun Abu Helu, Davide Barbieri, **Daniele Braga**, Marco Brambilla, Alessandro Bozzon, Alessandro Campi, Davide Chicco, Emanuele Della Valle, Piero Fraternali, Nicola Gatti, Giorgio Ghisalberti, Davide Martinenghi, Marco Masseroli, Maristella Matera, Chiara Pasini, Silvia Quarteroni, Marco Tagliasacchi, Luca Tettamanti, **Salvatore Vadacca**, Serge Zagorac

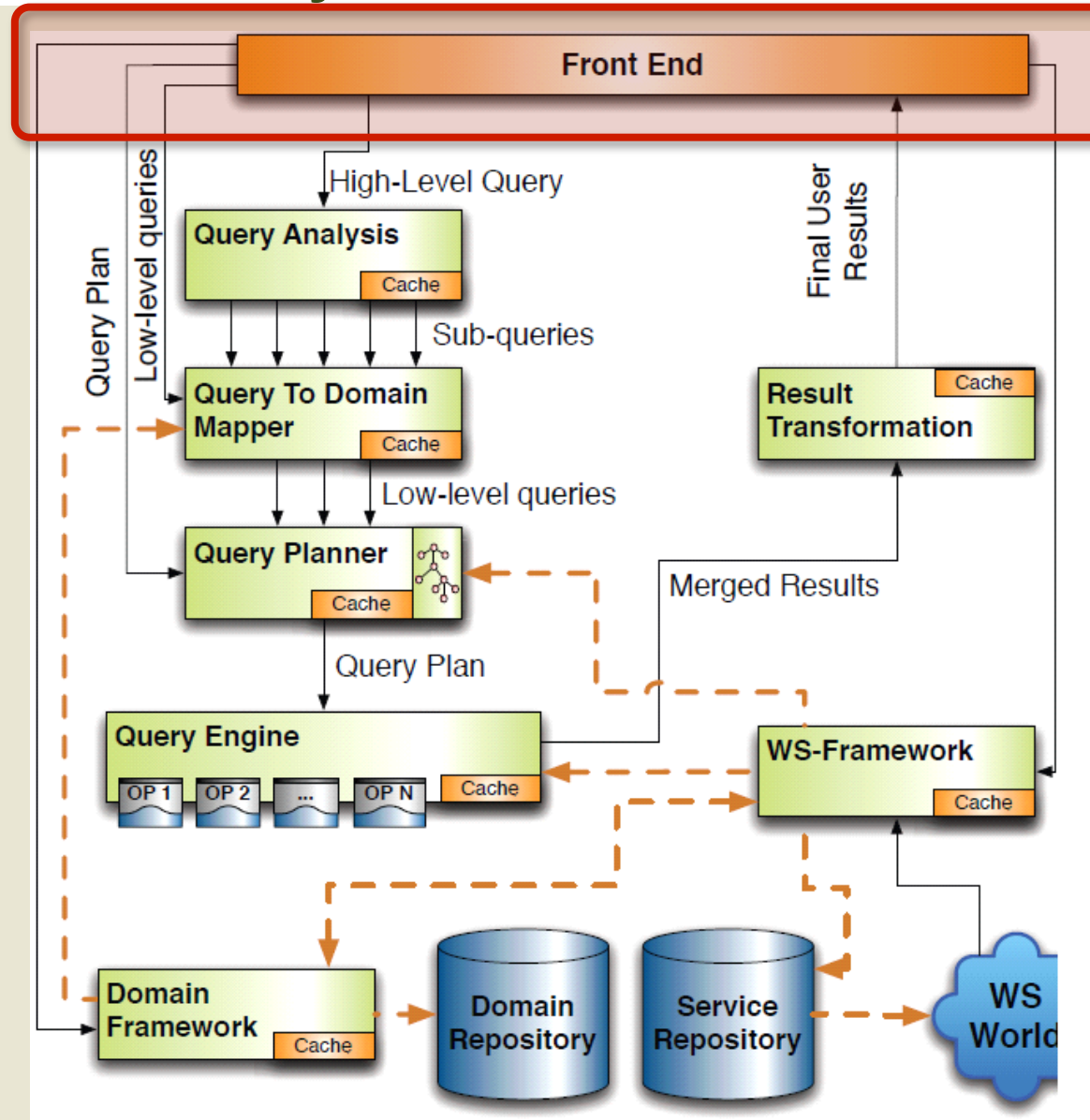
That was the Service Framework

2



On Monday, focus on the user Interface

3

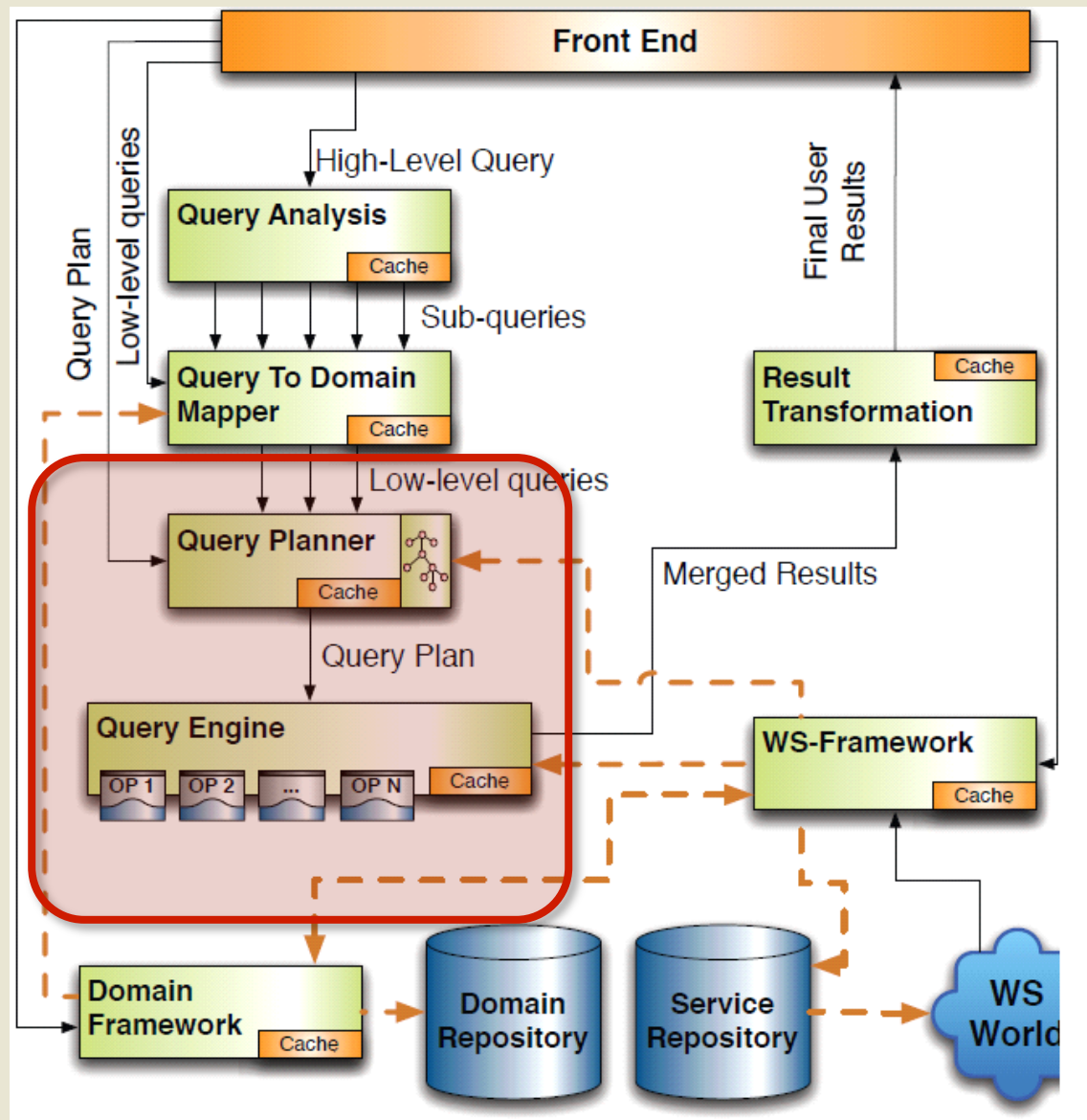


Liquid Queries

- **Intuitive formulation**
- **Exploratory search:**
search as a **process**
incremental formulation
- **Multi-dimensional visualization:**
multiple views for results

Now, focus on: SeCo queries and the Query Engine

4



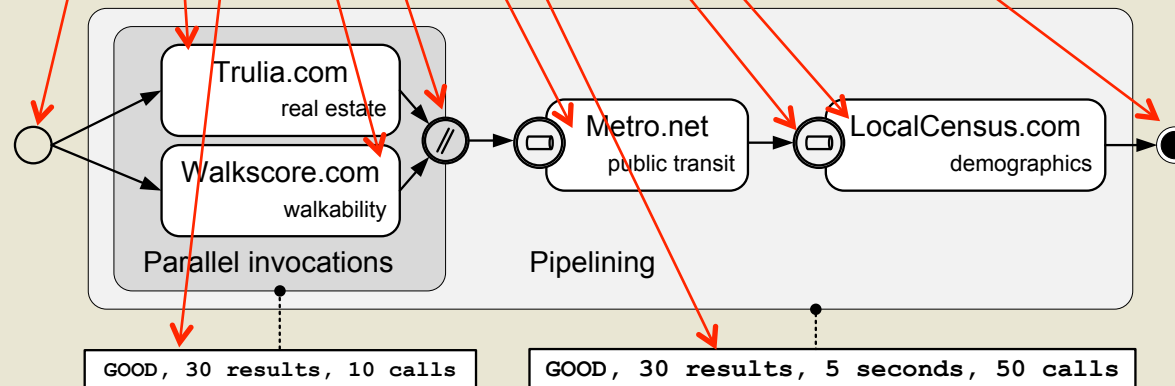
Query execution

- **Query planner:**
translates queries into query plans
- **Query engine:**
executes query plans according to strategies
 - Top-k
 - Optimistic (good-k)
 - Cost-driven

Search Computing = service composition “on demand” 5

- Composition **abstractions** should emphasize few aspects:
 - service invocations
 - fundamental operations (parallel invocations, joins, pipelining, ...)
 - global constraints on execution
- Data composition should be search-driven
 - aimed at producing **few top results** very fast

