

“Survivorship bias?
in my database workload?”

It's more likely than you think.

FREE DB CHECK!



CONTENTWATCH

Ryan Marcus, Jeffrey Tao, Peizhi Wu, Zijie Zhao
University of Pennsylvania

<https://rm.cab/cidr26>

Benchmarks are Key

- How do you evaluate a query optimizer, execution engine, scheduler, cache, etc.?
 - Option 1: Theory
 - Option 2: Benchmarks
- The benchmark is only good when it is *representative*
 - The key assumption: if I improve the benchmark, I'll improve the user's workload!

TPC-*

Join Order Benchmark,
DSB

SQLStorm, SQLBarber,
Adversarial Query Gen



Actual Workload



Survivorship Bias

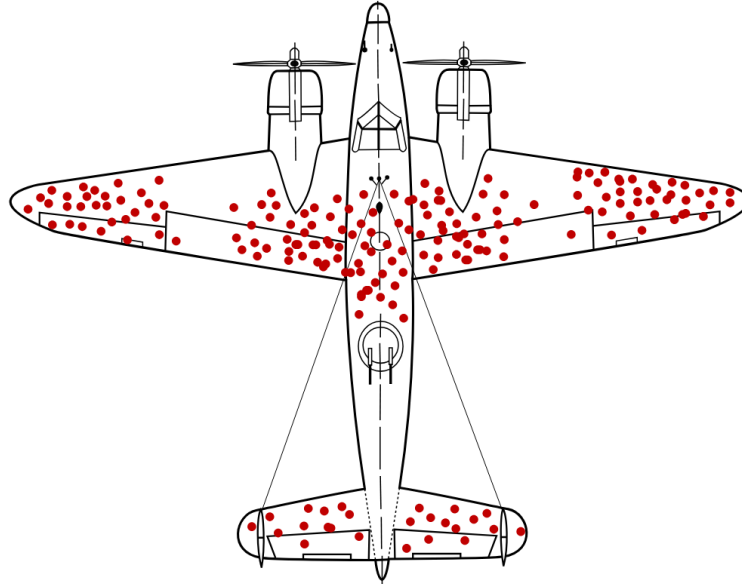


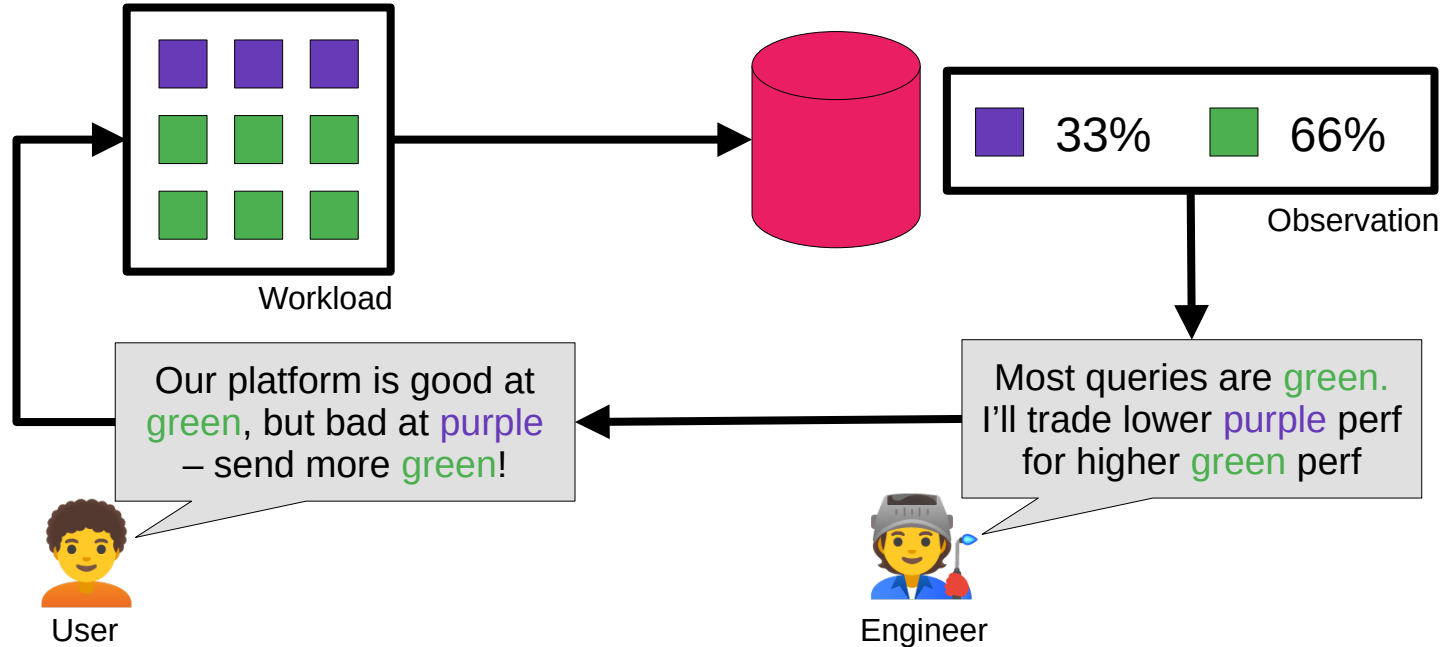
Figure 1: A diagram of bullet holes on returning WWII fighter planes. Since these planes all returned, the sample is biased towards “survivors.” One might think that additional armor should be added to the shot areas, but a more effective strategy is to add armor to the areas that appear unscathed.

