



Matt Butrovich Carnegie Mellon University

Tastes Great! Less Filling!

High Performance and Accurate
Training Data Collection for Self-
Driving Database Management
Systems

Training Data for Self-Driving DBMSs

What is a self-driving DBMS?

Goal: Automate onerous tuning and optimization tasks for DBMSs.

Given an objective (e.g., throughput, latency) a *self-driving DBMS* deploys **actions** that it deems will help the application's future performance for that objective.

Actions control:

- Physical design
- Knob configuration
- Hardware resources

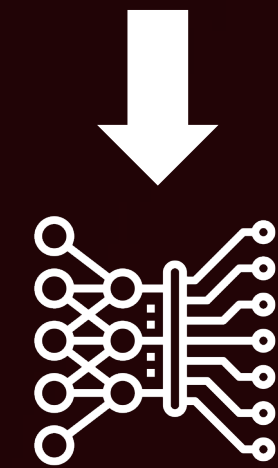


Behavior Models

```
SELECT * FROM foo WHERE balance < 500;
```

Input Features

Operation	Relation	# Filters	Execution Mode	Cost	...
Sequential scan	13	1	Compiled	15445.0	...



Sequential Scan Behavior Model

Output Metrics

CPU Cycles	Memory Bytes	Network Bytes Read	Disk Bytes Read	...	Elapsed Time
12131989	1208640	0	65536	...	35419



Input Sources

External features:

- Execute SQL queries (e.g., EXPLAIN) or other public APIs.
- QPPNet (Marcus et al., VLDB 2019)

Internal features:

- Modify DBMS source code to capture state.
- MB2 (Ma et al., SIGMOD 2021)



Output Sources

User-space metrics:

- Operating system APIs (e.g., perf, getrusage)
- Scrape kernel file system (e.g., /proc)

Kernel-space metrics:

- Kernel data structures and privileged APIs.
- Efficient RCU-synchronized data structures.

BPF

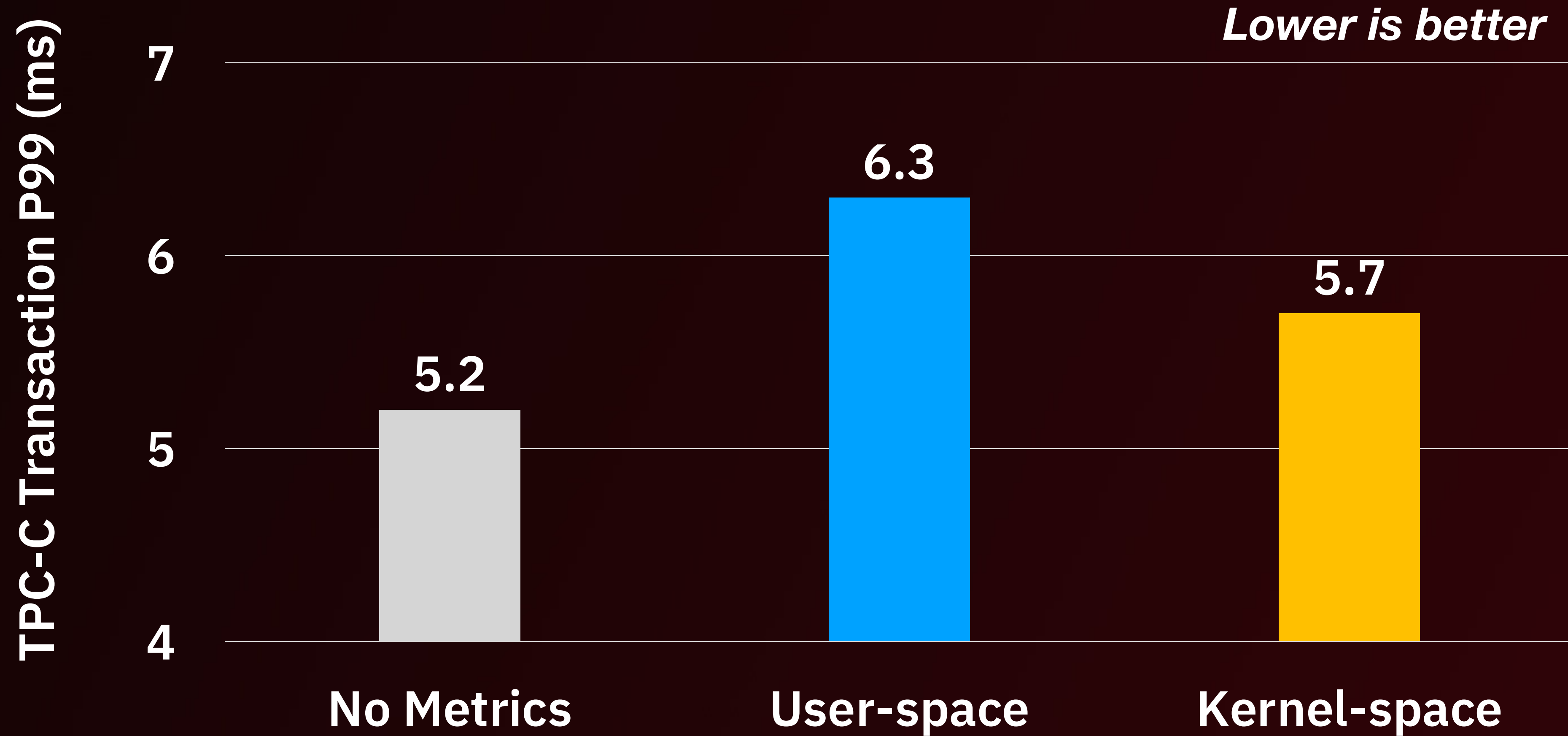
Berkeley Packet Filter in 1992, *Extended Berkeley Packet Filter* (eBPF) since 2014, but we'll just say *BPF*.