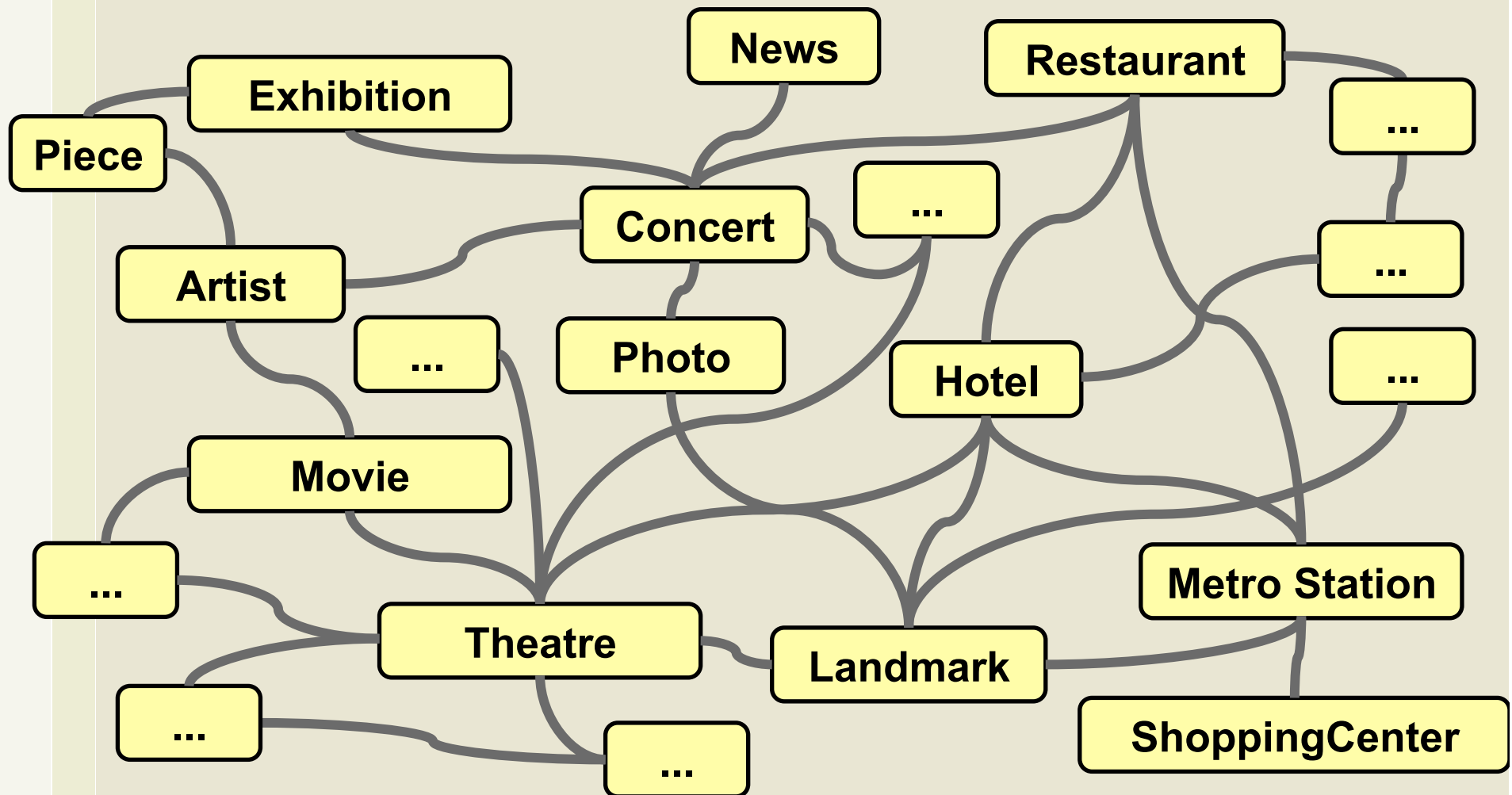
A house in a walkable area, close to public transportation and located in a pleasant neighborhood



SeCo queries

6

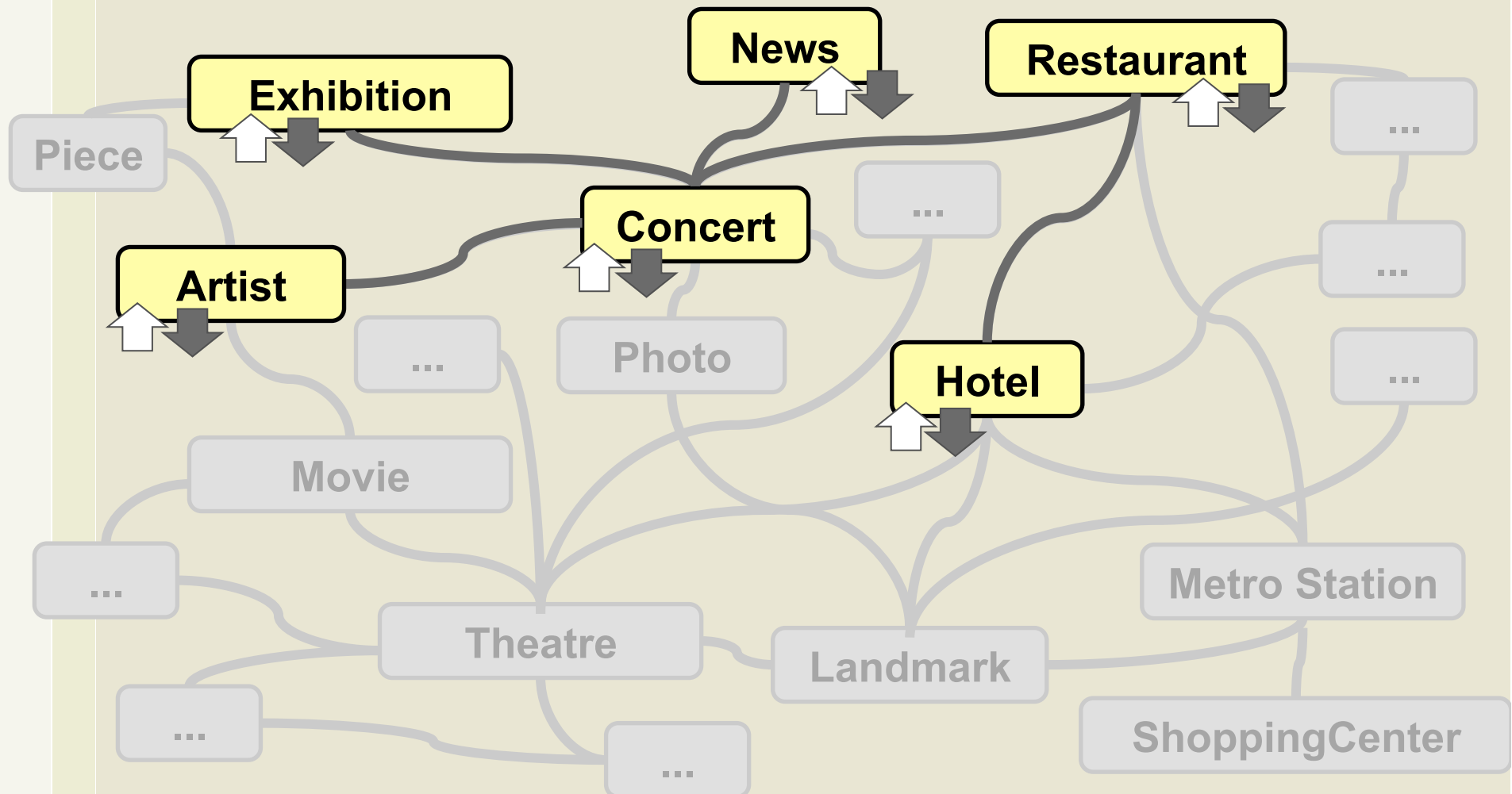
- Navigation of a graph of connected data sources



Query formulation

7

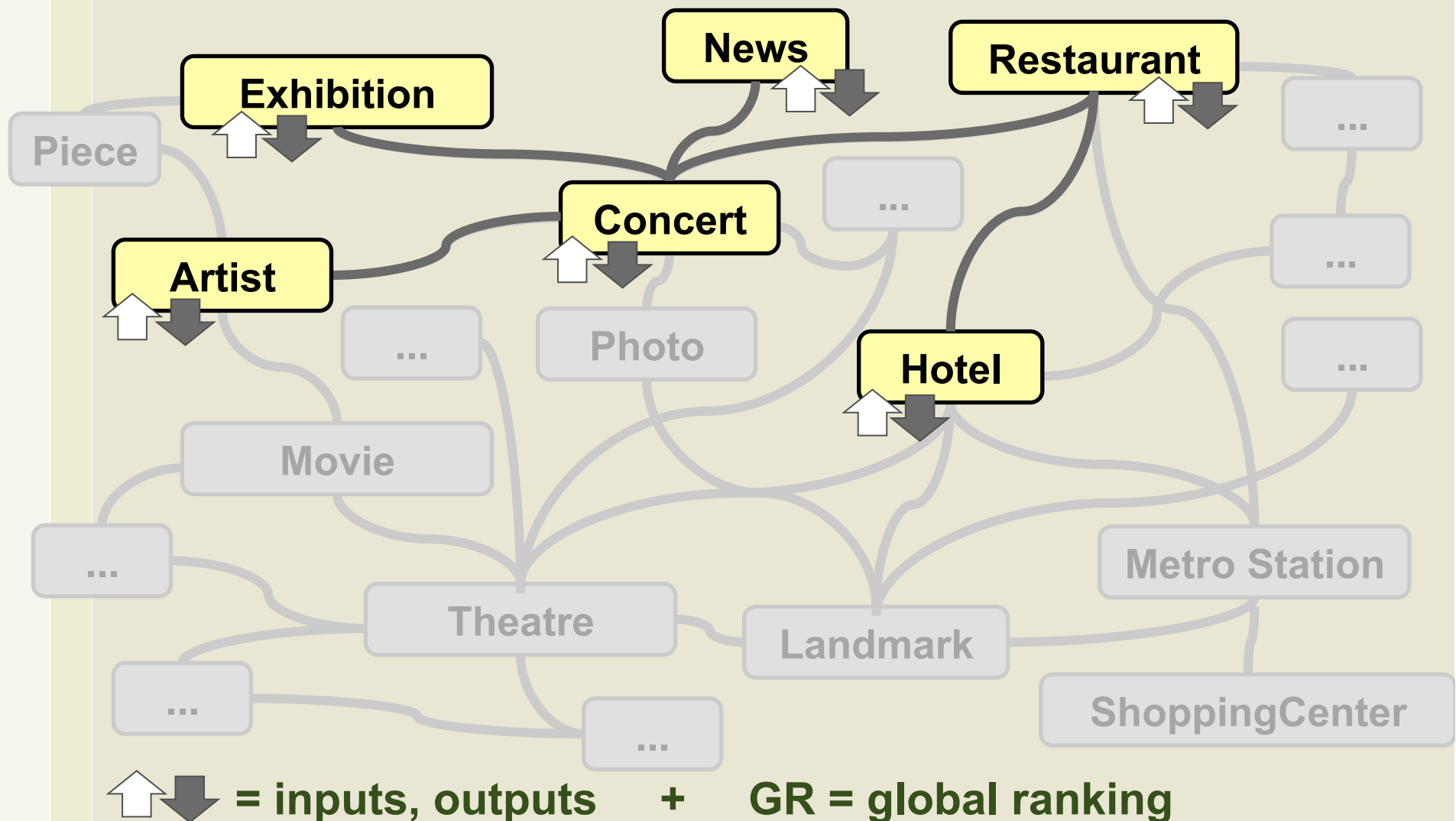
Consists in subsetting and parametrizing the graph



