# Apollo:
## Learning Database Query Correlations for Predictive Caching

Brad Glasbergen, Khuzaima Daudjee, Michael Abebe, Scott Foggo, Anil Pacaci
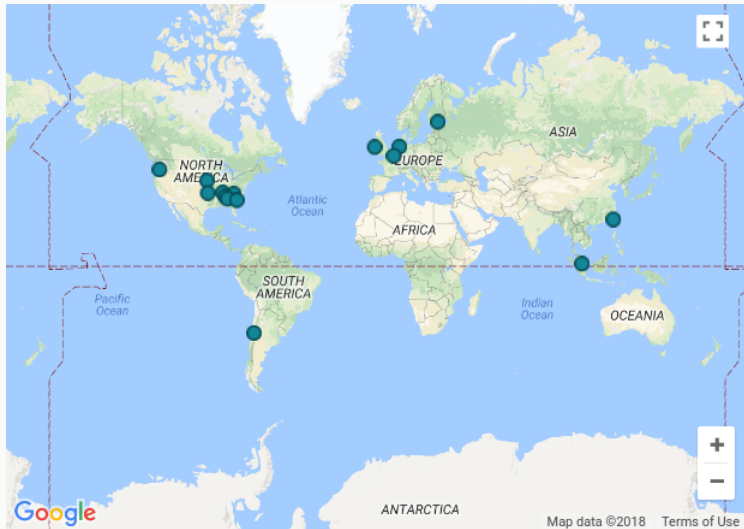
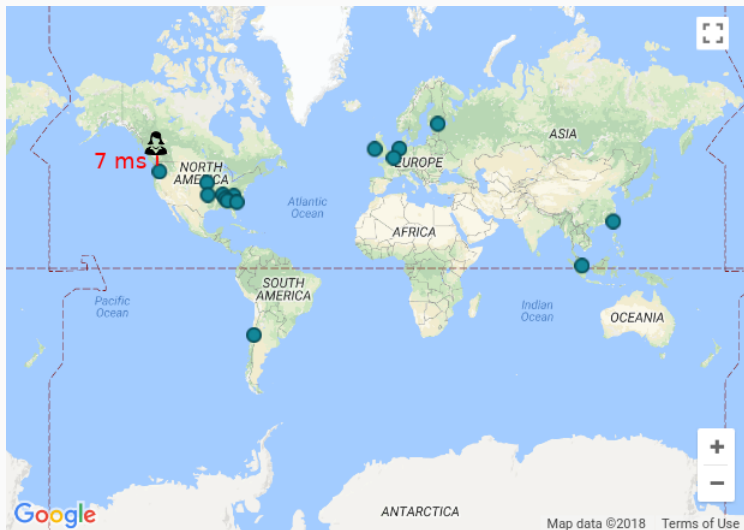UNIVERSITY OF WATERLOO

@bglasber

1

# Simple Web Application Architecture



Client

Database

Data Center

# Worldwide Client/Data Center Distribution

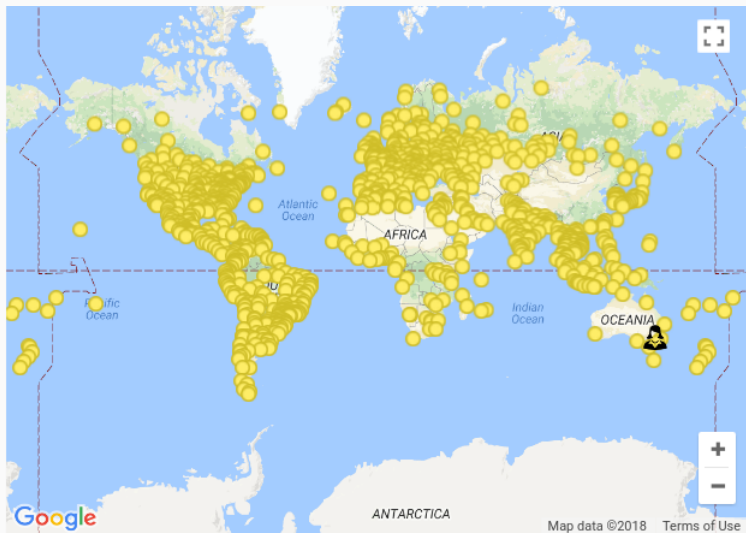# Worldwide Client/Data Center Distribution

# Latency Effects on Clients

Increased latency reduces user engagement, and consequently revenue!