VM to run code in privileged kernel mode without writing kernel modules.

Strict constraints:

- Number of instructions
- Boundedness
- Memory safety

Metrics Collection Overhead



Training Data Wish List

DBMS Internal Features

- More information about current operation.
- Learn interactions with background tasks.

Kernel-space Metrics

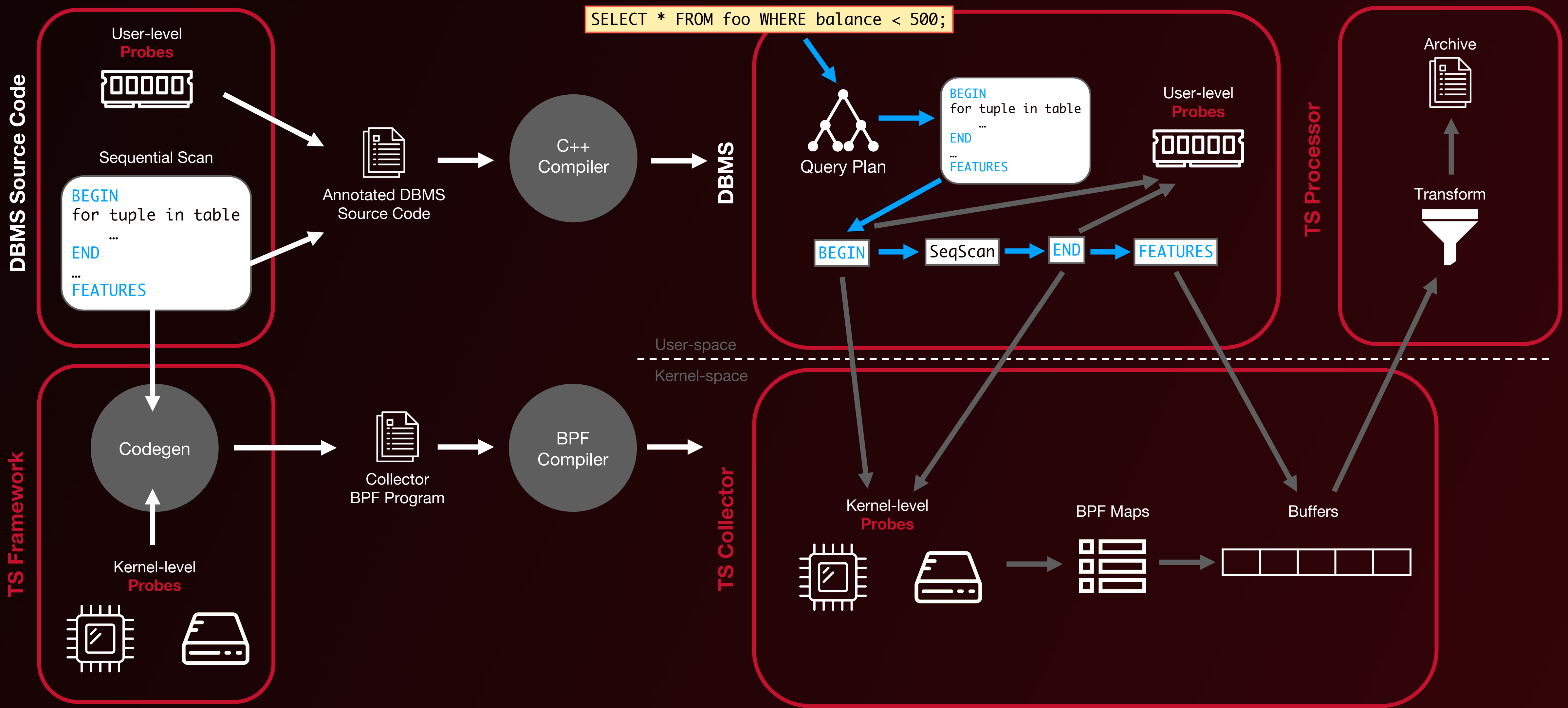
- Low overhead instrumentation.
- Inspect kernel data structures.