↑ ↓ = inputs, outputs + GR = global ranking

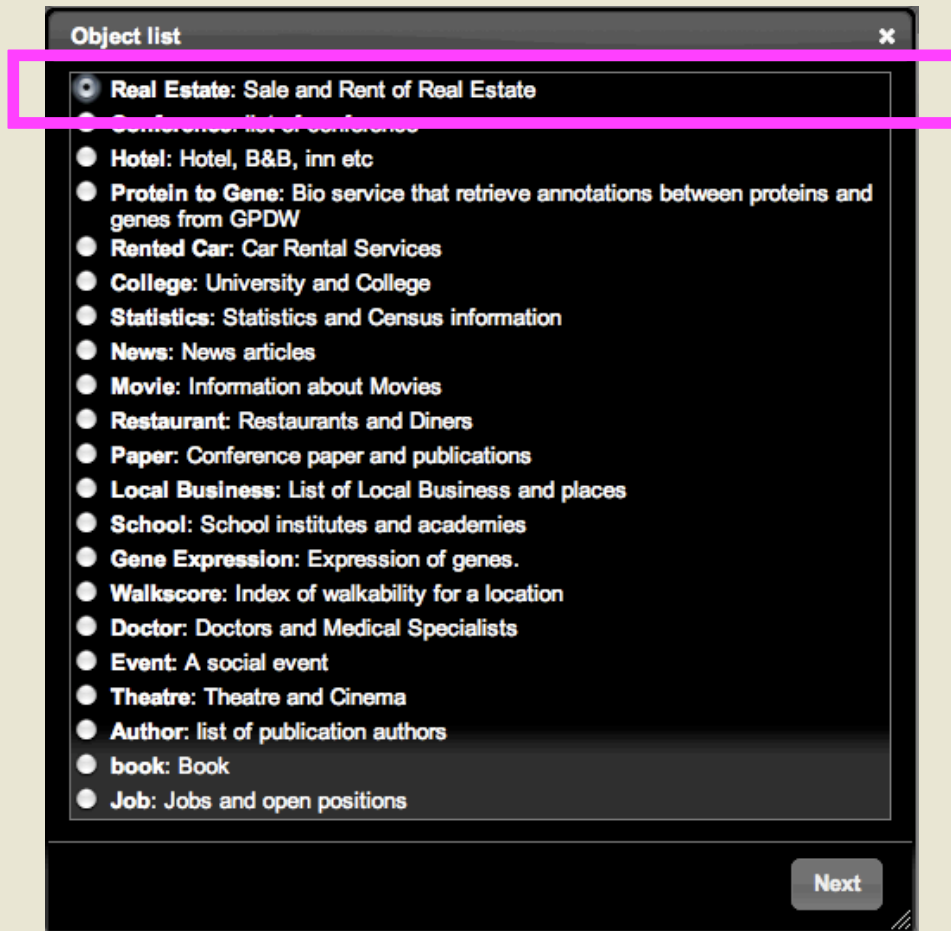
Query formulation

Consists in subsetting and parametrizing the graph...



Exploration of the Service Space

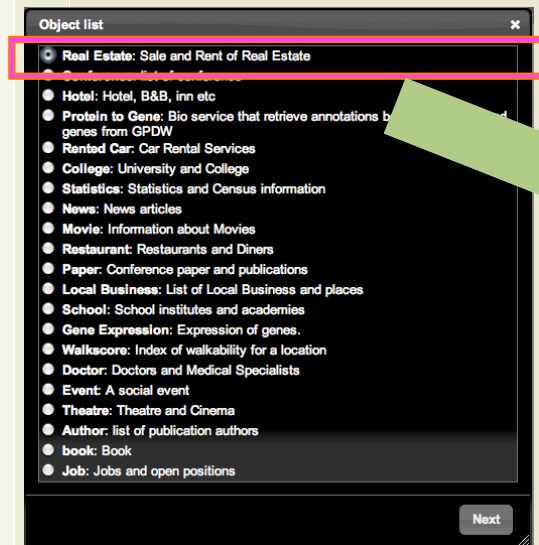
9



Entity Selection (service mart)

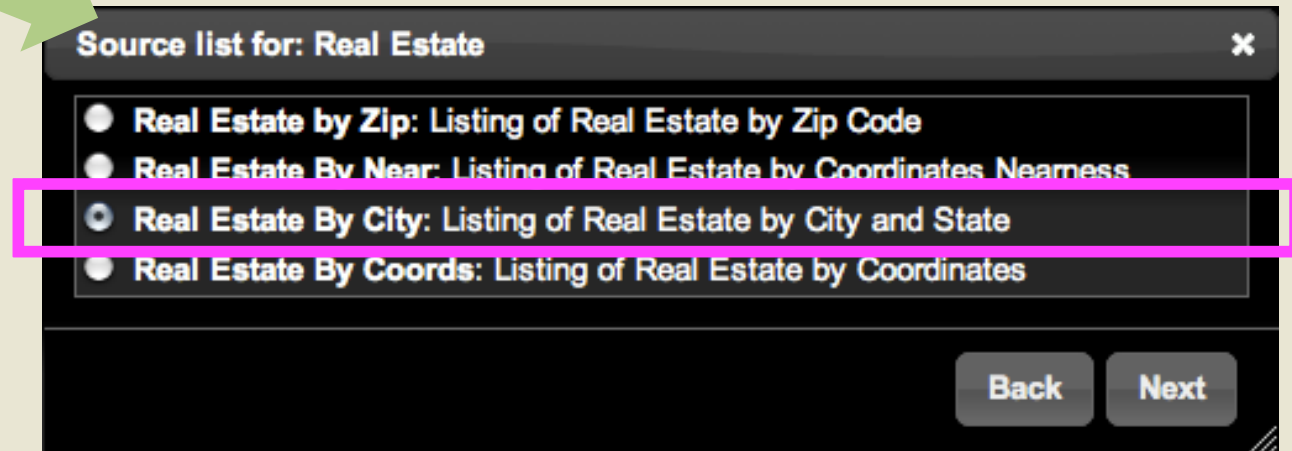
Exploration of the Service Space

10



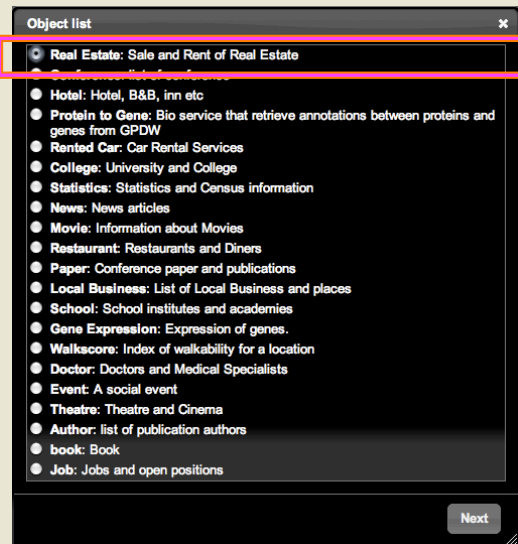
Entity Selection

Service Selection
(access pattern)

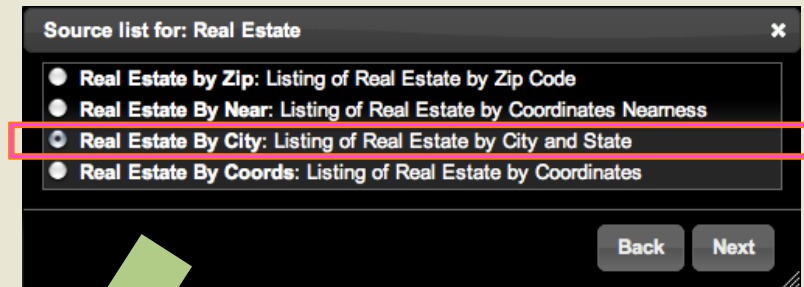


Exploration of the Service Space

11



Entity Selection



Service Selection

Input parameters

Zillow by Location

State

NY

City

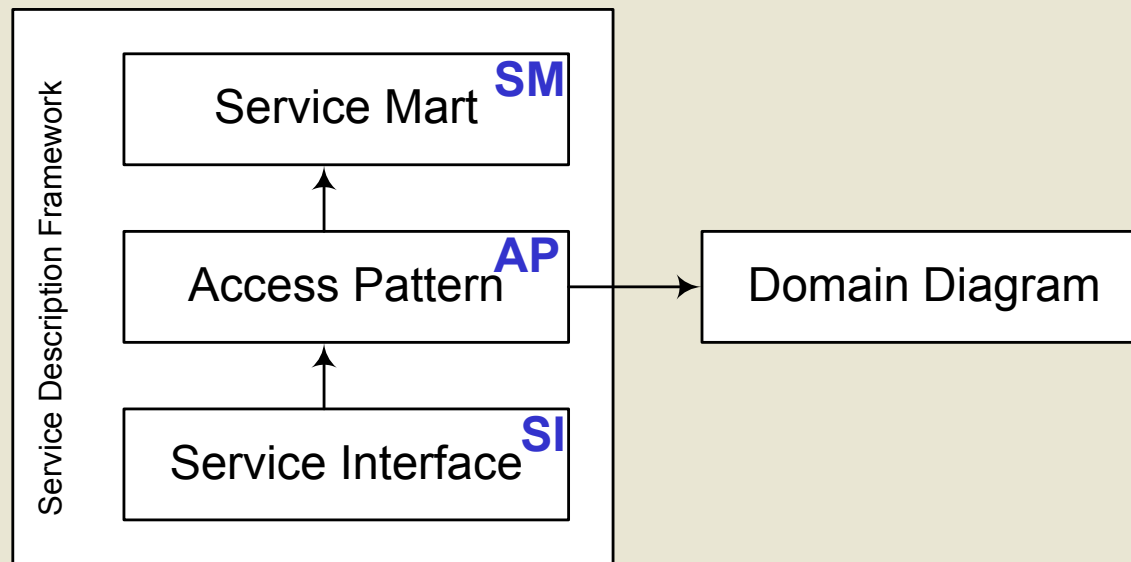
Brooklyn

Type

all

Back Submit

Submit
Query



Formulation of SeCo Queries

Logical and Physical Query plans

The Panta Rhei execution framework

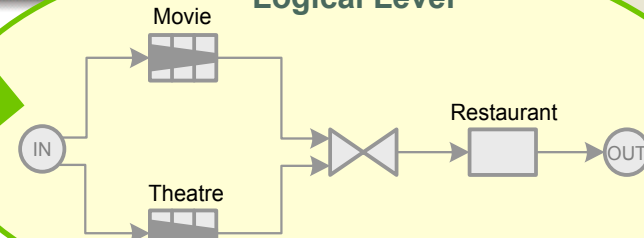
Query processing: Declarative, Logical and Physical formulation

