# Data Management Research @ UW Seattle



uwdb.io

Magdalena Balazinska

Alvin Cheung

Dan Suciu

## UW Database Group

Data management systems, cloud services, probabilistic databases, and data pricing in Computer Science & Engineering at the University of Washington in Seattle.
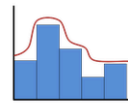
### Current Projects

**MYRIA**
Big Data as a Service

**VISUALCLOUD**
A DBMS for Virtual Reality

**ENTROPYDB**
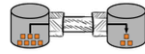EntropyDB for Data Exploration

**SQLSHARE**
Database-as-a-Service for High-Variety Data

**QURO**
Query reordering in OLTP transactions

**PIPEGEN**
Data Pipe Generation for Hybrid Analytics

**COSETTE**
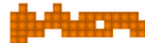An Automated SQL Solver

**ZALIQL**
A Declarative Framework for Drawing Causal Inference from Big Data

**LARA**
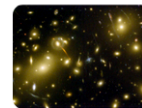A Key-Value Algebra underlying Arrays and Relations

**UW BRANCH OF SCIDB**
Parallel distributed array database engine

**DATA ECO\$Y\$TEM**
Data management and pricing in the cloud

**ASTRODB**
An inter-disciplinary collaboration for new methods and tools for Big Data Astronomy

## W PAUL G. ALLEN SCHOOL
### OF COMPUTER SCIENCE & ENGINEERING
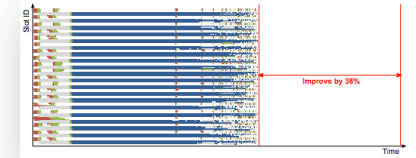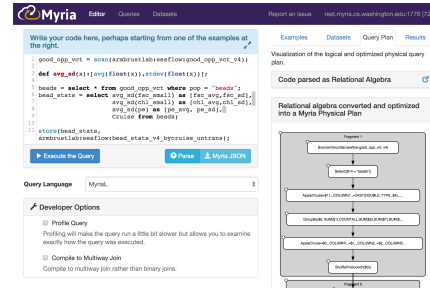
`http://uwdb.io/`

**Research in database systems, theory, and programming languages**

~15 students + postdocs

# Research Areas



## Big data processing in the cloud

- **Theory:** optimal query processing **P** Walter Cai
- **Systems**: Myria, efficient & complex processing at scale, image analytics, DBMS+NN, data summarization **P** Jenny Ortiz
- **Usability**: Cloud SLAs, performance tuning, viz analytics **P** Leilani Battle

## New Types of DBMSs

- Open World DBMS
- Image & video DBMS
- LightDB: VR/AR/MR DBMS **P** Brandon Haynes

## Scientific data management

- Collaborations with scientists & deep involvement with eScience Institute

## Databases and programming languages

- DBMS & app co-optimization

## Probabilistic Databases **P** Laurel Orr

## Causality

# Towards Application-Specific Databases

uwdb.io

PAUL G. ALLEN SCHOOL
OF COMPUTER SCIENCE & ENGINEERING

PLSE
uwplse.org

Scidb — Scientific Workloads

Column Stores — OLAP

Storm — Streams

Main Memory DB — OLTP

SparkSQL — Analytics

Specialization

WHY SO MANY STORES?

Can we generate customized data stores from application code?