Online Environments

- Train on the target workload.
- Learn about deployment hardware.

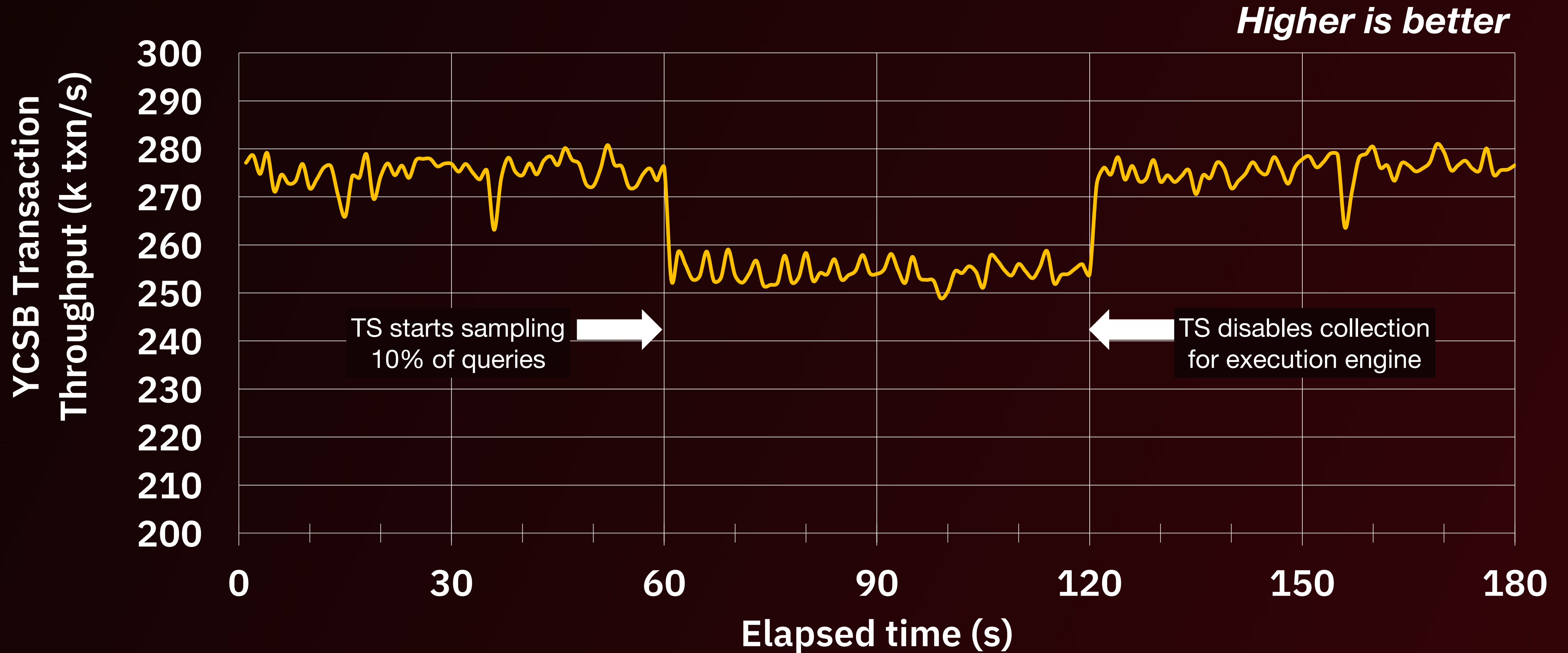
TScout Training Data Collection Framework