This Talk

- **A theory of survivorship bias in workloads**
- Some tests of that theory
- Analysis and insights

① User submits their workload to the system

② Engineers observe properties of the workload



④ Users optimize their workloads based on their platform

③ Engineers identify hotspots and optimize their system

**The queries you run impact what the engineers build.
What the engineers build impacts the queries you run!**

This Talk

- A theory of survivorship bias in workloads
- **Some tests of that theory**
- Analysis and insights Does this loop actually happen?

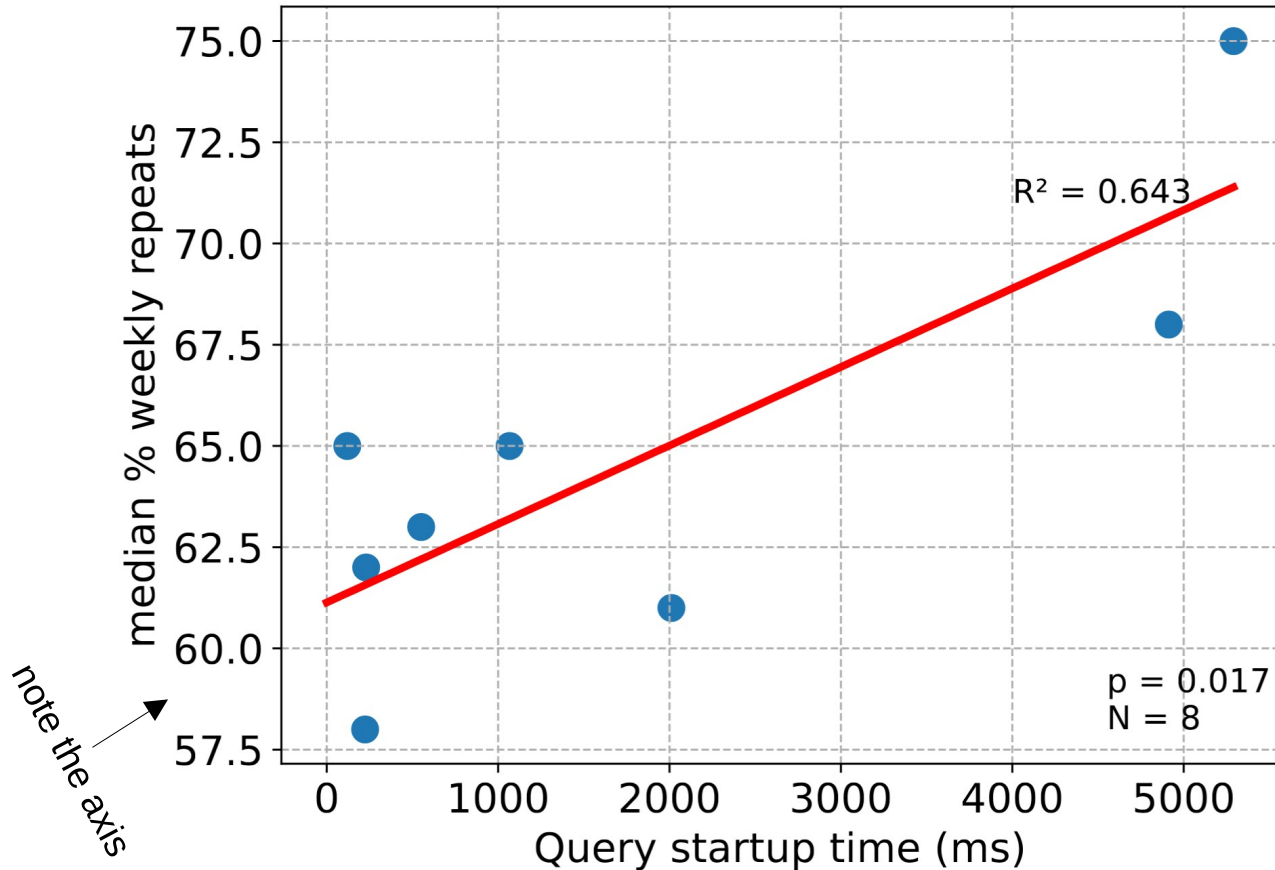
In a perfect world...

- A vendor could artificially slow down a certain type of query to see if their customers stop issuing them.
- A vendor could invest time optimizing something their customers don't currently do.
- Seeking volunteers.

We'll settle for...

- Natural experiments!
- Vendors that are better at processing ad-hoc queries should see more of them.
- A vendor that introduces a fast feature should see uptake from users