Schurman et al., "Performance Related Changes and Their User Impact". *Velocity*, 2009.

# Edge Caching (Content Delivery Networks)



Client → Edge Node (Data Cache) → Data Center (Database)
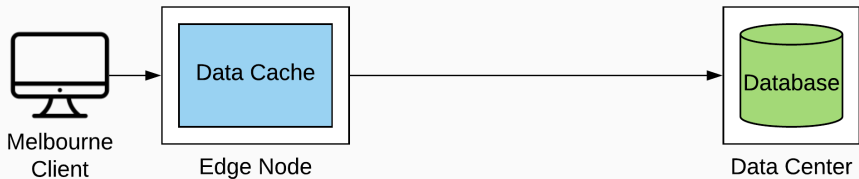
# Worldwide Client/Edge Node Distribution



Map of metros where at least one Edge node (GGC) is present.

- Limited support for non-static data!
- Can we extend support for dynamic data?

Melbourne
Client

Data Cache

Edge Node

Database

Data Center

# Extending Edge Cache Support



Melbourne Client → Edge Node (Cache Manager / Data Cache) → Data Center (Database)

# Extending Edge Cache Support



Melbourne Client → Edge Node [Cache Manager / Data Cache] → Data Center [Database]

Prefetch relevant data from DB to cache

```
1. SELECT  C_ID   FROM CUSTOMER WHERE
C_UNAME = @C_UN and C_PASSWD = @C_PAS

2. SELECT MAX(O_ID) FROM ORDERS WHERE
O_C_ID =  @C_ID
```
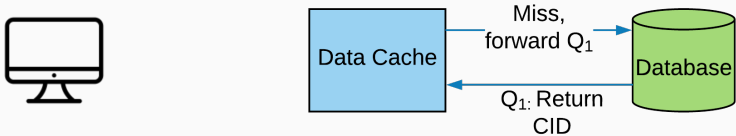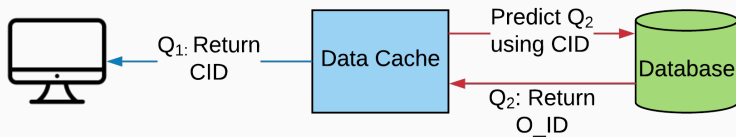
# Predictive Caching



**Q$_1$**: Look up customer ID
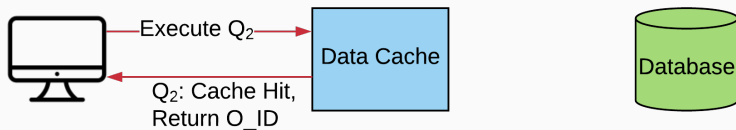**Q$_2$**: Look up last order for customer ID

Miss, forward $Q_1$

Data Cache

Database

$Q_1$: Return CID

**$Q_1$**: Look up customer ID
**$Q_2$**: Look up last order for customer ID

# Predictive Caching



$Q_1$: Return CID

Data Cache

Predict $Q_2$ using CID

Database

$Q_2$: Return O_ID

**$Q_1$**: Look up customer ID
**$Q_2$**: Look up last order for customer ID

# Predictive Caching



Execute $Q_2$

Data Cache

$Q_2$: Cache Hit, Return O_ID

Database

**$Q_1$**: Look up customer ID
**$Q_2$**: Look up last order for customer ID

## Apollo: Caching Dynamic Data

We developed Apollo, a middleware system that:

We developed Apollo, a middleware system that:

- Uses online learning to discover client query patterns.

