

Human-powered Sorts and Joins

Adam Marcus, Eugene Wu, David Karger, Samuel Madden, Robert Miller

Presenter: Xiaoyi Fan

2016. 10. 12

Outline

- Introduction
- Overview
- Join Operator
- Sort Operator
- Conclusion

Overview of Qurk

- Declarative “workflow” system integrating human intelligence
 - User-defined functions
 - Human Intelligence Tasks (HIT)
 - Capture operations outside relational algebra
 - Typically external API calls

Overview of Qurk

Query Model:

SQL

Overview of Qurk



```
Query =SELECT * FROM photos WHERE  
isFemale(photos.picture);
```

UDF

Overview of Qurk

- Instead of writing code for UDFs, can be described at a high level using Tasks
- Tasks = High level-templates for commonly occurring crowd-operations and or algorithms
- Filter, Generate, Sort, Joins

Overview of Qurk

● Generative Tasks

```
SELECT c.name  
FROM celeb AS c  
WHERE isFemale(c)
```

isFemale defined as follows:

```
TASK isFemale(field) TYPE Filter:  
  Prompt: "<table><tr> \  
    <td><img src='%s'></td> \  
    <td>Is the person in the image a woman?</td> \  
  </tr></table>", tuple[field]  
  YesText: "Yes"  
  NoText: "No"  
  Combiner: MajorityVote
```

Overview of Qurk

- Interface



Yes

No



Yes

No

Overview of Qurk

• Join Operator

```
SELECT c.name  
FROM celeb c JOIN photos p  
ON samePerson(c.img,p.img)
```

```
TASK samePerson(f1, f2) TYPE EquiJoin:  
  SingularName: "celebrity"  
  PluralName: "celebrities"  
  LeftPreview: "<img src='%s' class=smImg>",tuple1[f1]  
  LeftNormal: "<img src='%s' class=lgImg>",tuple1[f1]  
  RightPreview: "<img src='%s' class=smImg>",tuple2[f2]  
  RightNormal: "<img src='%s' class=lgImg>",tuple2[f2]  
  Combiner: MajorityVote
```

Overview of Qurk

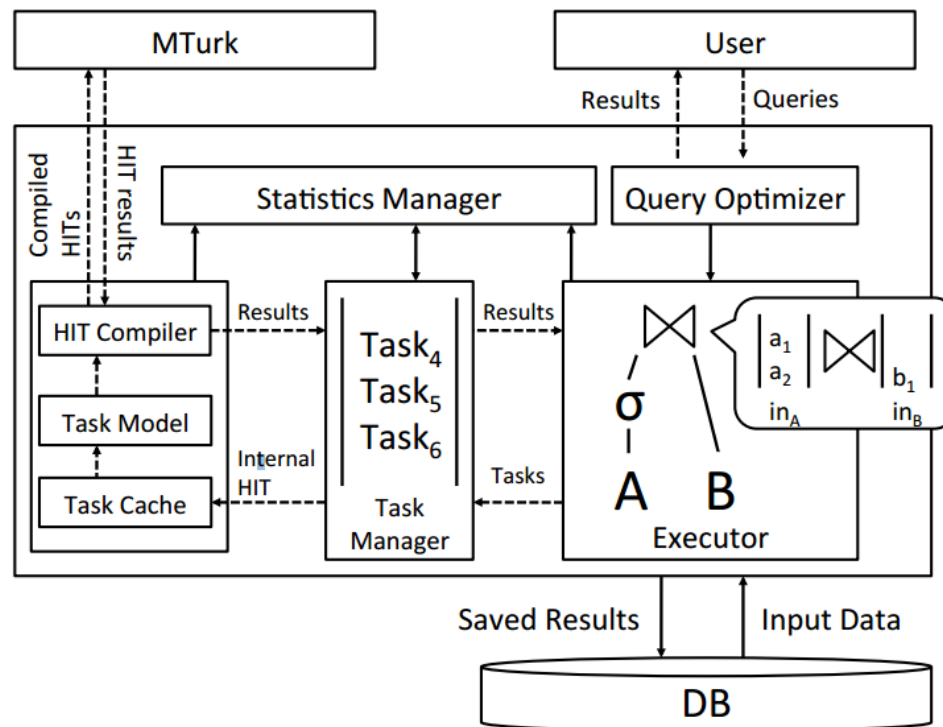
- Sort Operator

```
SELECT squares.label  
FROM squares  
ORDER BY squareSorter(img)
```

```
TASK squareSorter(field) TYPE Rank:  
  SingularName: "square"  
  PluralName: "squares"  
  OrderDimensionName: "area"  
  LeastName: "smallest"  
  MostName: "largest"  
  Html: "<img src='%s' class=lgImg>", tuple[field]
```

