

Scalable Package Queries in Relational Database Systems

Matteo Brucato



Juan F. Beltran



Azza Abouzied



Alexandra Meliou



Package Queries

- An important class of **combinatorial optimization queries**
- Largely **unsupported** by existing technology
- We make them ***first-class citizens*** in relational databases:
 - The Package Query Language (PaQL)
 - PaQL to Integer Linear Programming (ILP)
- Scalable evaluation through the SKETCHREFINE algorithm

Example: Meal Planning

A dietitian wants to build a **Meal Plan** for a patient

Meal Plan

- 3 gluten-free recipes (breakfast, lunch, dinner)
- At least 2 Cal **in total**
- Lowest **total** fat intake

A Meal Plan is a “*Package*”

Meal Plan

- 3 gluten-free recipes (breakfast, lunch, dinner)
- At least 2 Cal
- Lowest total fat intake

Base Constraint (selection predicate)

- **All recipes** in the meal plan are **gluten-free**

Cardinality Constraint

- **3** recipes

Summation Constraint

- **≥ 2** Calorie in **total**

Objective Criterion

- **Lowest total** fat intake

Global:
True of the entire
set of recipes

Using SQL

Meal Plan

- All gluten-free recipes (selection)
- Exactly 3 recipes (cardinality)
- At least 2 Cal in total (summation)
- Lowest total fat intake (objective)

```
SELECT *  
FROM Recipes R  
WHERE R.gluten = 0
```

Using SQL

Meal Plan

- All gluten-free recipes (selection)
- Exactly 3 recipes (cardinality)
- At least 2 Cal in total (summation)
- Lowest total fat intake (objective)

```
SELECT *  
FROM Recipes R  
WHERE R.gluten = 0  
LIMIT 3
```

Using SQL

Meal Plan

- All gluten-free recipes (selection)
- Exactly 3 recipes (cardinality)
- At least 2 Cal in total (summation)
- Lowest total fat intake (objective)

SELECT *

FROM Recipes R1, Recipes R2, Recipes R3

WHERE

R1.gluten = 0 AND R2.gluten = 0 AND R3.gluten = 0
AND R1.kcal + R2.kcal + R3.kcal >= 2.0

Using SQL

Meal Plan

- All gluten-free recipes (selection)
- Exactly 3 recipes (cardinality)
- At least 2 Cal in total (summation)
- Lowest total fat intake (objective)

SELECT *

FROM Recipes R1, Recipes R2, Recipes R3

WHERE

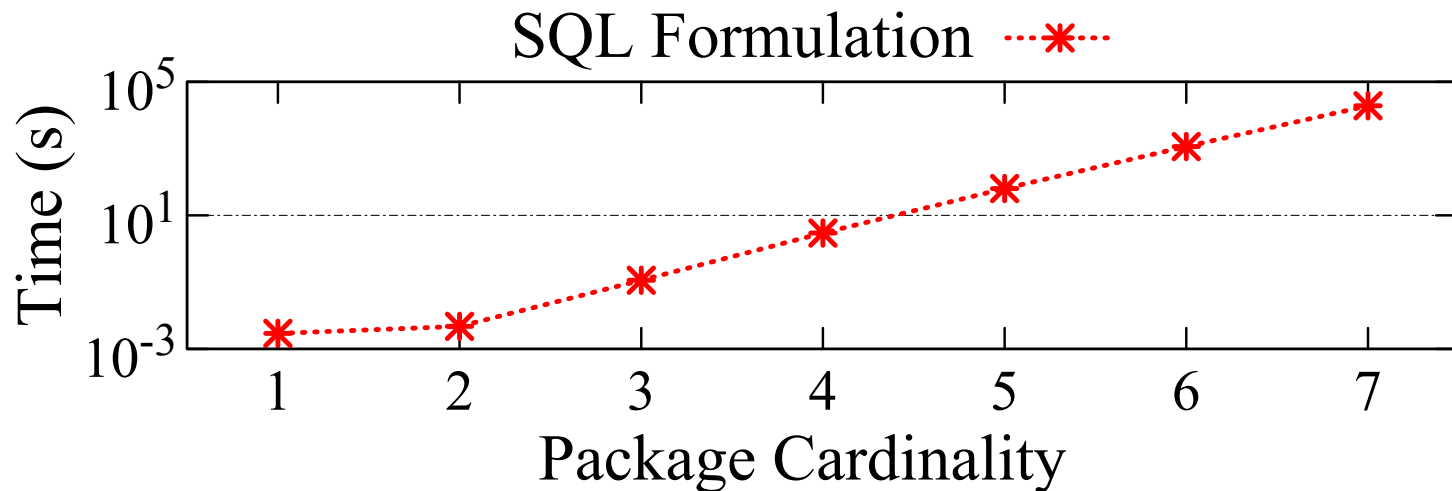
R1.gluten = 0 AND R2.gluten = 0 AND R3.gluten = 0
AND R1.kcal + R2.kcal + R3.kcal >= 2.0

ORDER BY R1.fat + R2.fat + R3.fat

Self-joins: Too Expensive!