Cong Yan

## Application Inefficiencies

- Code translated to inefficient queries
- Misplaced computation
- Redundant data loads
- Issuing queries with known results
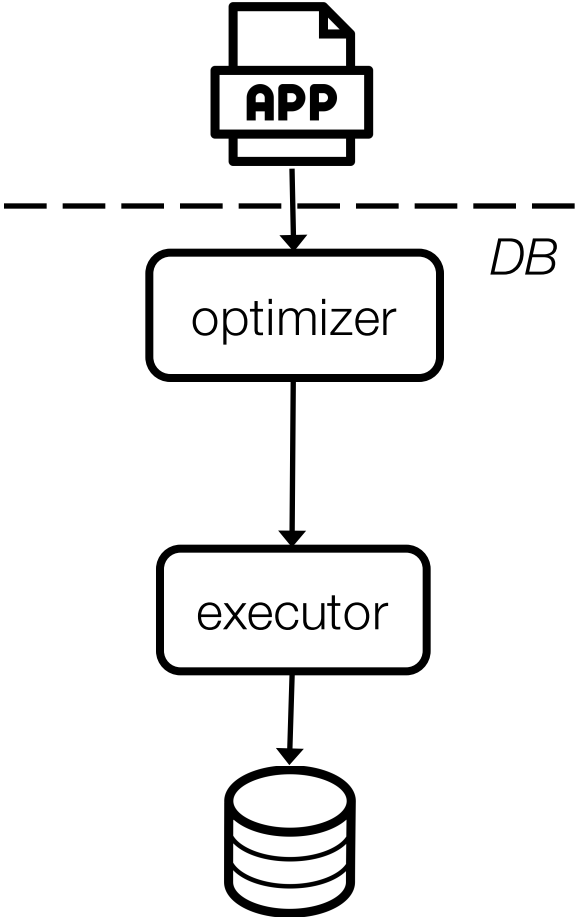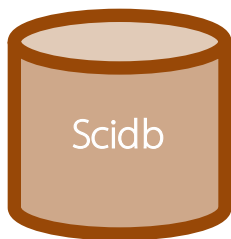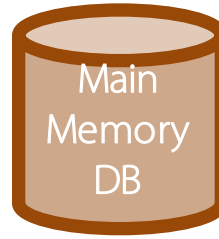- Loading unused data
- Missing indexes

**78% of fixes took fewer than 5 lines**
**Max app speedup: 39x**

| # stars | Application | # issues |
|---------|-------------|----------|
| 22k | Discourse (forum) | 85 |
| 1k | Lobster (forum) | 45 |
| 49k | Gitlab (collaboration) | 23 |
| 13k | Redmine (collaboration) | 59 |
| 17k | Spree (E-commerce) | 20 |
| 1.7k | ROR Ecommerce | 11 |
| 697 | Fulcrum (task mgmt) | 2 |
| 3.5k | Tracks (task mgmt) | 30 |
| 18k | Diaspora (social network) | 57 |
| 1.2k | Onebody (social network) | 76 |
| 8k | Openstreetmap (map) | 4 |
| 1.1k | Fallingfruit (map) | 16 |
| **Total** | | **428** |

| Image Rotate | → | Blur |

| Hash Partitioning | → | Join |

# SEARCH

Target code                    Proof of translation

SEARCH

# PROGRAM SYNTHESIS

Target code                    Proof of translation

# Verified Lifting: Casper

Maaz Ahmad

### 1. Define semantics of map and reduce

```
SumXY  = reduce(map(points, fm),
fr)
fm(x,y)   = x * y
fr(v1,v2) = v1 + v2
```

### 3. Retarget spec to Hadoop

codegen

```
void map(Object key, Point [] value)
{   for(Point p : points)
      emit("sumxy", SumXY); }
void reduce(Text key, int [] vs)
{   int SumXY = 0;
    for (Integer val : vs)
      SumXY = SumXY + val;
    emit(key, SumXY); }
```

### 2. Synthesizer infers spec from source

```
// sequential implementation
void regress(Point [] points)
{
    int SumXY = 0;
    for(Point p : points){
     SumXY += p.x * p.y;
    }
    return SumXY;
}
```

Lifted code can be optimized by Hadoop max 32x speedup

Q1

Q2

$\forall$ D . Q1(D) = Q2(D)
$\exists$ D . Q1(D) $\neq$ Q2(D) **?**

Query Optimizers

Autograders

Application Caches

Boris Trakhtenbrot

Deciding the equality of two arbitrary relational queries is undecidable.