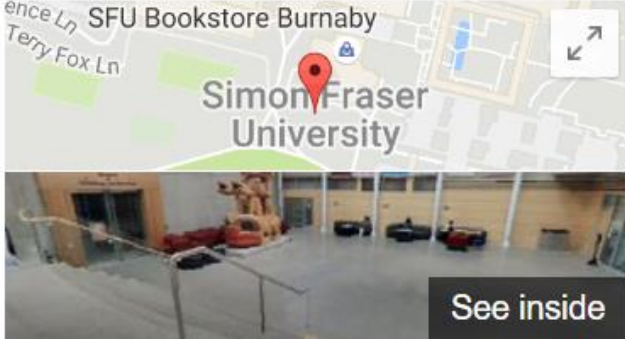

Overview of Qurk

- Qurk system architecture



Join Operator

Simon Fraser
University
==
SFU



Simon Fraser University ★

[Website](#) [Directions](#)

Public university in Burnaby, British Columbia · 400.0 m

Simon Fraser University, commonly referred to as SFU, is a public research university in British Columbia, Canada with campuses in Burnaby, Vancouver and Surrey. [Wikipedia](#)

Address: 8888 University Dr, Burnaby, BC V5A 1S6

Total enrollment: 34,990 (2015)

Mascot: McFogg the Dog

Phone: (778) 782-3111

President: [Andrew Petter](#)

Join Operator

- Matching Celebrities



Join Operator

- Simple Join $O(nm)$

Is the same celebrity in the image on the left and the image on the right?

Yes

No



Join Operator

- Naive Batching $O(nm/b)$

Is the same celebrity in the image on the left and the image on the right?

☐ Yes ☐ No



☐ Yes ☐ No



Submit

Join Operator

• Smart Batching $O(nm/b^2)$

Find pairs of images with the same celebrity

- To select pairs, click on an image on the left and an image on the right. Selected pairs will appear in the **Matched Celebrities** list on the left.
- To magnify a picture, hover your pointer above it.
- To unselect a selected pair, click on the pair in the list on the left.
- If none of the celebrities match, check the **I did not find any pairs** checkbox.
- There may be multiple matches per page.

Matched Celebrities
To remove a pair added in error, click on the pair in the list below.

☐ I did not find any pairs

Submit

4-10x reduction
in cost

Join Operator

- Experiments setup:

