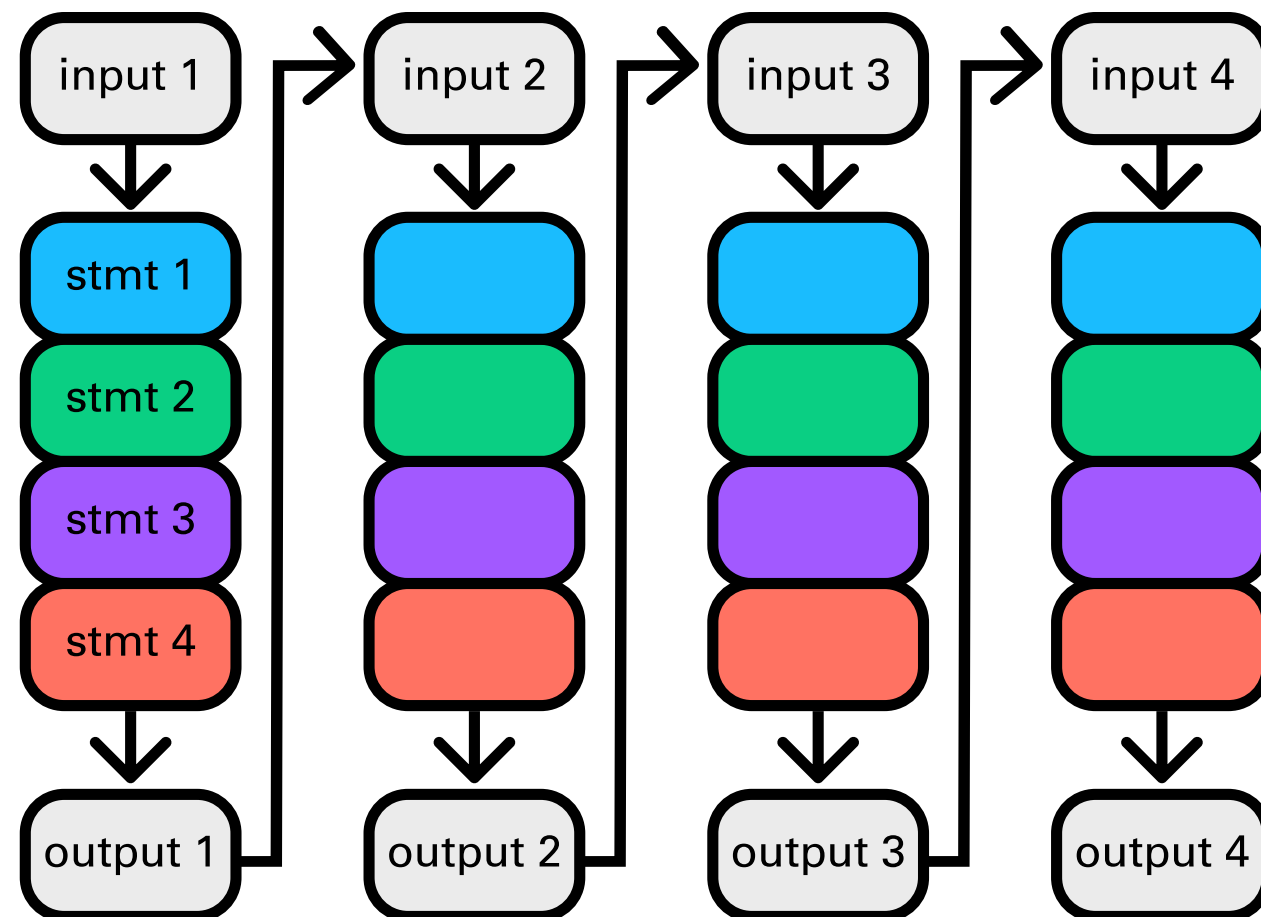# Efficiently Executing UDFs via Batching
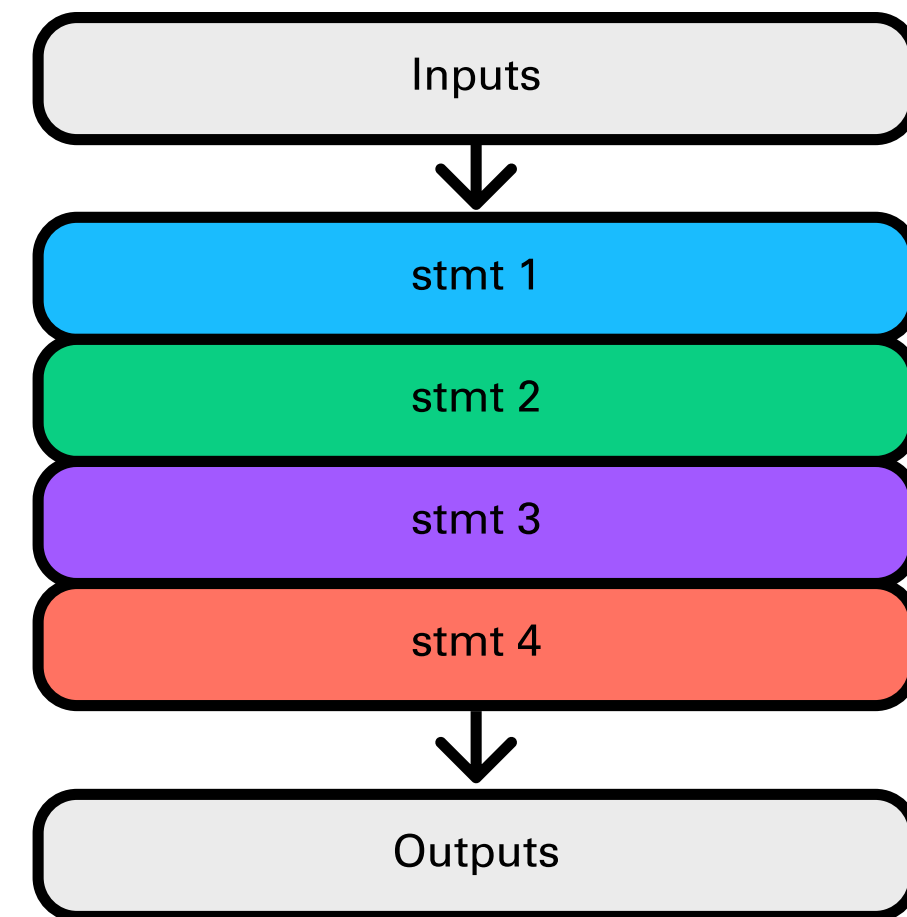
Kai Franz