Exp. 1: Repetitive Queries

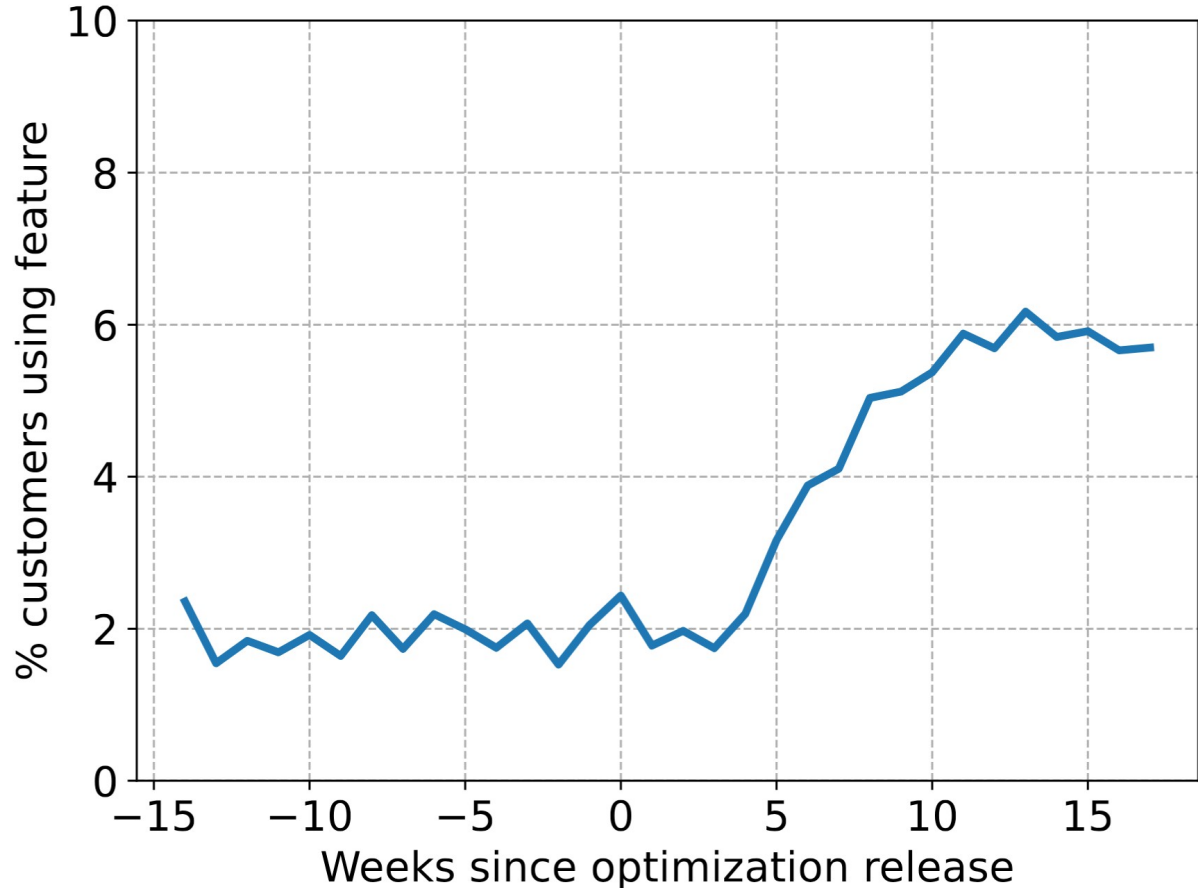


**Cold query startup time vs.
% repeated queries**

Data platforms with longer
query startup see more
repetitive queries.

Could be platform choice or
optimization over time.

Exp. 2: Build it and they will come



% users of a particular feature after optimization

Only 2% usage before optimization, but 6% usage afterwards!

Could be “release notes” teaching about the feature