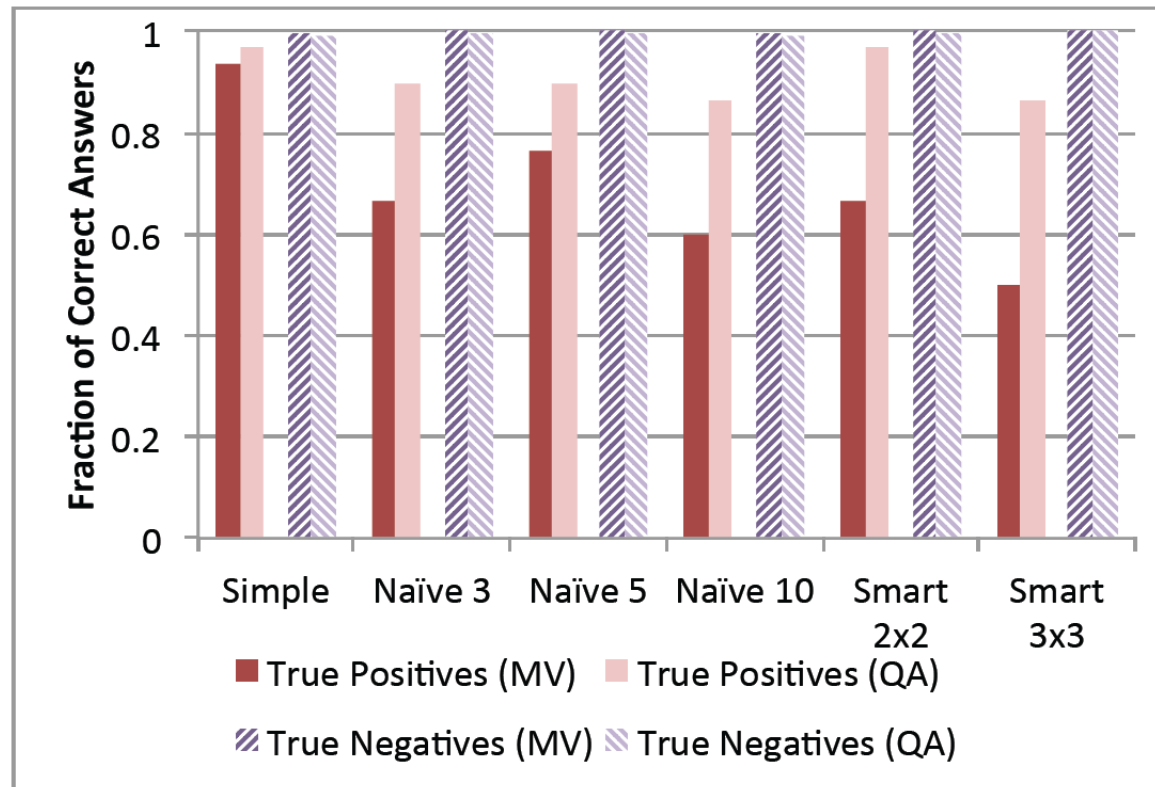
- Compare *IDEAL*, *Simple*, *Naïve*, and *Smart* schemes
- Dataset: *celebrity join dataset*

Implementation	True Pos. (MV)	True Pos. (QA)	True Neg (MV)	True Neg (QA)
IDEAL	20	20	380	380
Simple	19	20	379	376
Naive	19	19	380	379
Smart	20	20	380	379