Full decision procedure exists for conjunctive queries

Simple heuristics can already prove many common cases

**Repeat**

HTML Data

Images

Regex → Filter → Join

CNN
Conv → ... → Conv

RNN

Output Model

**Generate Training Labels**

**Train a caption-generating model**

Many regex and join algorithms to choose from!

Likewise for convolution

# Cuttlefish: A Lightweight Primitive for Online Tuning

Tomer Kaftan

```
def loopConvolve(image, filters): ...
def fftConvolve(image, filters): ...
def mmConvolve(image, filters): ...


for image, filters in convolutions:

  start = now()
  result = convolve(image, filters)
  elapsedTime = now() - start

  output result, elapsedTime
```

# Cuttlefish: A Lightweight Primitive for Online Tuning

Tomer Kaftan

```python
def loopConvolve(image, filters): ...
def fftConvolve(image, filters): ...
def mmConvolve(image, filters): ...
tuner = Tuner([loopConvolve, fftConvolve, mmConvolve])

for image, filters in convolutions:

  start = now()
  result = convolve(image, filters)
  elapsedTime = now() - start

  output result, elapsedTime
```

# Cuttlefish: A Lightweight Primitive for Online Tuning

Tomer Kaftan

```python
def loopConvolve(image, filters): ...
def fftConvolve(image, filters): ...
def mmConvolve(image, filters): ...
tuner = Tuner([loopConvolve, fftConvolve, mmConvolve])

for image, filters in convolutions:
    convolve, token = tuner.choose()
    start = now()
    result = convolve(image, filters)
    elapsedTime = now() - start

    output result, elapsedTime
```
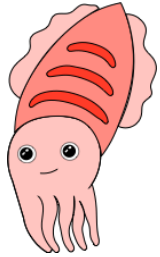
# Cuttlefish: A Lightweight Primitive for Online Tuning
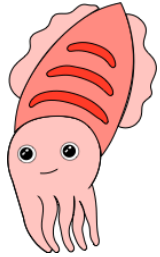
Tomer Kaftan

```python
def loopConvolve(image, filters): ...
def fftConvolve(image, filters): ...
def mmConvolve(image, filters): ...
tuner = Tuner([loopConvolve, fftConvolve, mmConvolve])

for image, filters in convolutions:
    convolve, token = tuner.choose()
    start = now()
    result = convolve(image, filters)
    elapsedTime = now() - start
    tuner.observe(token, elapsedTime)
    output result, elapsedTime
```
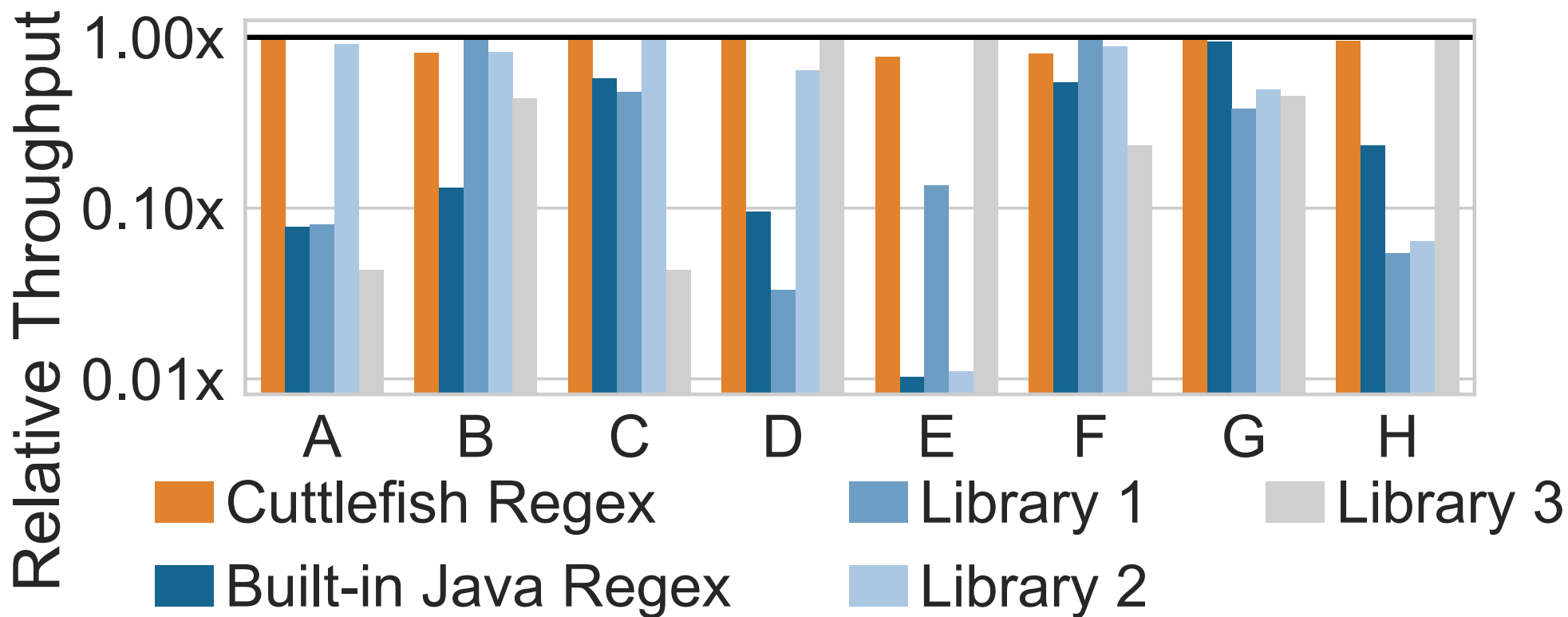
Note: Y-axis is Log-scale