14

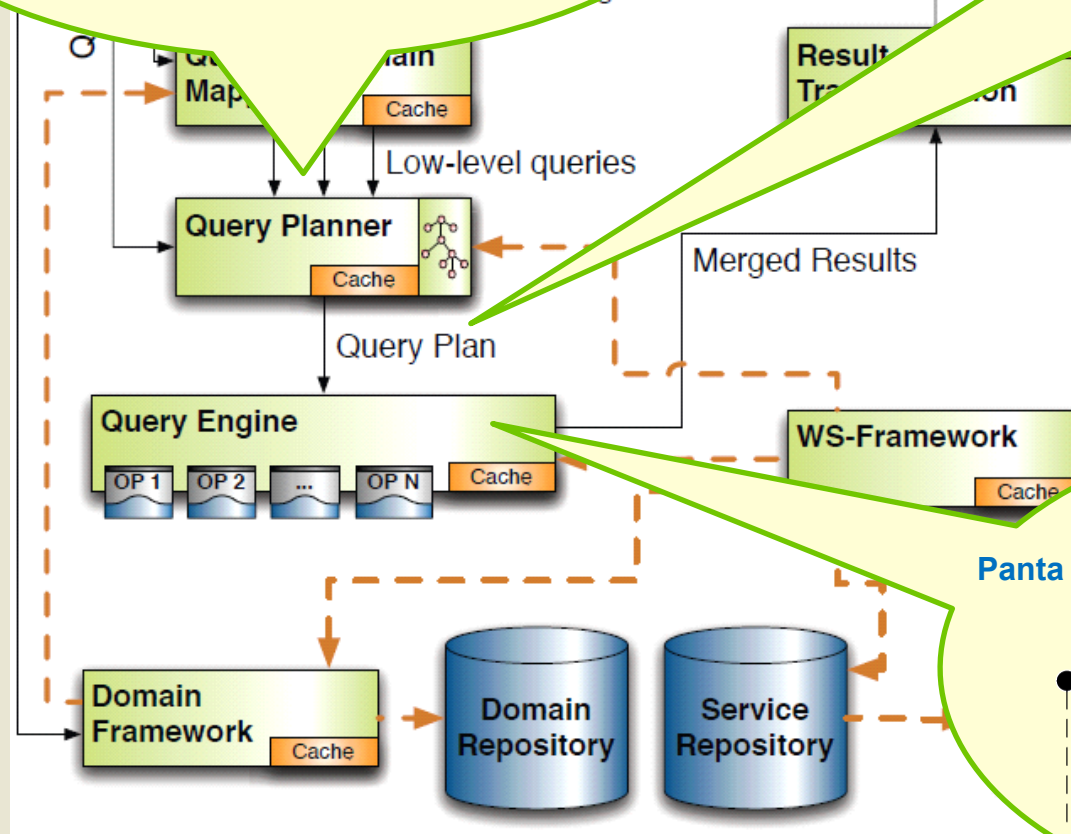
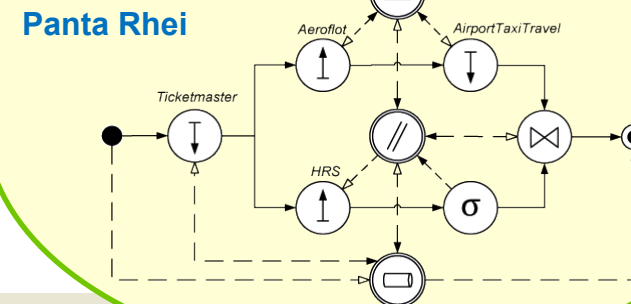
SeCoQL (declarative)

```
DEFINE QUERY NightOut ( $Address : String , $City : String , $Country : String ) AS  
SELECT *  
FROM theatre_by_address ( iAddress : $Address , iCity : $City , iCountry : $Country ) AS T USING Google_Theatre_by_Addr  
JOIN Movies_by_Genre ( iGenre : Drama ) AS M USING IMDB_Movies_by_Genre  
ON T.title = M.title  
JOIN Restaurant_by_Address ( iAddress : T.address , iCity : $City , iCountry : $Country ) AS R USING QMaps_Rest_by_Addr  
WHERE M.year < 2008  
RANK BY ( T = 0.4 , R = 0.5 , M = 0.1 )  
LIMIT 5 TUPLES AND 4 CALLS
```

Logical Level



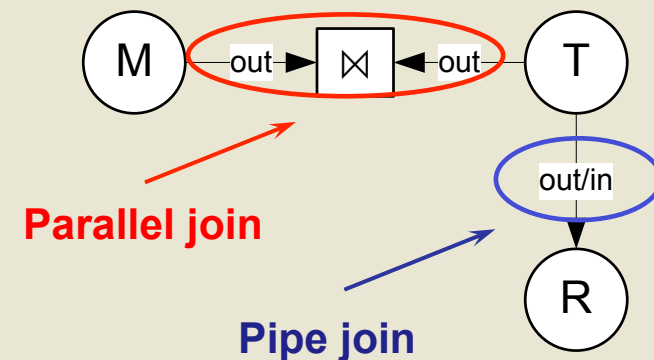
Physical Level



A good old drama movie showing tonight in a theatre close to a good restaurant

NightOut (“via Golgi”, “Milan”, “Italy”)

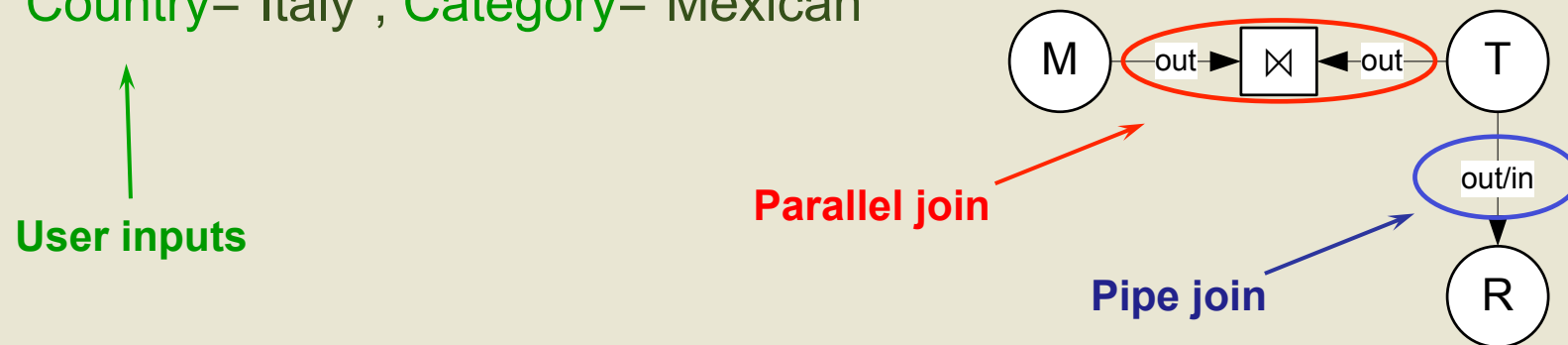
```
DEFINE QUERY NightOut ( $Address : String , $City : String , $Country : String ) AS
  SELECT *
  FROM theatre_by_address ( iAddress : $Address , iCity : $City , iCountry : $Country ) AS T USING Google_Theatre_by_Addr
  JOIN Movies_by_Genre ( iGenre : Drama ) AS M USING IMDB_Movies_by_Genre
  ON I.title = M.title
  JOIN Restaurant_by_Address ( iAddress : T.address, iCity : $City , iCountry : $Country ) AS R USING GMaps_Rest_by_Addr
  WHERE M.year < 2008
  RANK BY ( T = 0.4 , R = 0.5 , M = 0.1 )
  LIMIT 5 TUPLES AND 4 CALLS
```



First step: from conjunctive queries to logical plans

A good American action movie showing tonight in a theatre in Milan, close to a good Mexican restaurant

results(Title, TAddress, RName, RPhone, RAddress, ShowTimes.st) :—
Movie(Genre^I, MCountry^I, Title^O, Score^O, Genres.Genre^O, Actors.Act^O)
Theatre(Country^I, City^I, TName^O, TAddress^O, TPhone^O, Titles.Title^O,
RunningTimes.rt^O, ShowTimes.st^O)
Restaurant(Country^I, City^I, Category^I, TAddress^I, RName^O,
RAddress^O, RPhone^O, RScore^O, Categories.cat^O)
Genre="Action", MCountry="USA", City="Milan"
Country="Italy", Category="Mexican"



Query Processor Overview: Two-Step Approach

17

- Logical level
 - **Input:** a **declarative conjunctive** query
 - Access Patterns and Service Interfaces have already been chosen in higher layers, and join types are therefore known
 - **Output:** a **logical plan**, a specification of a workflow with quantitative estimates of the size of partial results
 - The planner exploits the available degrees of freedom to fix the topology, the number and sequence of service invocations, the join strategies, ...
- Physical level
 - From the **query planner** level to **run-time enactment**
 - **Output:** the **physical plan** in the **Panta Rhei** language
 - Directly executable by the engine
 - **Distribution, parallelization and replication** are considered

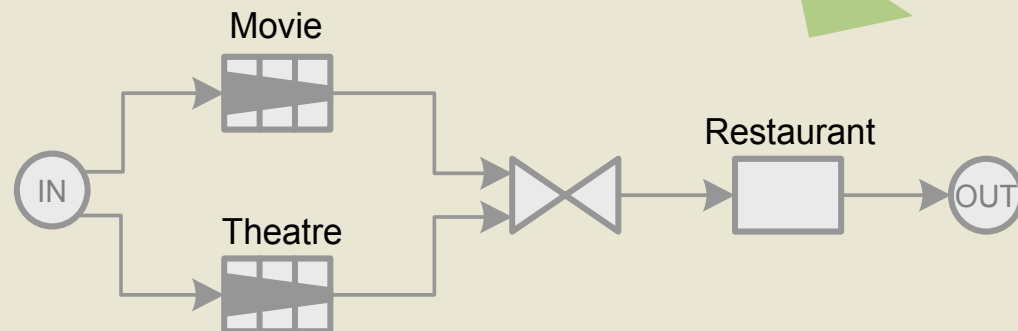
First step: from conjunctive queries to logical plans

18

```
DEFINE QUERY NightOut ( $Address : String , $City : String , $Country : String ) AS
  SELECT *
  FROM theatre_by_address ( iAddress : $Address , iCity : $City , iCountry : $Country ) AS T USING Google_Theatre_by_Addr
  JOIN Movies_by_Genre ( iGenre : Drama ) AS M USING IMDB_Movies_by_Genre
    ON T.title = M.title
  JOIN Restaurant_by_Address ( iAddress : T.address, iCity : $City , iCountry : $Country ) AS R USING GMaps_Rest_by_Addr
  WHERE M.year < 2008
  RANK BY ( T = 0.4 , R = 0.5 , M = 0.1 )
  LIMIT 5 TUPLES AND 4 CALLS
```