Survivorship Bias

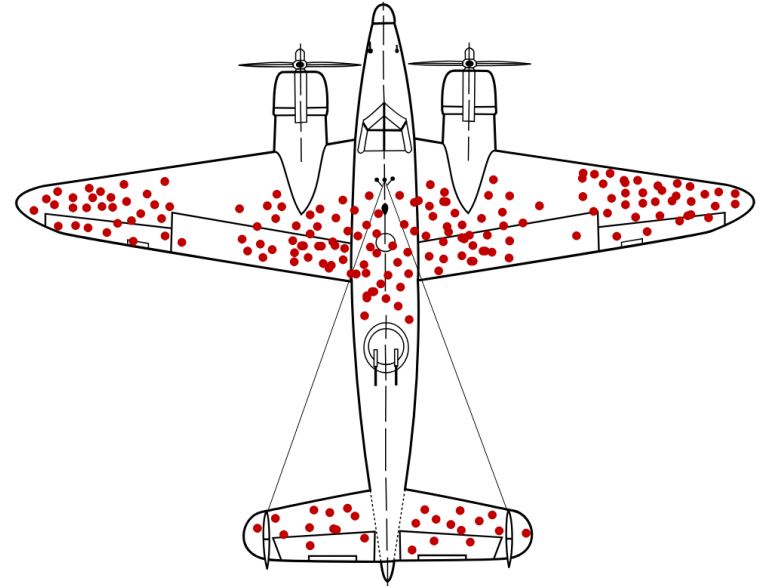
- The DBMS influences what queries the user sends
 - Users use fast features, don't use slow ones
- The user's queries influence the DBMS
 - Engineers optimize for what the user sends
- *Workload traces only capture queries that “survive” this negotiation!*

This Talk

- A theory of survivorship bias in workloads
- Some tests of that theory
- **Analysis and insights**

Are we doing it wrong?