# Scythe

Chenglong Wang

Input tables

| id | date |
|----|-------|
| 1 | 12/25 |
| 2 | 11/21 |
| 4 | 12/24 |
| ... | ... |

Output tables

| id | date | max |
|----|-------|-----|
| 1 | 12/25 | 30 |
| 2 | 11/21 | 10 |
| 4 | 12/24 | 20 |
| ... | ... | ... |

Search for abstract queries

```
Select *
From    (Select *
         From T1
         Where □)
Join    (Select id,
                 Max(val)
         From   T2
         Where □
         Group By oid
         Having □) T3
On            □
```

Instantiate abstract queries

Stored using specialized data structures

```
Select *
From (Select *
      From T1
      Where True)
Join (Select id,
             Max(val)
      From T2
      Where val < 50
      Group By oid
      Having True) T3
On T3.oid = T1.uid
```

Prune query skeletons

Rank results based on simplicity

# Scythe

Chenglong Wang

Supported features
- SPJ
- Grouping
- Aggregation
- Subqueries
- Outer join
- Exists
- Union

Benchmark: 193

Scythe: 143

Enum: 92

59% can be answered within 20 seconds

34x faster on avg.

Scythe
Enum

Time (seconds)

# Is there something equivalent to argmax in SQL?

In a more general sense: is there a function that will allow me to find the entire row where a value in Column X is the max value of the column?