Meal Plan

- All gluten-free recipes (selection)
- Exactly 3 recipes (cardinality)
- At least 2 Cal in total (summation)
- Lowest total fat intake (objective)



Using SQL

Meal Plan

- All gluten-free recipes (selection)
- Exactly 3 recipes (cardinality)
- At least 2 Cal in total (summation)
- Lowest total fat intake (objective)

Self-joins are able to express **cardinality** and **linear** constraints.

But what if there is **no cardinality constraints**?...

Using SQL

Meal Plan

- All gluten-free recipes (selection)
- Exactly 3 recipes (cardinality)
- At least 2 Cal in total (summation)
- Lowest total fat intake (objective)

How many self-joins?...

It needs SQL's full recursive power to build all candidate subsets!

Recap: Why is SQL insufficient?

1. With known cardinality, **self-joins** can express package queries, but they are **too expensive to compute**
2. Without known cardinality, self-joins are **not powerful enough**
3. The self-join solution **loses in declarativeness**

Packages as first-class citizens in database systems

- **Packages** are **not unique** to each application
- Nor are the **algorithms** for constructing them
- The **data** typically resides in a database

We introduce:

- A declarative package query language: **PaQL**
- PaQL to ILP **translation** method
- **Scalable** PaQL evaluation method

PaQL: The Package Query Language

Meal Plan

- All gluten-free recipes (selection)
- Exactly 3 recipes (cardinality)
- At least 2 Cal in total (summation)
- Lowest total fat intake (objective)

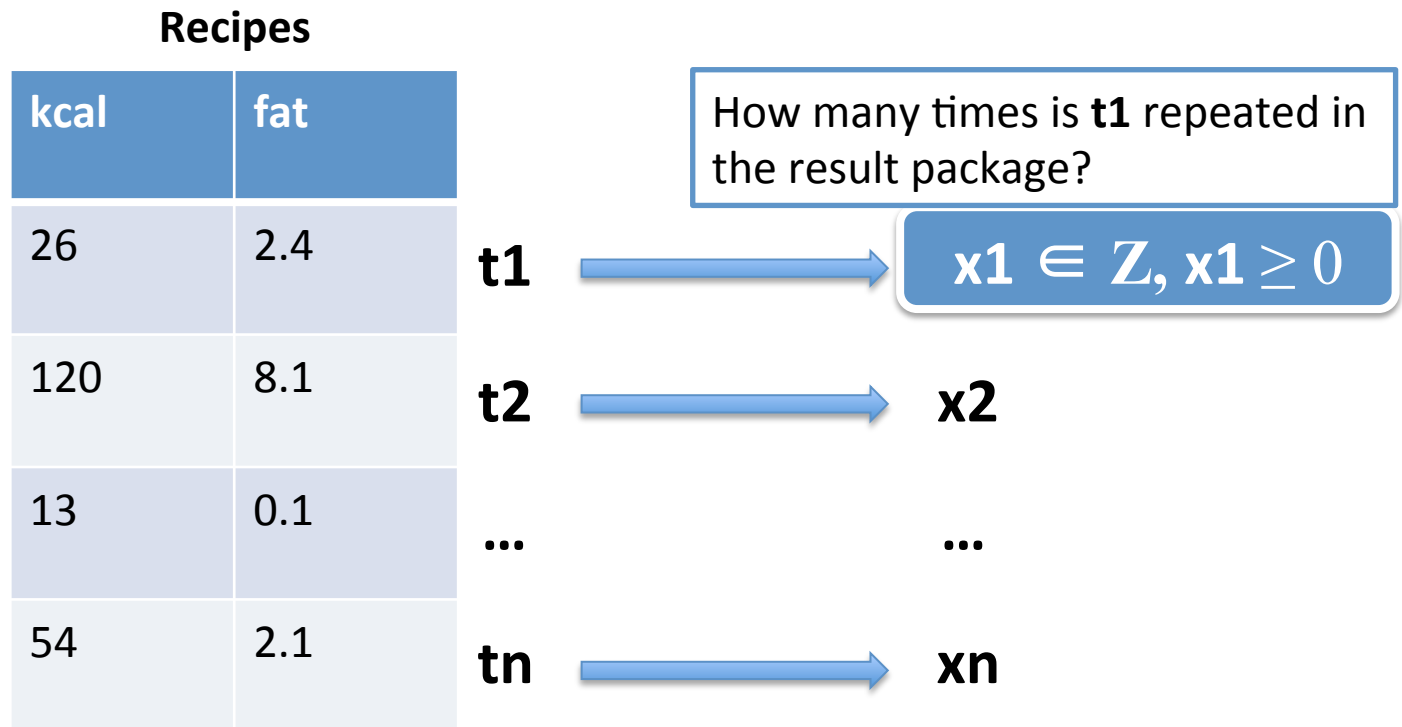
```
SELECT PACKAGE (R)
  FROM Recipes R
  WHERE R.gluten=0
  SUCH THAT COUNT(*)=3 AND
            SUM(kcal) >=2
  MINIMIZE SUM(fat)
```