# Batching UDFs

- Naive UDF execution:
  process one row at a time

- Batching: apply each statement to all rows at the same time

- DBMS is good at set-oriented processing

# Parameters & Local Variables

- Create temp table with a column for each variable

```
CREATE OR REPLACE FUNCTION maxPurchaseChannel(ckey INT, fromDateSk INT, toDateSk INT)
    RETURNS VARCHAR(50)
    LANGUAGE plpgsql
AS
$$
DECLARE
    numSalesFromStore   INT;
    numSalesFromCatalog INT;
    numSalesFromWeb     INT;
    maxChannel          VARCHAR(50);

...
```

```
CREATE TABLE #temp
(
    ckey                INTEGER,
    fromDateSk          INTEGER,
    toDateSk            INTEGER,
    numSalesFromStore   INTEGER,
    numSalesFromCatalog INTEGER,
    numSalesFromWeb     INTEGER,
    maxChannel          VARCHAR(50),
    maxPurchaseChannel  VARCHAR(50),
);
```

# Writing to Local Variables

- SET, SELECT INTO, etc. become UPDATE

```
SET spending =
        (SELECT SUM(ws_net_paid_inc_ship_tax)
            FROM web_sales_history,
                 date_dim
          WHERE d_date_sk = ws_sold_date_sk
            AND d_year = 2000
            AND ws_bill_customer_sk = cust_sk);
```

→

```
UPDATE #temp
   SET spending =
           (SELECT SUM(ws_net_paid_inc_ship_tax)
               FROM web_sales_history,
                    date_dim
             WHERE d_date_sk = ws_sold_date_sk
               AND d_year = 2000
               AND ws_bill_customer_sk = cust_sk);
```

# Branching

- Froid-style boolean column for each branch predicate
- Statements inside branch guarded with WHERE

```
IF (numSalesFromStore > numSalesFromCatalog) THEN
        maxChannel := 'Store';
END IF;
```

→

```
CREATE TABLE #temp
(
    ...
    p1                 BIT,
    p2                 BIT,
    p3                 BIT,
    ...
);

UPDATE #temp
    SET p1 = CASE WHEN numSalesFromStore > numSalesFromCatalog THEN 1 ELSE 0 END;

UPDATE #temp
    SET maxChannel = 'Store'
 WHERE p1 = 1;
```