16

sql

share improve this question
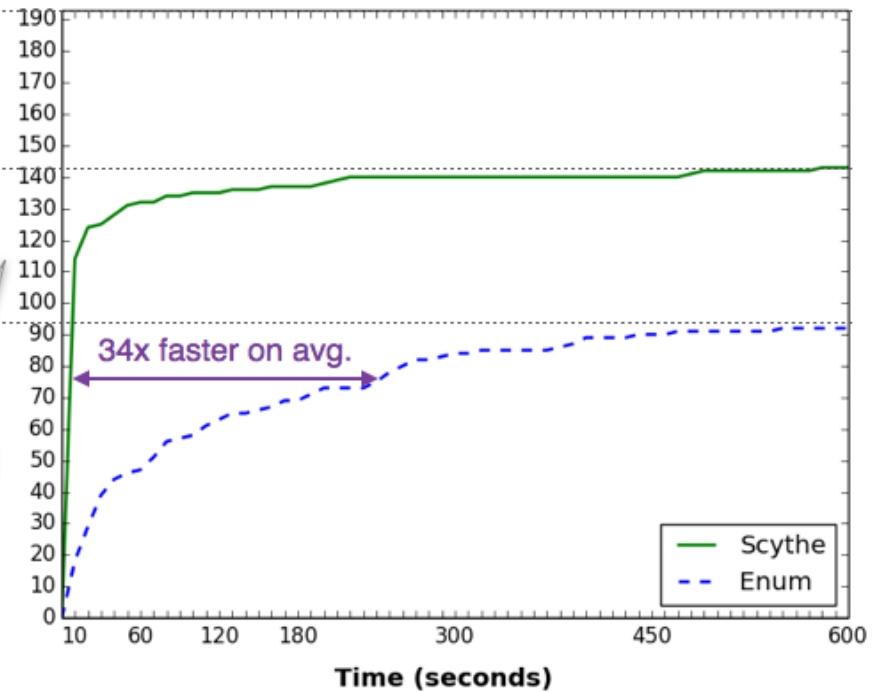
7

If I'm reading your question correctly, the following query should do it (assuming that our column names are a,b, and c and that a is the column that we're maximizing):

```
select a,b,c
from table
where a=(select max(a) from table);
```

Of course, if you have more than one row where the column a attains its maximum, then you'll get more than one row back from the query. If you want a unique row back, you can add something like "order by b,c limit 1", or use some other way to rank the rows in which a attains its max.

Titles summarize post 80% of the time
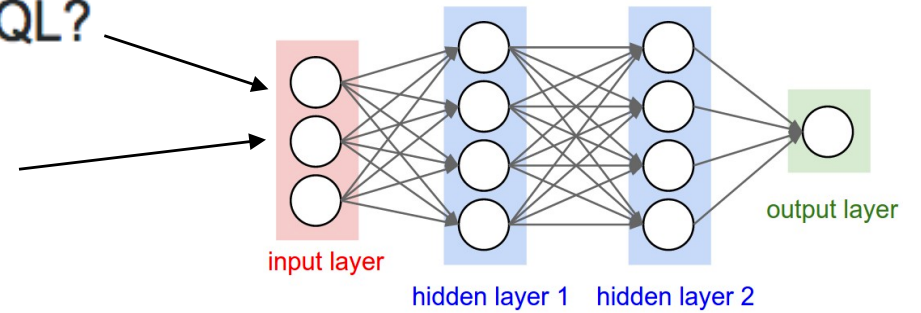
Stackoverflow dataset
- Posts tagged with #sql, #oracle, #database (430k)
- Posts containing an accepted answer in SQL
- Results: 41k (title, query) pairs

Filtered away titles
- My query doesn't work!
- Why is my query slow?
- I hate SQL!

Is there something equivalent to argmax in SQL?

```
select a,b,c
from table
where a=(select max(a) from table);
```



input layer
hidden layer 1   hidden layer 2
output layer

| Model | Naturalness | Informativeness |
|-------|-------------|-----------------|
| Code-NN (Ours) | 2.6 | 1.55 |
| Nearest neighbor | 1.9 | 1.55 |
| MOSES | 1.76 | 1.36 |
| ATTEN | 2.82 | 0.93 |

Srini Iyer

# UWDB Collaborators

**UW**

- Bill Howe (iSchool)
- Andrew Connolly (Astronomy)
- Aaron Lee (Ophtalmology)
- Ariel Rokem (eScience)
- Emilio Zagheni (Sociology)
- Prog Lang & SW Eng group

**Industry**

- Adobe
- Huawei
- Intel
- Microsoft
- Teradata