- By optimizing existing user workloads, are we putting armor on the wrong part of the plane?
- A customer continuously issuing a query is a certificate of “good enough.”



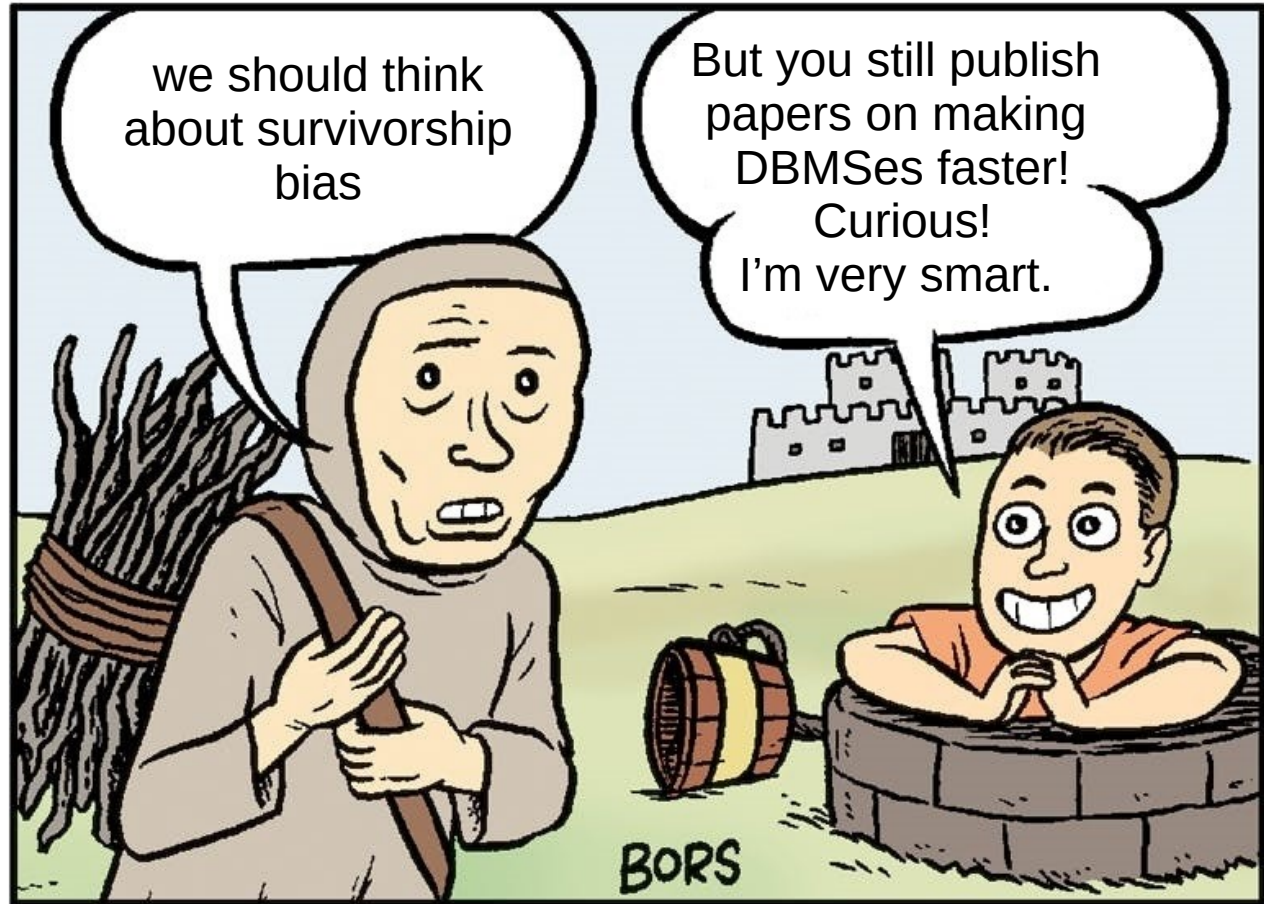
Thought experiment

- Imagine column stores were never invented, and you run a cloud database.
- How many long, OLAP style queries do you think you'd see compared to Redshift?