Baseline comparison of three join algorithms with no batching enabled

Join Operator

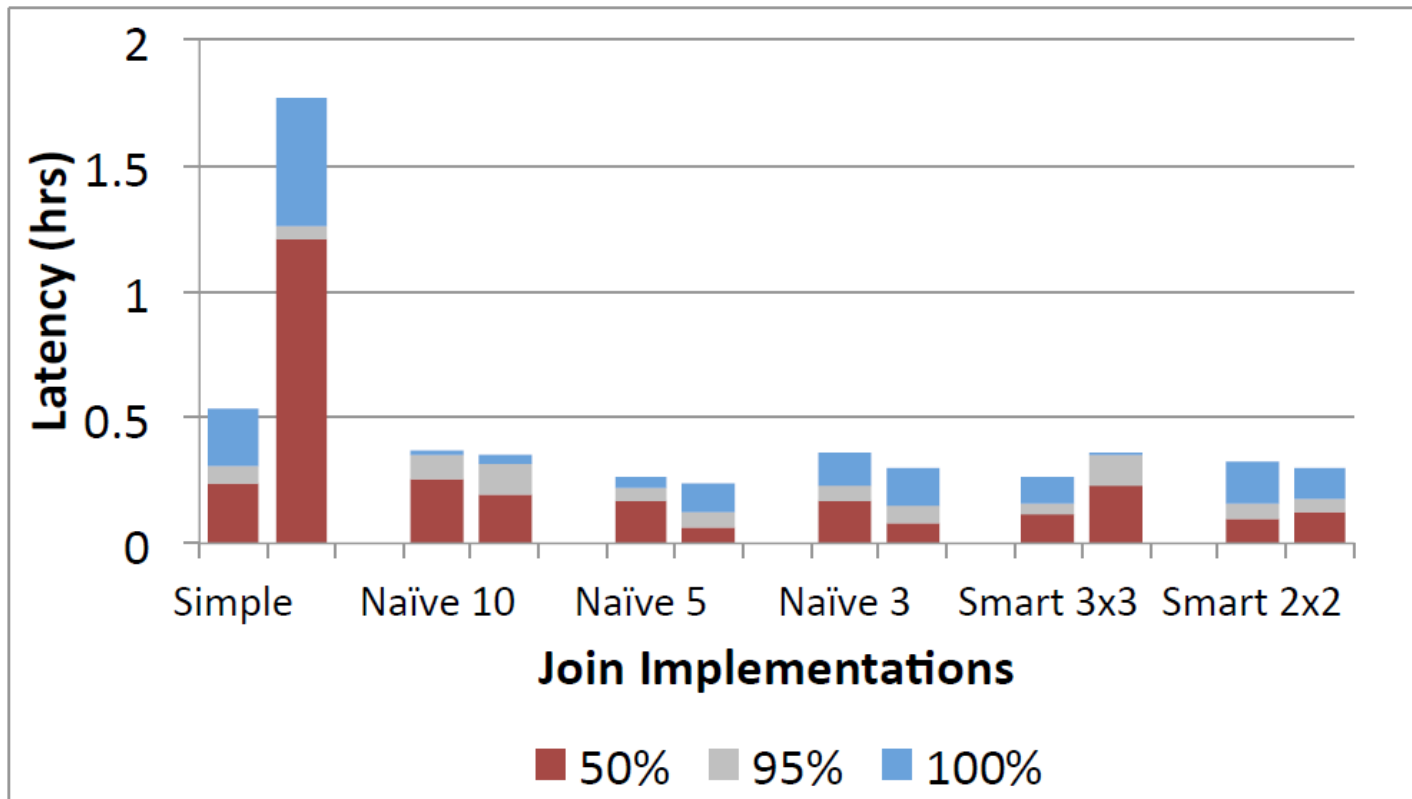
Experiments



Fraction of correct answers on celebrity join for different batching approaches

Join Operator

Experiments



Completion time in hours for variants of celebrity join on two tables

Sort Operator

- Sort application with human intelligence

The image shows two overlapping web pages. The background page is a Yelp profile for 'Club Ilia', which is marked as 'Claimed' and has 21 reviews. It lists 'Italian, Burgers, Pizza' and shows a map location at 8902 University High Street, Burnaby, BC V5A 4Y6. The foreground page is an Amazon product listing for 'Fundamentals of Database Systems, Global Edition'. It features a star rating of 3.4 out of 5 stars based on 15 customer reviews. A detailed breakdown of the star ratings is shown in a pop-up:

Star Rating	Percentage
5 star	40%
4 star	14%
3 star	13%
2 star	13%
1 star	20%

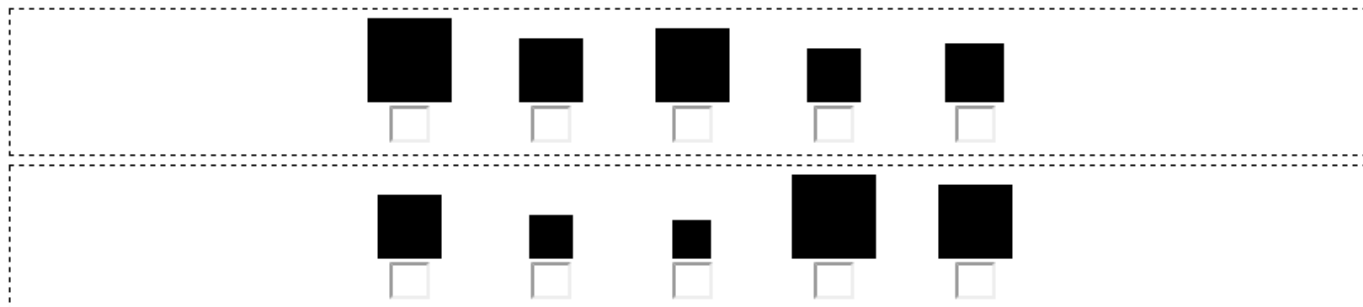
Below the star breakdown, there is a link to 'See all verified purchase reviews'. The Amazon page also includes sections for 'Refine by', 'Shipping' (with a 'Free shipping' checkbox), 'Condition' (with a 'Learn more' link), and 'Delivery' (with a note 'Arrives between 13.' and a link to 'Shipping rates').