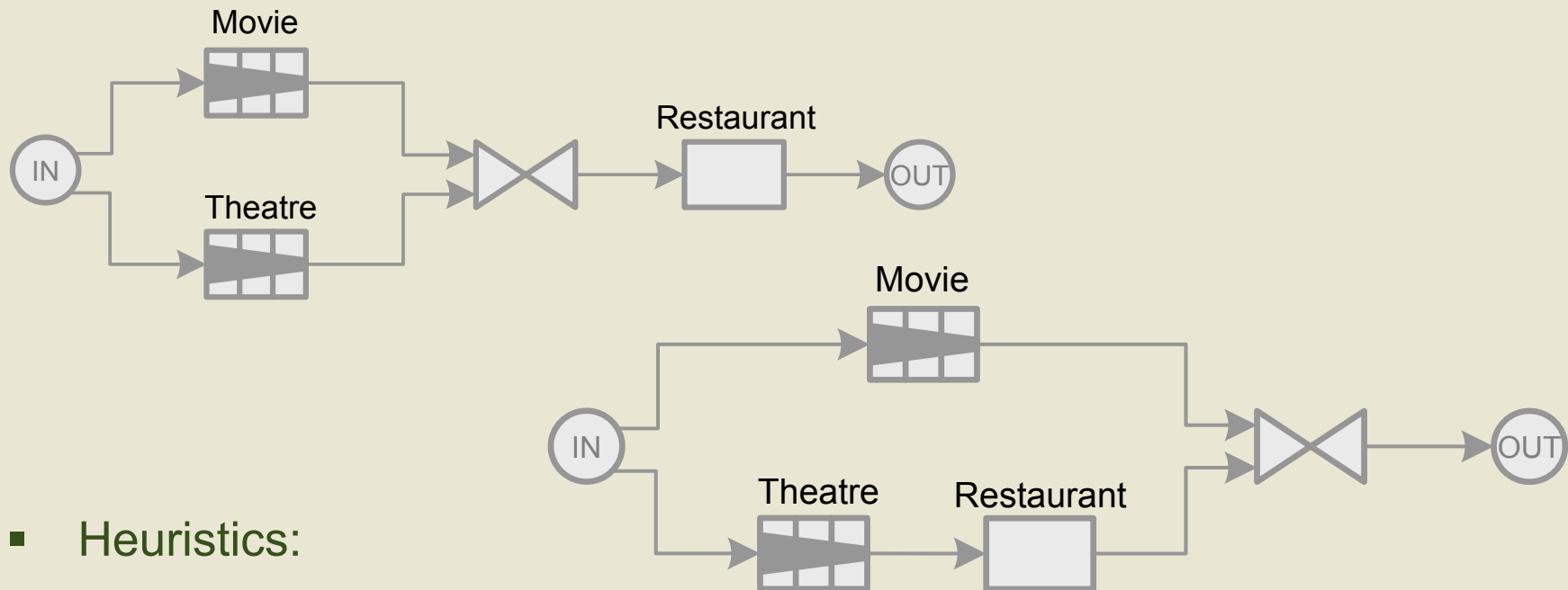
Generation of a logical plan

*Estimated number
of invocations and
cardinality of
inputs and outputs*



Choice of topology for the logical plan

19



■ Heuristics:

– Selective first

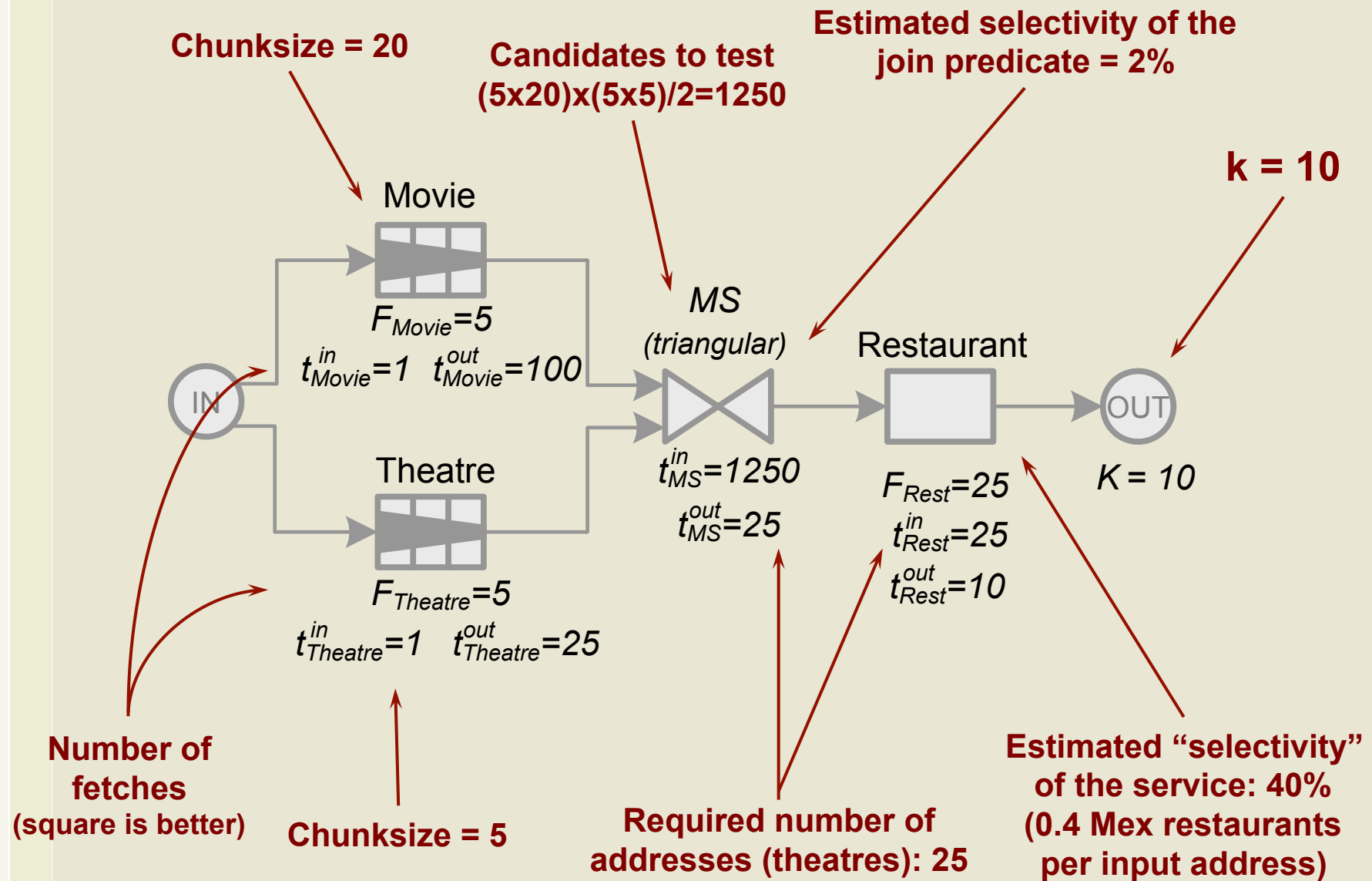
- Anticipating selective invocations minimizes intermediate results
- Expensive services can be invoked later, with qualified input only

– Parallel is better

- Invocations with no i/o dependencies should always be performed in parallel, to minimize latency

Estimating the number of fetches

20

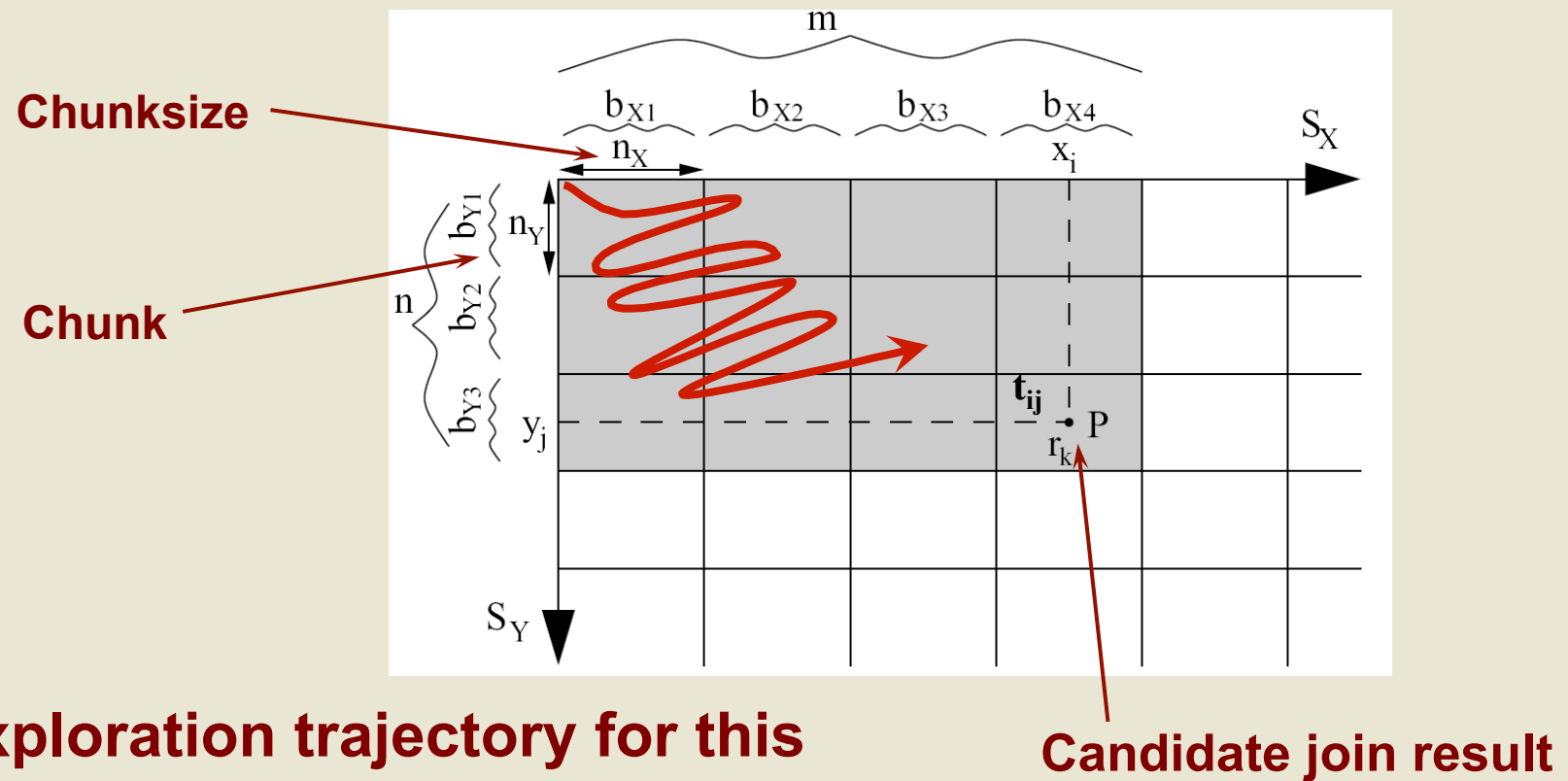


Choice of the join strategies (search space)