DIRECT: PaQL to ILP Translation

- An ILP problem consists of:
 - A set of **integer variables**
 - A set of **linear constraints** over the variables
 - A **linear objective function** over the variables
- An ILP solution consists of:
 - An **assignment** to each of the integer variables

PaQL to ILP: Variables

- Integer Variables: one for each tuple

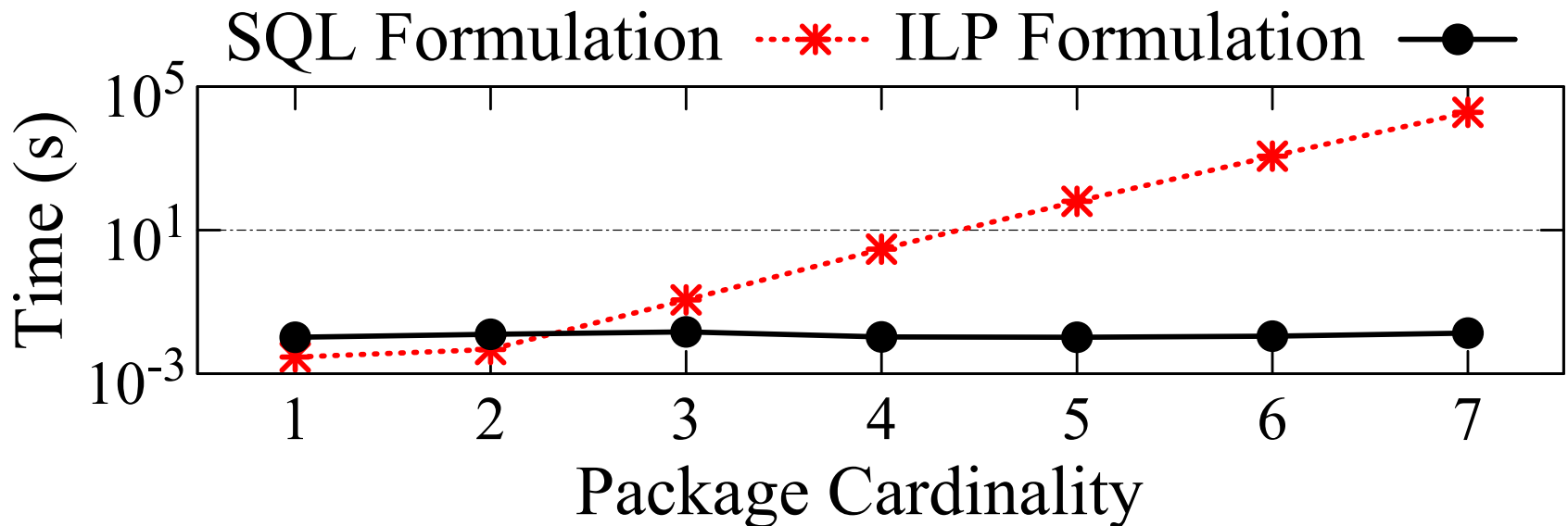


PaQL to ILP: Global Conditions

- **SUCH THAT**
 - $\text{COUNT}(\ast) = 3$
 $\sum_i x_i = 3$
 - $\text{SUM}(\text{kcal}) \geq 2$
 $\sum_i (t_i.\text{kcal} \cdot x_i) \geq 2$
- Objective Criterion
 - MINIMIZE SUM(fat)
minimize $\sum_i (t_i.\text{fat} \cdot x_i)$