Sort Operator

● Comparison Sort

There are 2 groups of squares. We want to order the squares in each group from smallest to largest.

- Each group is surrounded by a dotted line. Only compare the squares within a group.
- Within each group, assign a number from 1 to 7 to each square, so that:
 - 1 represents the smallest square, and 7 represents the largest.
 - We do not care about the specific value of each square, only the relative order of the squares.
 - Some groups may have less than 7 squares. That is OK: use less than 7 numbers, and make sure they are ordered according to size.
 - If two squares in a group are the same size, you should assign them the same number.



Submit

Sort Operator

● Rating Sort

There are 2 squares below. We want to rate squares by their size.

- For each square, assign it a number from 1 (smallest) to 7 (largest) indicating its size.
- For perspective, here is a small number of other randomly picked squares:



<div><div></div><div>smallest 1 2 3 4 5 6 7 largest</div></div>
<div><div></div><div>smallest 1 2 3 4 5 6 7 largest</div></div>

Submit

Hybrid Schemes

- Initially use the rating-based sort
- Use comparisons, in one of three flavors:
 - **Random**: pick S items, compare
 - **Confidence-based**: pick most confusing “window”, compare
 - **Sliding-window**: for all windows, compare

Sort Operator

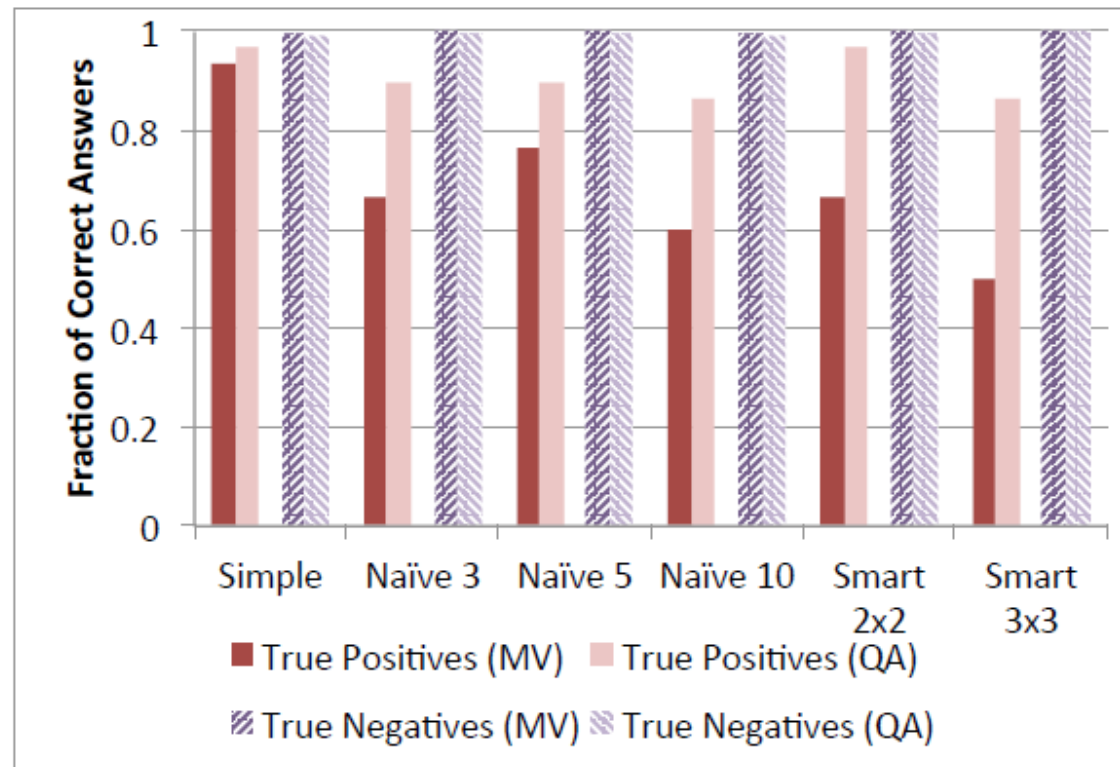
- Experiments setup:
 - Compare *IDEAL*, *Simple*, *Naive* and *Smart* schemes
 - Dataset: *celebrity join dataset*