21

■ The **join search space**

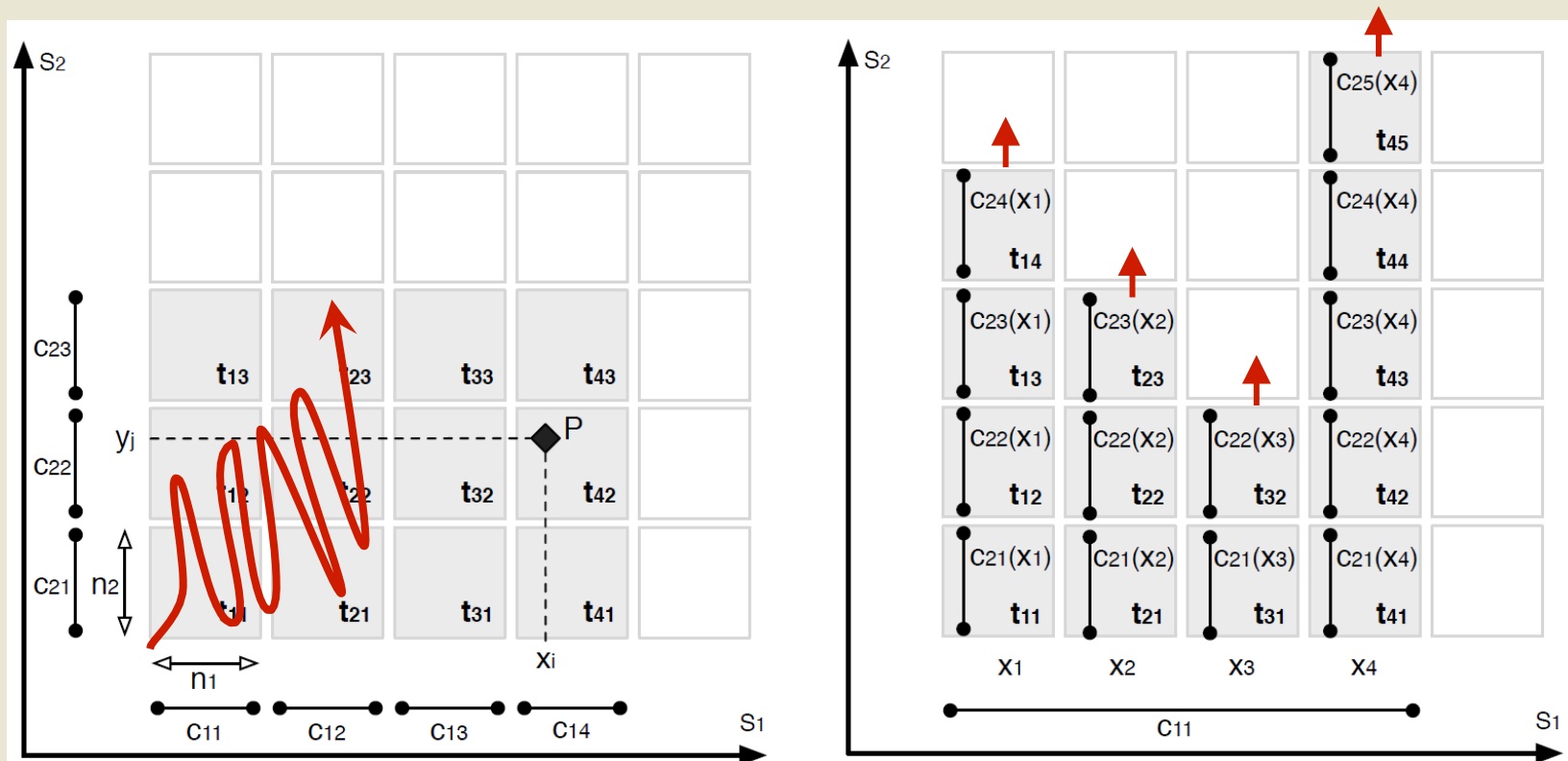
- Different explorations for different joins methods under different assumptions and with different guarantees



Any exploration trajectory for this space is a parallel join strategy

Search space: parallel join vs pipe join

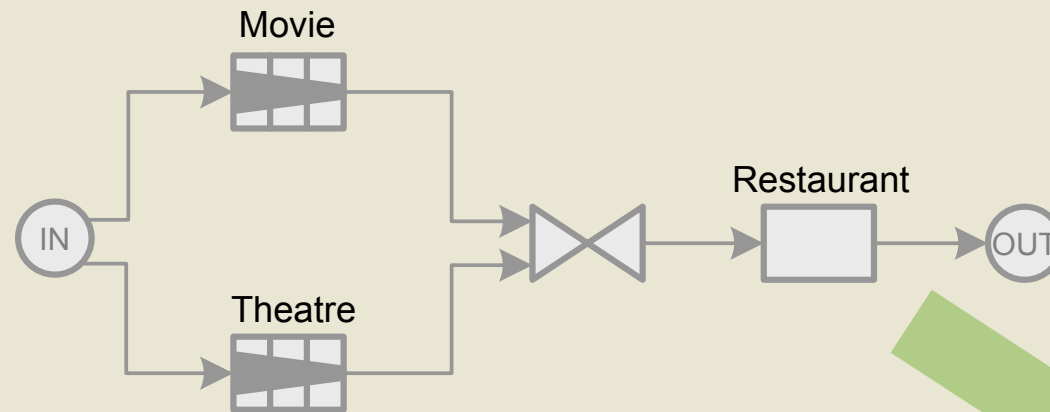
22



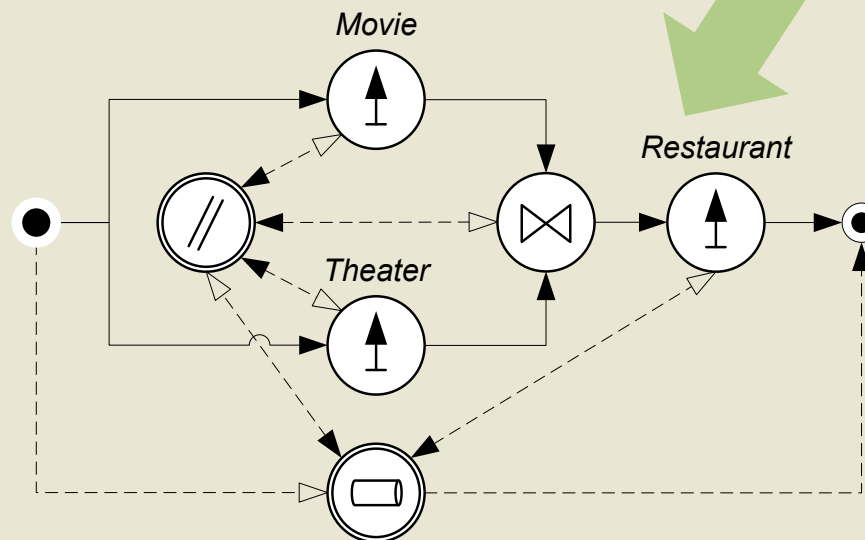
In the case of pipe joins, a join strategy is a pulling strategy deciding at each step from which invocation the next chunk is to be fetched

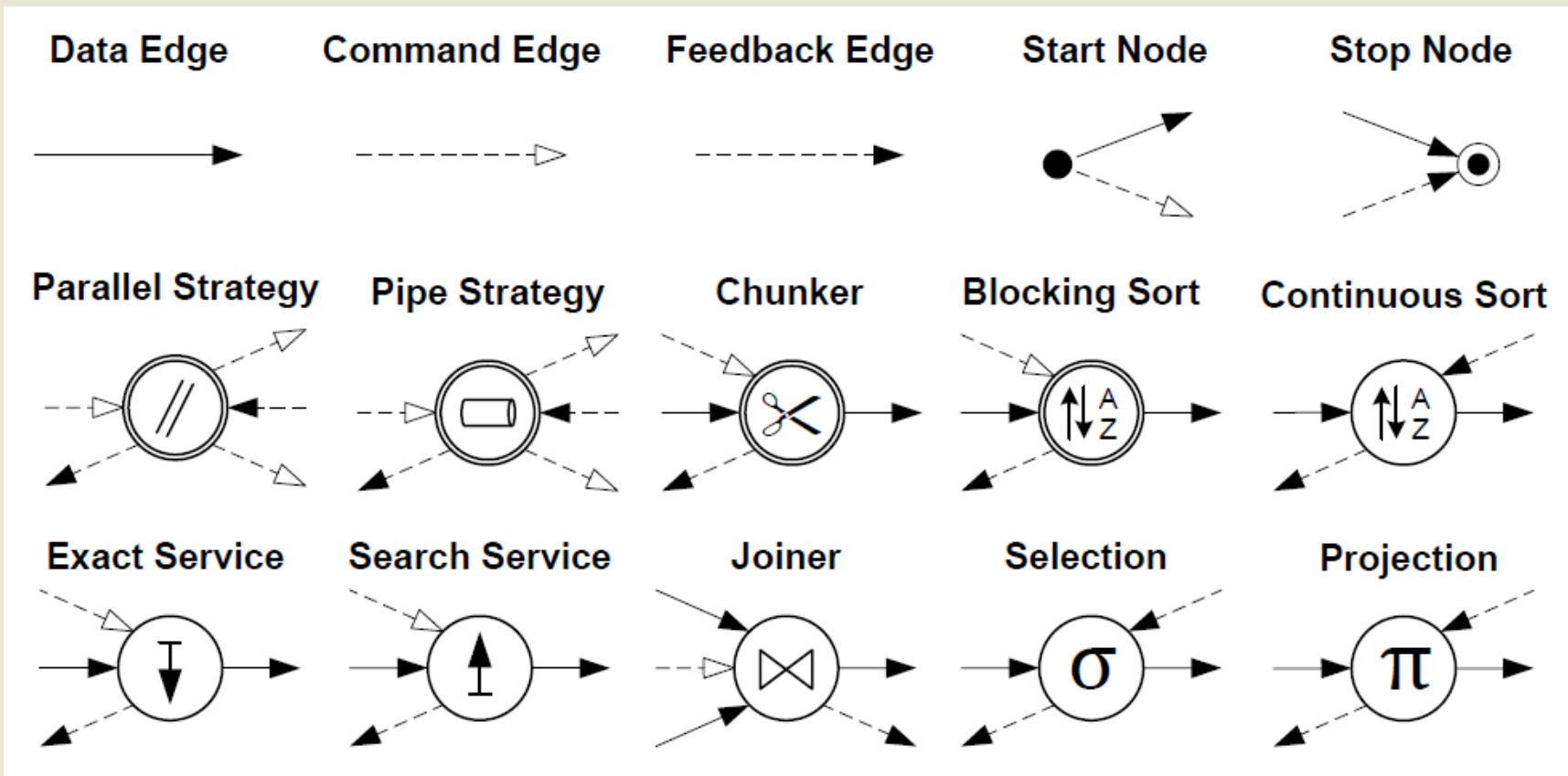
Second step: from logical to physical query plans

23



Then the planner generates a physical, executable query plan, expressed in *Panta Rhei*







- A query execution plan (QEP) accepts an incoming **data and control flow** edge and has a data and control flow edge in output
 - The incoming data edge buffers chunks of tuples in input, while the outgoing data edge transmits chunks of tuples in output to a downstream query execution plan or to the stop node.
 - The commands that are received through the incoming control edge express how the tuples in input relate to chunks in output. The outgoing control edge transmits both new commands generated by the query execution plan and feedback data about the execution of the plan.