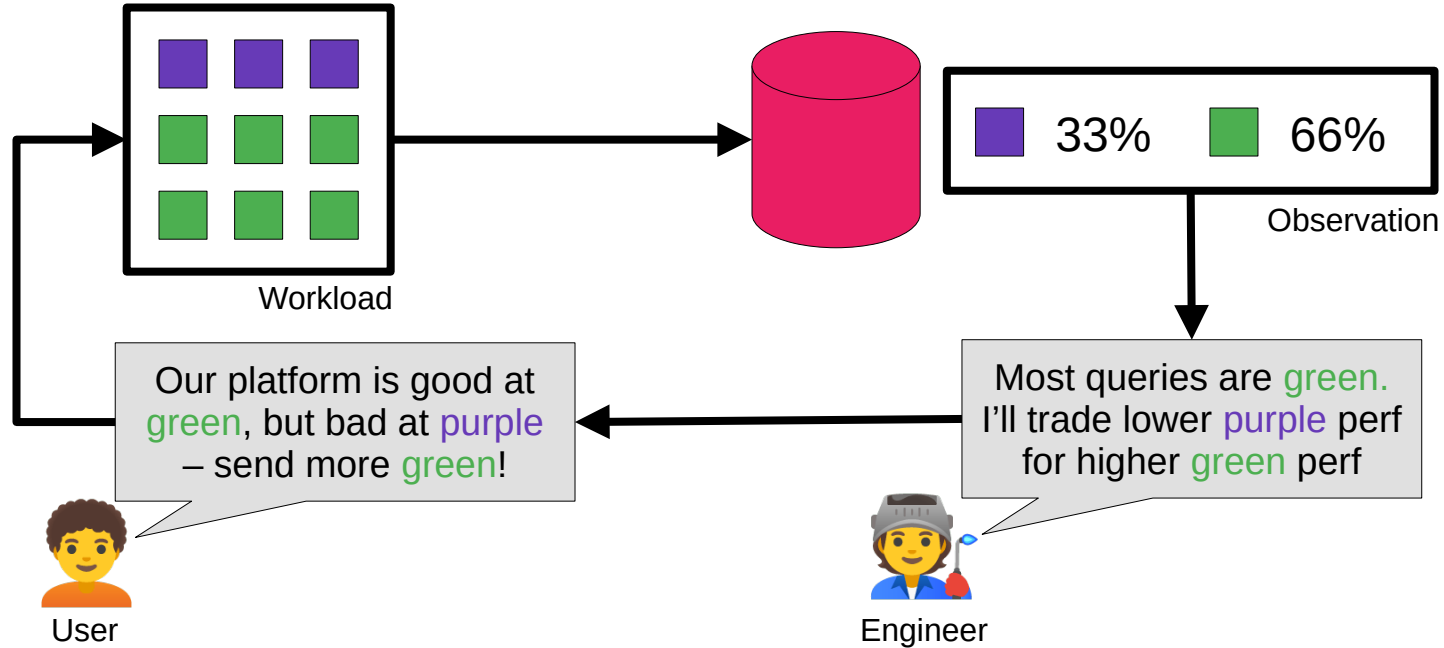
Faster DBMSes are still good

- Faster DBMSes *are still good*
- ... basing that opt work on *real traces* is **WAY** better than making stuff up



① User submits their workload to the system

② Engineers observe properties of the workload



④ Users optimize their workloads based on their platform

③ Engineers identify hotspots and optimize their system

How can we optimize the query the user *wants* to run, but can't?