Implementation	True Pos. (MV)	True Pos. (QA)	True Neg (MV)	True Neg (QA)
IDEAL	20	20	380	380
Simple	19	20	379	376
Naive	19	19	380	379
Smart	20	20	380	379

Baseline comparison of three join algorithms with no batching enabled

Sort Operator

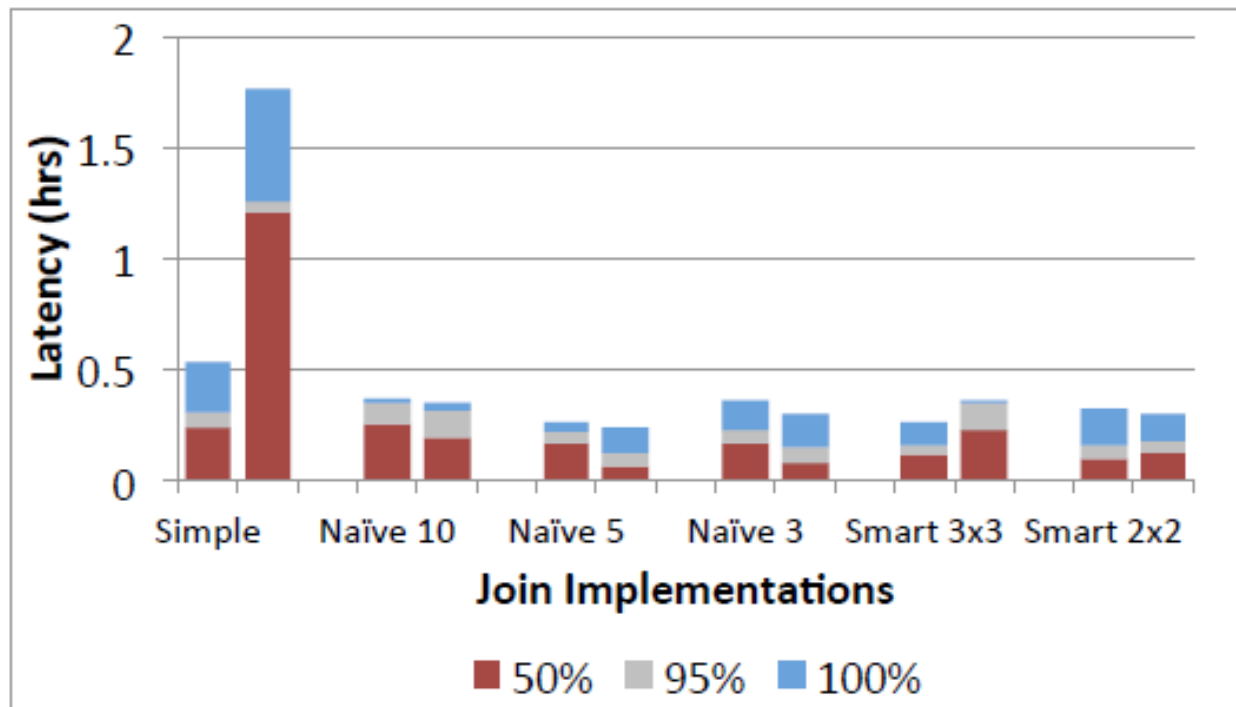
Experiments



Fraction of correct answers on celebrity join for different batching approaches

Sort Operator

- Experiments



Completion time in hours for variants of celebrity join on two tables

Conclusion

- A declarative workflow engine, Qurk
- Integrated crowd intelligence
- Sorting and joining operation
- Optimization in task batching
- Reduce the overall cost from \$67 to \$3

Questions and Comments?

Thank you!