Composition Rules (QEPs as basic building blocks)

26

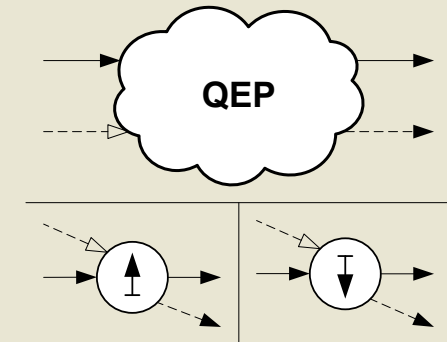
$$(1) \quad QEP := \{ \uparrow, \downarrow \}$$

$$(2) \quad QEP := QEP \cup \{ \otimes, \otimes^A_Z, \otimes^A_Z, \sigma, \pi \}$$

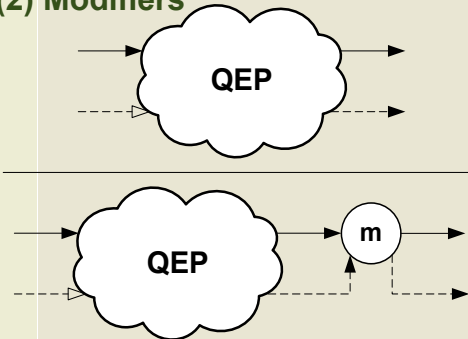
$$(3) \quad QEP := QEP \bowtie_{pipe} QEP$$

$$(4) \quad QEP := QEP \bowtie_{parallel} QEP$$

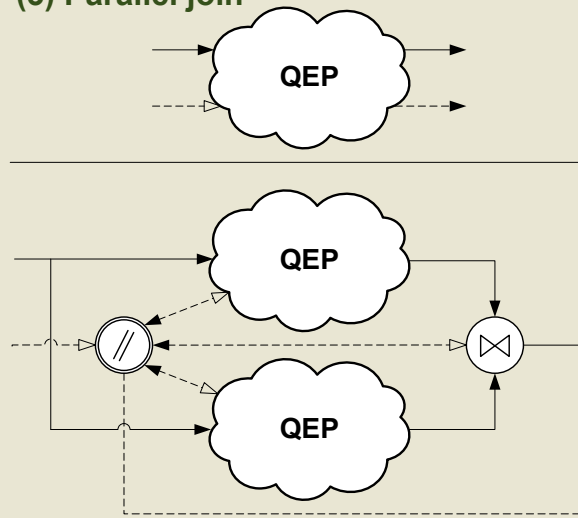
(1) Service invocations



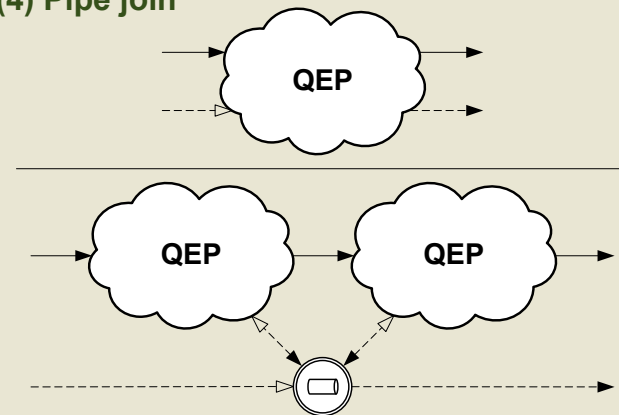
(2) Modifiers



(3) Parallel join

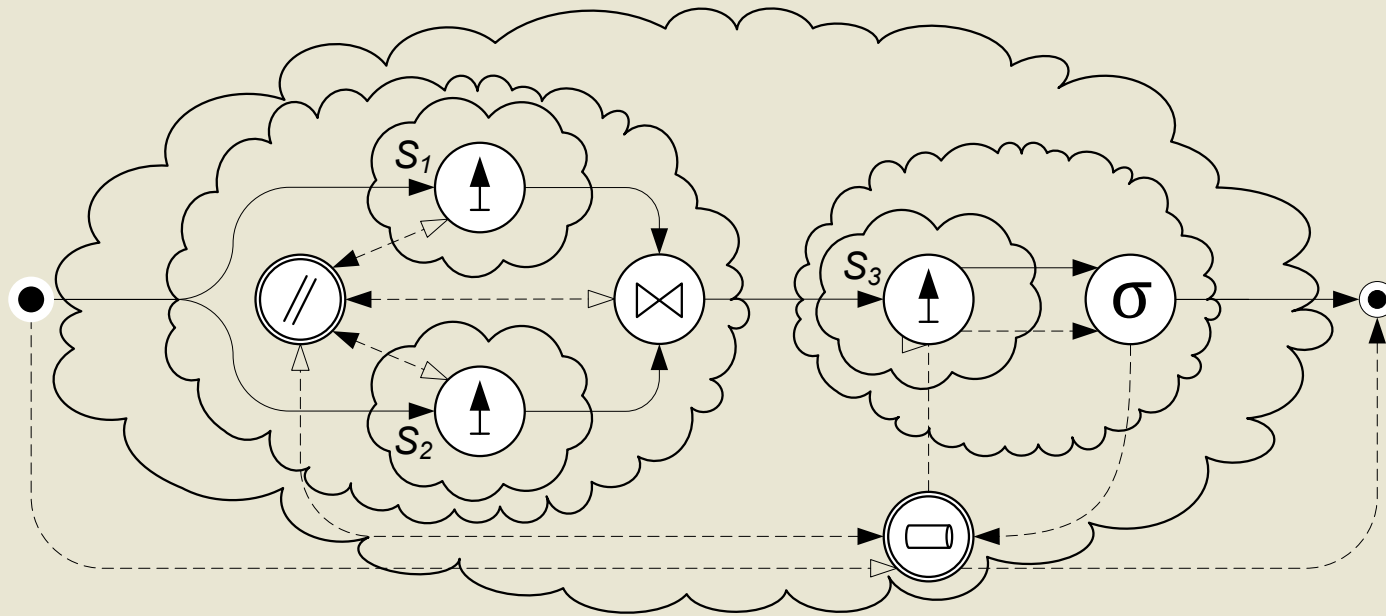


(4) Pipe join



Query Plan Compilation: Topology

27



- Traversal of logical plan to identify query execution plans
 - logical plan is processed “backwards”
 - returns an ordered list of query execution plans
- Processing of query execution plans
 - insertion of unit configurations and parameters
 - consecutive query execution plans are combined with pipe joins
- Right associativity
 - $QEP \bowtie_{\text{pipe}} QEP \bowtie_{\text{pipe}} QEP = QEP \bowtie_{\text{pipe}} (QEP \bowtie_{\text{pipe}} QEP)$

Support for Forward and Backward Control

28

- Strategy depends on whether execution plan is scheduled in forward or backward manner
 - **Forward:** the strategy is entirely determined by scheduler at compile-time and configured in strategy units. The strategy defines which input tuple should be used in the next execution of a query execution plan.
 - **Backward:** there is no predefined strategy, but an algorithm that is encapsulated in the strategy units controls the sending of commands to the dependent query execution plan
 - **“Middle Way” (adaptive):** the strategy unit has a predefined strategy but takes feedback into consideration to allow a bounded deviation from the plan or react to user input.

- Operation is entirely defined by the inputs and can be distributed, parallelized and replicated



Search Service Invocation: Completes the tuples in input by invoking a search service. Per service invocation, it returns a chunk of completed tuples in output.



Joiner: Receives chunks of tuples on two different data edges and joins them into chunks of tuples according to a predicate.



Selection: Filters chunks of tuples according to a selection predicate. Since the selection unit does not re-chunk the tuples, the chunk size can decrease in the order of the selectivity of the given predicate.

- Operation depends on the inputs and the internal state of the unit.



Parallel Strategy: Controls two query execution plans that are composed in a parallel join. It breaks down the commands it receives in input and distributes them to the two sub-plans and, optionally receives feedback from the query execution plans.



Pipe Strategy: Controls two query execution plans that are composed in a pipe join. It breaks down the commands it receives in input and distributes them to the two sub-plans and, optionally receives feedback from the query execution plans.

Good vs Top-k joins

Pluggable join strategies

Demo

System architecture

Guarantee of top-k-ness

32

- Good-k strategies
 - Rely on the sorting of data sources and produce query output as soon as it is available

- Top-k strategies
 - The output is ordered according to a score function
 - Results are sent to the user when there is a guarantee that no result with higher score can be produced
 - Rely on the sorting of data sources
 - Slower strategies
 - Higher time-to-screen

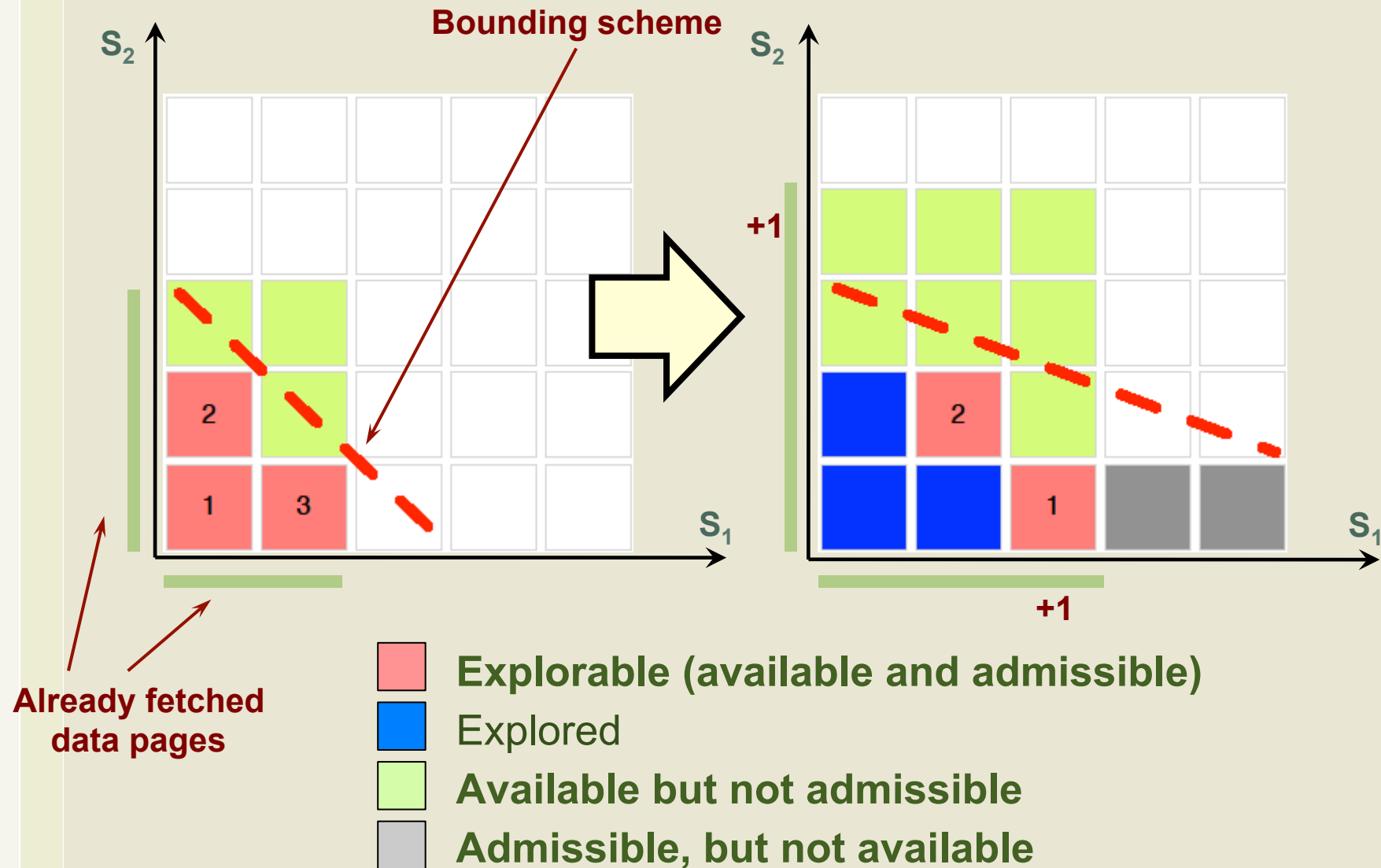
Offline good strategies

33

- Sequences of steps
- Pulling strategy
 - The number of invocations to be sent to each service
- Bounding scheme
 - Determines the tiles which will be joined
- No adaptation
 - Pulling strategy and bounding scheme are defined at compile-time