DIRECT: Query Evaluation

1. Apply base predicate (selections)
2. Formulate ILP
3. Solve ILP



Drawbacks of DIRECT

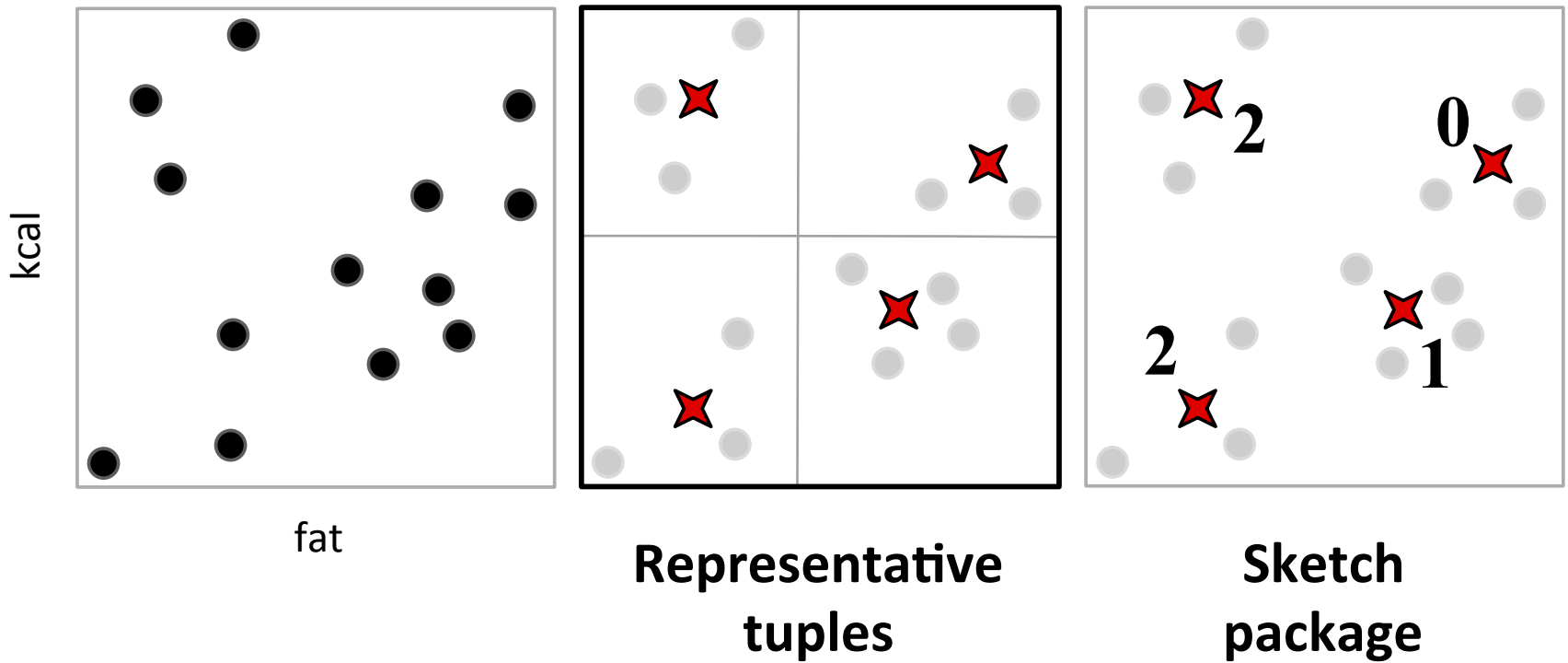
- Only applicable if data **fits** entirely in main memory
- It **may fail** due to the **complexity** of the ILP problem