# Return Statements

- Boolean column tracks if row has returned yet

```
IF netProfit > 0 THEN
        RETURN 1;
ELSE
        RETURN 0;
END IF;
```

→

```
UPDATE #temp
    SET profitablemanager = 1,
        returned          = 1
  WHERE p1 = 1
    AND returned = 0;

UPDATE #temp
    SET profitablemanager = 0,
        returned          = 1
  WHERE p1 = 0
    AND returned = 0;
```

# Experiments

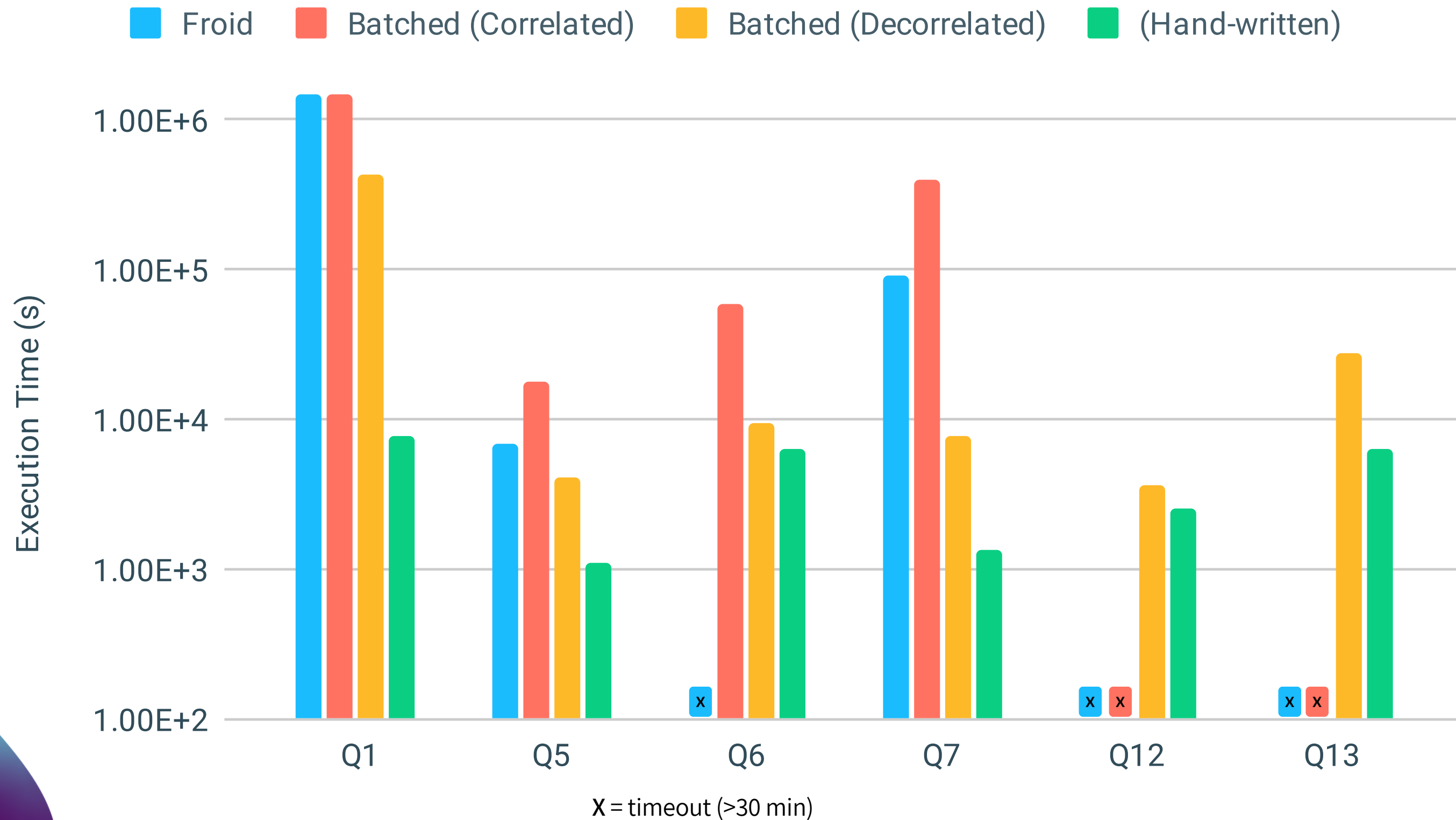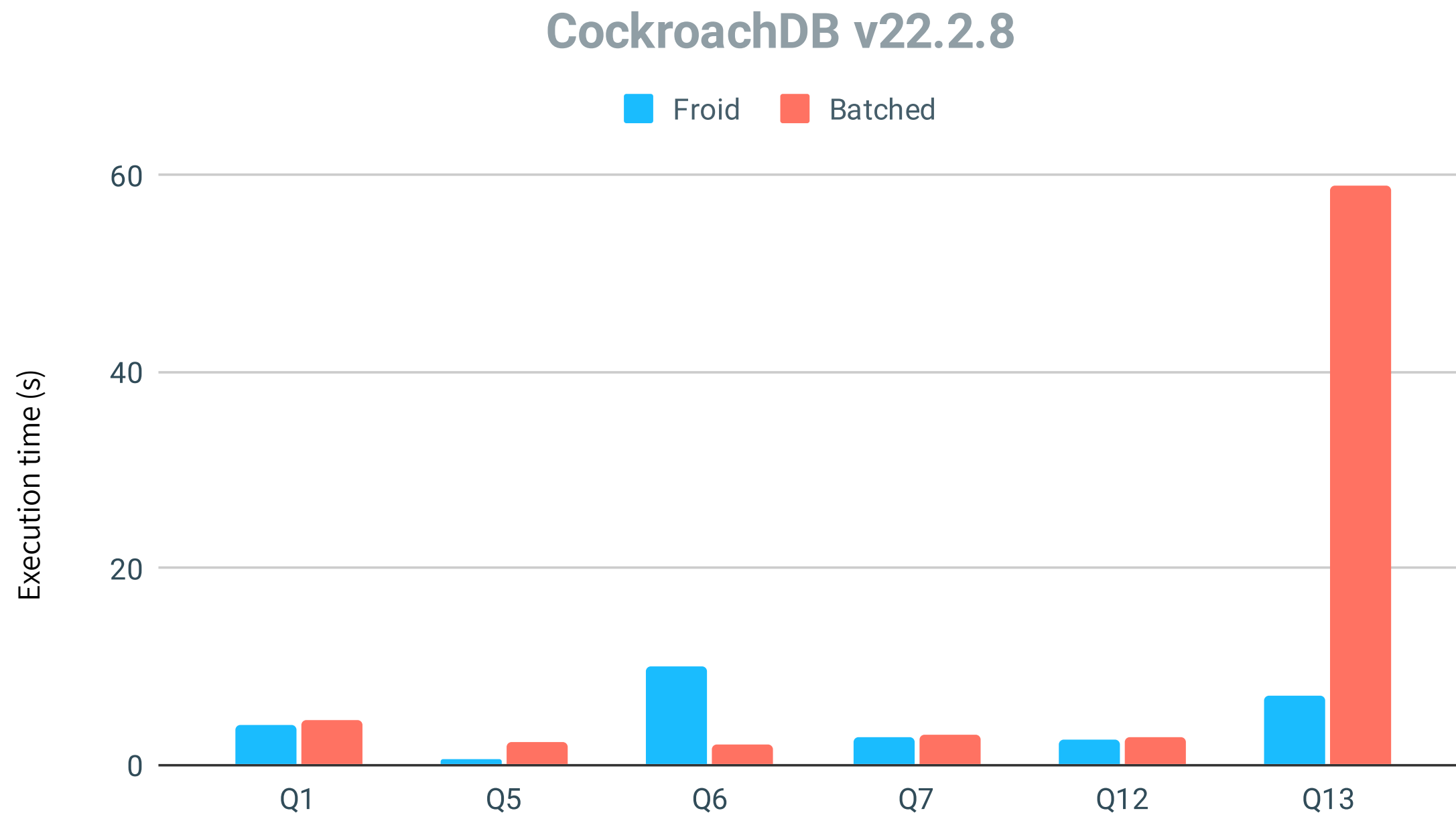- ProcBench benchmark, scale factor 1
- Compared against Froid



SQL Server 2019

# Experiments

- Postgres not good at decorrelation
- Batching works well with SQL->SQL decorrelation but unsure if this is a fair comparison



PostgreSQL 15

# Experiments

- CockroachDB is pretty good at handling Froid queries
- Temp tables still an experimental feature

**CockroachDB v22.2.8**

# Oracle Database 21c Enterprise



X = timeout (>30 min)

- Very inconsistent with Froid queries
- Batched queries have less variance

# Future Work

- Test on MySQL, DuckDB
- Optimizing procedural code
  using large language models

Thank you!