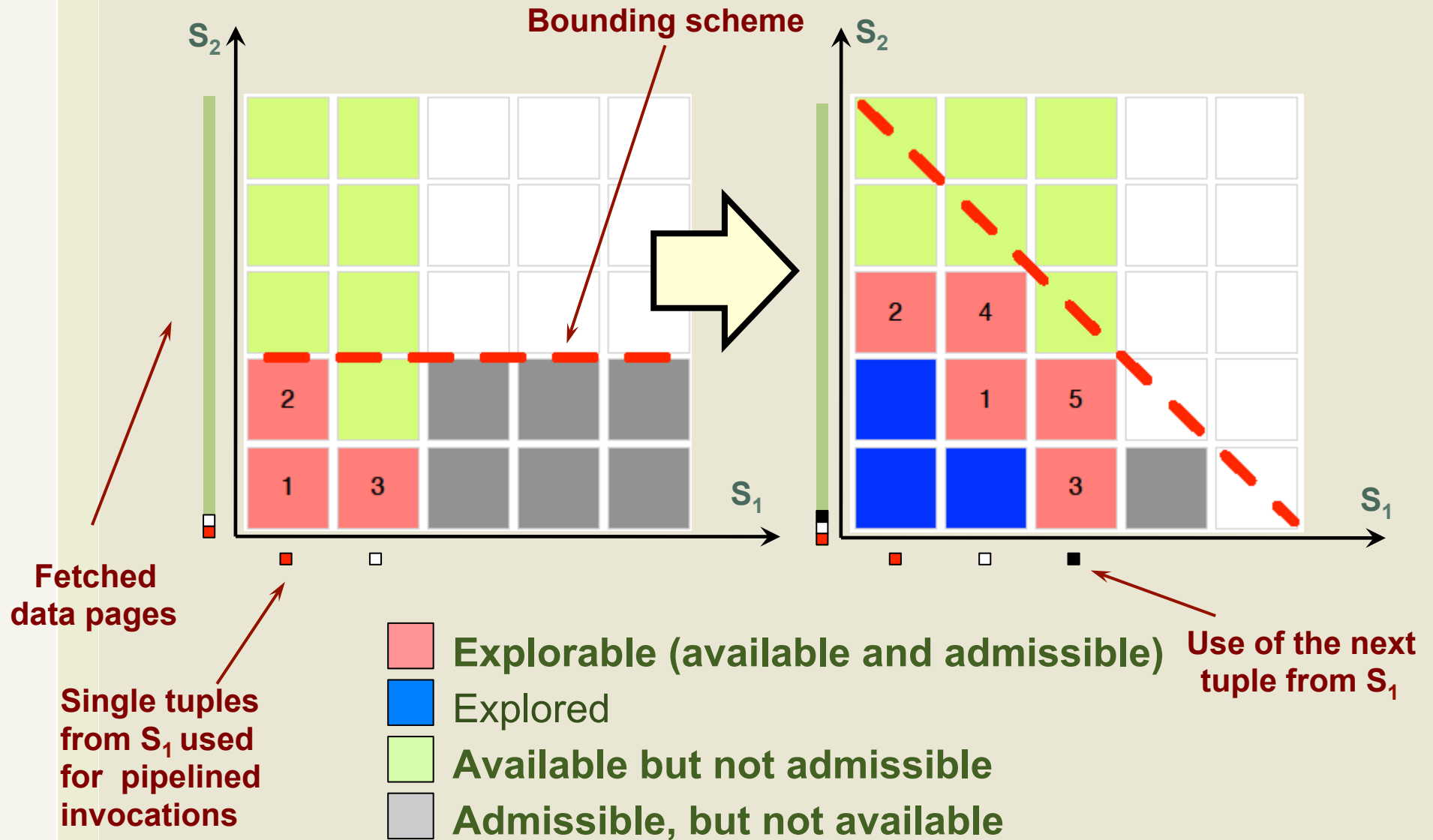
Parallel join: a strategy with a bounding scheme

34



Pipe join: a strategy with a bounding scheme

35



Offline good parallel join strategy

36

ALGORITHM 1: Offline good parallel join strategy

Input: $S = (S_1, S_2, \dots, S_n)$, set of n data sources to be joined

Input: $\mathcal{J} = (J_1, J_2, \dots, J_m)$, strategy; $J_i = \{P_i, b_i\}$

Input: $P_i = \{P_{1,i}, P_{2,i}, \dots, P_{n,i}\}$ pulling strategy at step i

Input: b_i bounding scheme at step i

Data: \mathcal{A} , set of admissible tiles of the join space

Data: \mathcal{R} , set of explored tiles of the join space

Data: \mathcal{E} , list of explorable tiles of the join space

```
1  $\mathcal{R} \leftarrow \emptyset;$ 
2 for  $j = 1$  to  $|\mathcal{J}|$  do
3    $\mathcal{E} \leftarrow \emptyset;$ 
4   // Pulling strategy
5   for  $i = 1$  to  $|S|$  do
6     for  $k = 1$  to  $P_{i,j}$  do
7        $\text{sendFetch}(S_i);$ 
8   // Bounding scheme
9    $\mathcal{A} \leftarrow \text{admissible}(b_j);$ 
10  foreach  $X \in \mathcal{A}$  do
11    if  $\text{isAvailable}(X) \wedge X \notin \mathcal{R}$  then
12       $\mathcal{E} \leftarrow \mathcal{E} \cup X;$ 
13   $\mathcal{E} \leftarrow \text{sortByDistance}(\mathcal{E}, b_j);$ 
14  foreach  $X \in \mathcal{E}$  do
15     $\text{sendJoinCommand}(X);$ 
16     $\mathcal{R} \leftarrow \mathcal{R} \cup X;$ 
```

Online top strategies

37

- Invocations are driven by scores
- The next service to be invoked is the most promising one
 - i.e., the one that could lead to the unseen combination with the highest score
- Rely on the feedback received by services and joiner to decide the next invocation to be performed

ALGORITHM 2: HRJN Parallel Strategy

Input: $S = \{S_1, \dots, S_n\}$, set of n data sources to be joined

Input: $C_{total}(s)$, total number of fetches assigned by the exploration strategy

Data: $F_{aggr}(S, r)$, function modeling a join combination $r \in \mathbb{N}_0^n$ in terms of its score (weighted by each source $s \in S$ and normalized to $[1 \dots 0]$)

Data: $C_{sent}(s)$, count of sent fetches for data source $s \in S$

Data: $t(s)$, threshold for data source $s \in S$, computed as $F_{aggr}(S, r)$

Data: $eof(s)$, end of file flag for data source $s \in S$

Data: $R = \emptyset$, set of join space points (x_1, \dots, x_n) for which a join command has been sent

Data: $Q = \emptyset$, set of join space points (x_1, \dots, x_n) for which a join combination has not been emitted

Data: $t = 1.0$, current score threshold (normalized to $[1 \dots 0]$)

```

1 begin
2   foreach  $s \in S$  do
3      $C_{sent}(s) \leftarrow 0$ ;
4      $t(s) \leftarrow 1.0$ ;
5      $eof(s) \leftarrow false$ ;
6   // invoke every service at least once
7   foreach  $s \in S$  do
8     if  $C_{total} > 0$  then
9        $eof(s) \leftarrow \text{SendFetch}(s)$ ;
10       $\text{UpdateSource}(s, \text{ReceiveFeedback}(s))$ ;
11       $C_{total} \leftarrow C_{total} - 1$ ;
12       $C_{sent}(s) \leftarrow C_{sent}(s) + 1$ ;

```

```
12  while  $C_{total} > 0$  do
    // choose next service to invoke and update its metadata
13   $hScore \leftarrow 0$ ;
14  foreach  $s \in S$  do
15      if  $\neg eof(s)$  and  $t(s) > hScore$  then
16           $hScore \leftarrow t(s)$ ;
17           $z \leftarrow s$ ;
18   $eof(z) \leftarrow \text{SendFetch}(z)$ ;
19   $\text{UpdateSource}(z, \text{ReceiveFeedback}(z))$ ;
20   $C_{total} \leftarrow C_{total} - 1$ ;
21   $C_{sent}(z) \leftarrow C_{sent}(z) + 1$ ;
22   $t \leftarrow \text{UpdateThreshold}()$ ;
    // create join combinations
23  foreach  $r \leftarrow (r_1, \dots, r_n) \in \mathbb{N}_0^n$  and  $r_i \in \{1, \dots, C_{sent}(s_i)\}$  do
24      if  $r \notin R$  then
25           $\text{SendJoinCommand}(r)$ ;
26           $R \leftarrow R \cup r$ ;
27           $Q \leftarrow Q \cup r$ ;
    // emit top-k tuples
28  foreach  $r \in Q$  do
29      if  $F_{aggr}(S, r) \geq t$  then
30           $Q \leftarrow Q \setminus r$ ;
31       $\text{Emit}(r)$ ;
```

Adaptive online strategies

40

- The feedback allows handling unexpected situations
 - Unexpected score values
 - Unexpected predicate selectivity and result cardinality
 - Unexpected service response time
 - Data source failures
 - Result diversification

Diversification parallel join strategy

41

ALGORITHM 3: Diversification parallel join strategy

Input: $S = (S_1, S_2, \dots, S_n)$, set of n data sources to be joined

Input: $\mathcal{J} = (J_1, J_2, \dots, J_m)$, strategy; $J_i = \{P_i, b_i\}$

Input: $P_i = \{P_{1,i}, P_{2,i}, \dots, P_{n,i}\}$ pulling strategy at step i

Input: b_i bounding scheme at step i

Data: \mathcal{A} , set of admissible tiles of the join space

Data: \mathcal{R} , set of explored tiles of the join space

Data: \mathcal{E} , list of explorable tiles of the join space

```
1  $\mathcal{R} \leftarrow \emptyset;$ 
2 for  $j = 1$  to  $|\mathcal{J}|$  do
3    $\mathcal{E} \leftarrow \emptyset;$ 
4   // Pulling strategy
5   for  $i = 1$  to  $|S|$  do
6     for  $k = 1$  to  $P_{i,j}$  do
7        $\text{sendFetch}(S_i);$ 
8   // Bounding scheme
9    $\mathcal{A} \leftarrow \text{admissible}(b_j);$ 
10  foreach  $X \in \mathcal{A}$  do
11    if  $\text{isAvailable}(X) \wedge X \notin \mathcal{R}$  then
12       $\mathcal{E} \leftarrow \mathcal{E} \cup X;$ 
13  for  $i = 1$  to  $|\mathcal{E}|$  do
14     $\mathcal{E} \leftarrow \text{sortByDiversification}(\mathcal{E}, \mathcal{R});$ 
15     $X \leftarrow \text{peek}(\mathcal{E});$ 
16     $\text{sendJoinCommand}(X);$ 
17     $\mathcal{R} \leftarrow \mathcal{R} \cup X;$ 
18     $\mathcal{E} \leftarrow \mathcal{E} \setminus X;$ 
```

System architecture