SKETCHREFINE: Scalable Evaluation

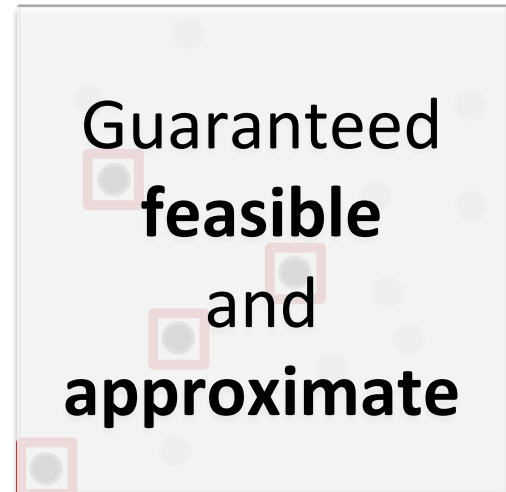
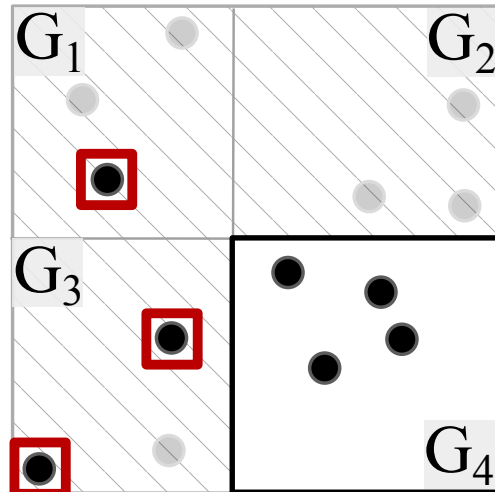
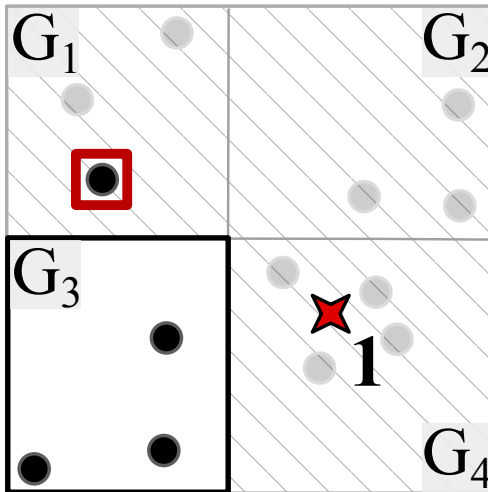
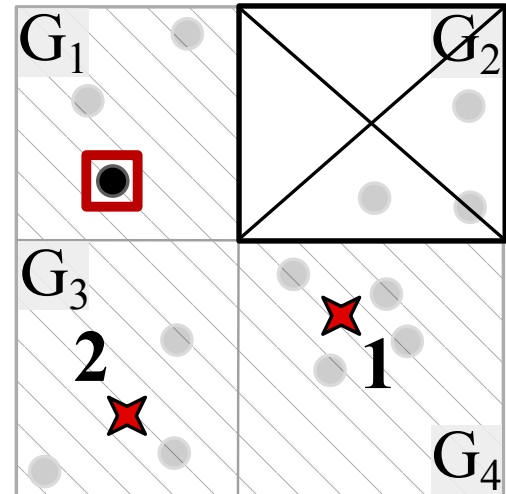
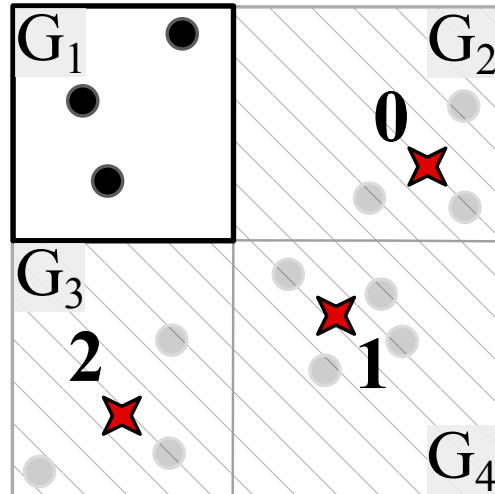
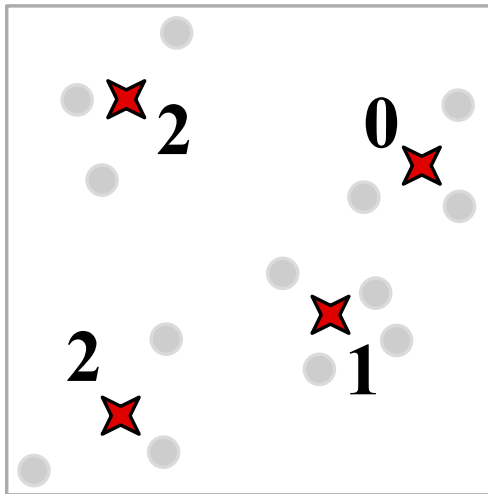
- Partition Data (offline)
 - Into groups of similar tuples
 - Elect a “representative” tuple for each group
- SKETCH
 - an initial package from the representatives
- REFINE
 - the initial package using real tuples