TScout Workflow

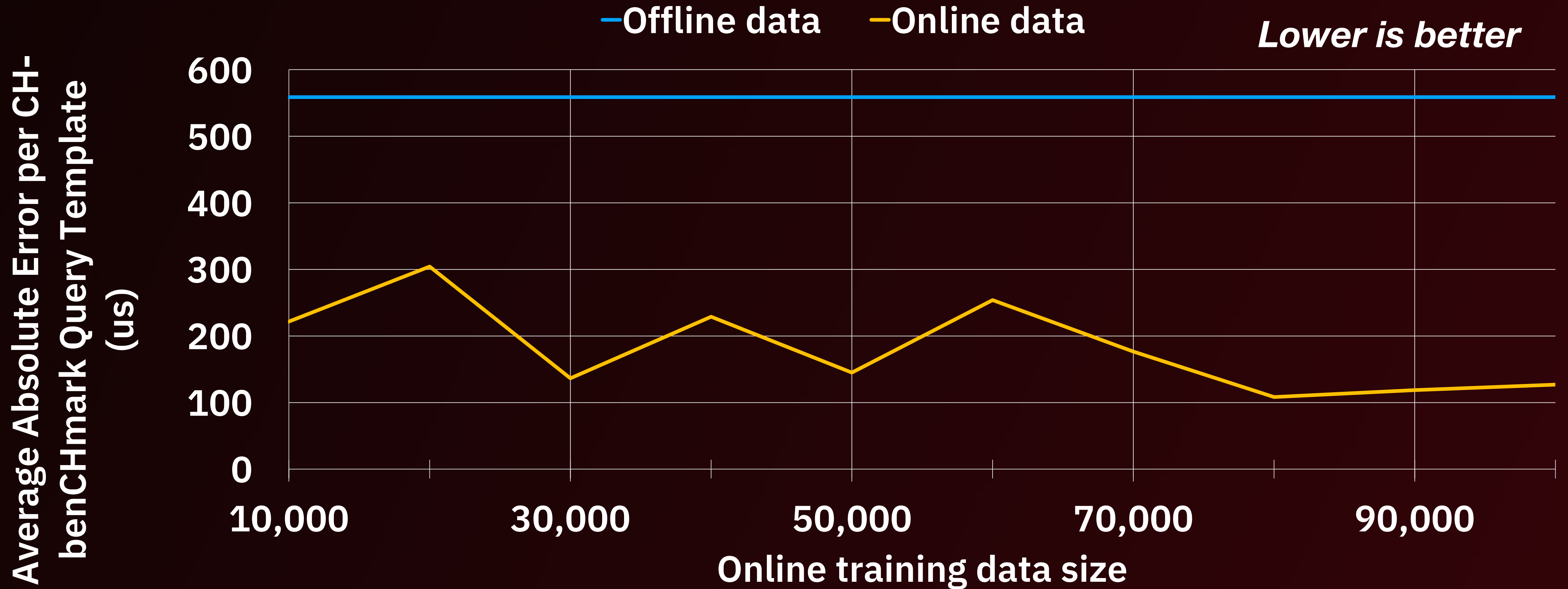


TScout Performance

High Performance Data Collection



Online Data Benefits



Training Data Wish List Revisited

✓ DBMS Internal Features

- TScout can read DBMS memory.

✓ Kernel-space Metrics

- TScout's Collector reduces round trips to kernel.
- TScout's Processor moves training data off critical path.

✓ Online Environments

- Low overhead data generation.
- Adjustable sampling.

END

<http://mattbutrovi.ch>

<https://noise.page>