Returns approximate, feasible, package

Partition and SKETCH



REFINE

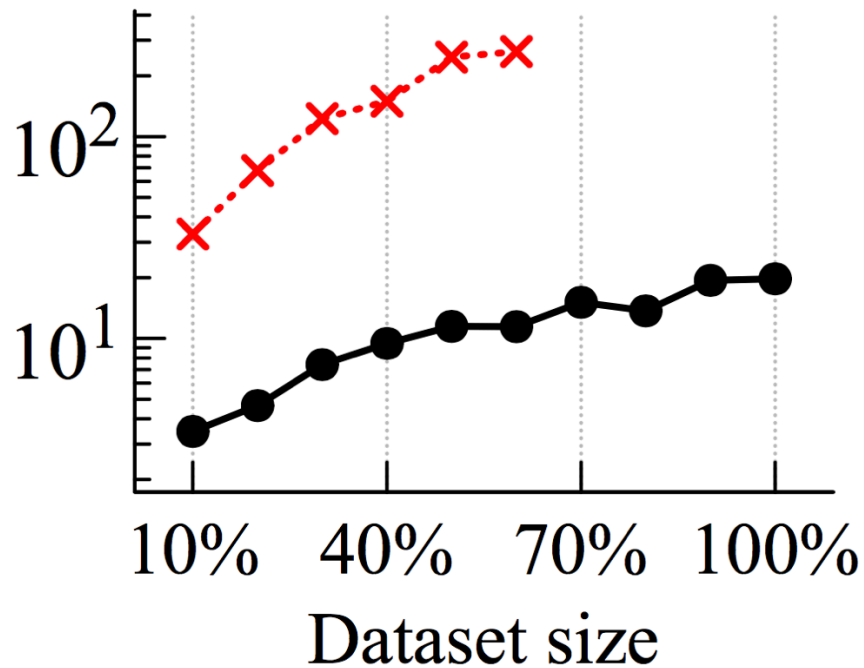


Approximation Guarantees

- SKETCHREFINE is a $(1 \pm \varepsilon)^6$ -approximation with respect to DIRECT
- ε limits the partitions radius
- The partition size (number of tuples per partition) does not affect the approximation guarantee

Scalability

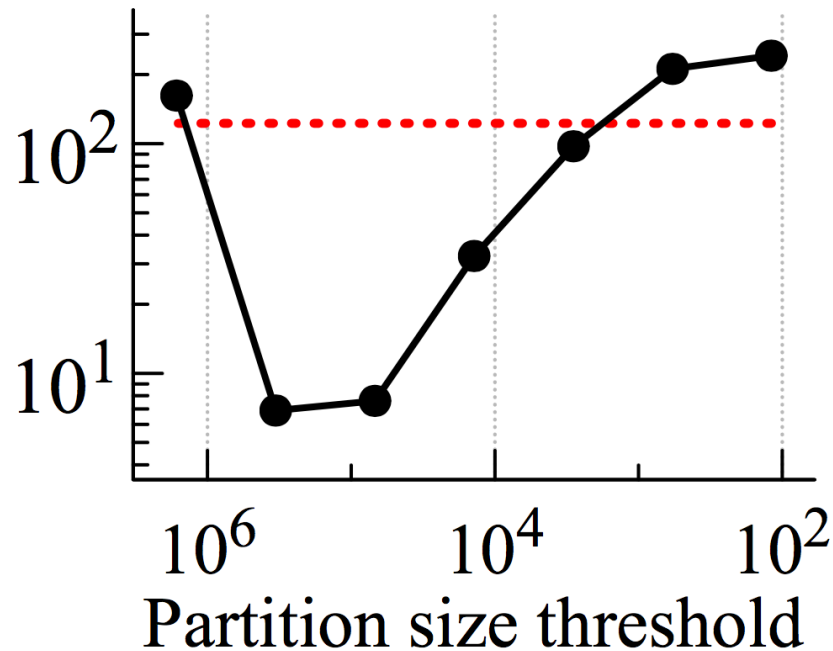
DIRECT x..... SKETCHREFINE ——●——



Approximation Ratio:
Mean: 1.01, Median: 1.00

Partition Size: Performance Impact

DIRECT✕..... SKETCHREFINE ——●——



Approximation Ratio:
Mean: 1.01, Median: 1.00

Thank You!

Matteo Brucato

matteo@cs.umass.edu