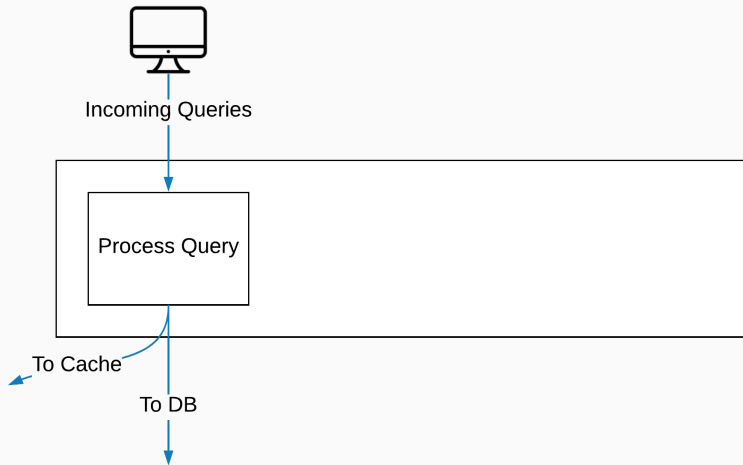# Apollo: Caching Dynamic Data

We developed Apollo, a middleware system that:

- Uses online learning to discover client query patterns.
- Predictively executes and caches query results using these patterns to reduce client response time.

# Apollo: Caching Dynamic Data

We developed Apollo, a middleware system that:

- Uses online learning to discover client query patterns.
- Predictively executes and caches query results using these patterns to reduce client response time.
- Employs a computationally efficient means of managing updates to cached data.

# Table of contents

# Predictive Query Model

# Apollo Overview

Incoming Queries

Process Query

Find Query Correlations

Find Parameter Mappings

To Cache

To DB

Predictively Execute and Cache

SELECT C_ID FROM CUSTOMER WHERE C_UNAME =
'Alice' and C_PASSWD = 'pass'

SELECT MAX(O_ID) FROM ORDERS WHERE O_C_ID = 3

# Query Templates

Two query instances, $Q_1$, $Q_2$ share the same query template if they have the same query text modulo *parameterizable constants*.

SELECT $\boxed{\text{C\_ID}}$ FROM CUSTOMER WHERE C_UNAME = **?**
and C_PASSWD = **?**

SELECT MAX(O_ID) FROM ORDERS WHERE O_C_ID = $\boxed{\textbf{?}}$

**Client Submits:** $Q_1$

# Client Query Streams

**Client Submits:** $Q_1$ →(Abstract)→ $Qt_1$

**Client Submits:**

$Q_2$  →  Abstract  →  $Qt_2$

Insert into Query Stream

| $Qt_1$ | $Qt_2$ |
|--------|--------|

# Client Query Streams

| Qt$_1$ **1** | Qt$_2$ **2** | Qt$_3$ **4** | Qt$_1$ **6** | ... |
|---|---|---|---|---|

| Qt$_1$ | Qt$_2$ | Qt$_3$ | Qt$_1$ |
|:---:|:---:|:---:|:---:|
| **1** | **2** | **4** | **6** |

...

# Client Query Streams

| $Qt_1$ | $Qt_2$ | $Qt_3$ | $Qt_1$ |
|---|---|---|---|
| **1** | **2** | **4** | **6** |

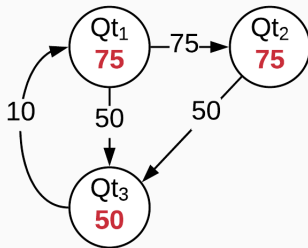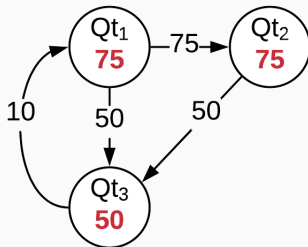**And so on...**

# Query Transition Graph

Probability of seeing $Qt_2$ within sliding window after we've seen $Qt_1$:
$P(Qt_2|Qt_1; T \leq \Delta t) = \frac{75}{75} = 1$

## Query Transition Graph



Probability of seeing $Qt_2$ within sliding window after we've seen $Qt_1$:

$P(Qt_2|Qt_1; T \leq \Delta t) = \frac{75}{75} = 1$

$P(Qt_1|Qt_3; T \leq \Delta t) = \frac{10}{50} = \frac{1}{5}$

SELECT C_ID FROM CUSTOMER WHERE C_UNAME = **?** and
C_PASSWD = **?**
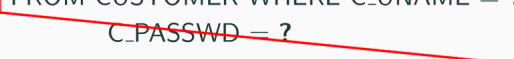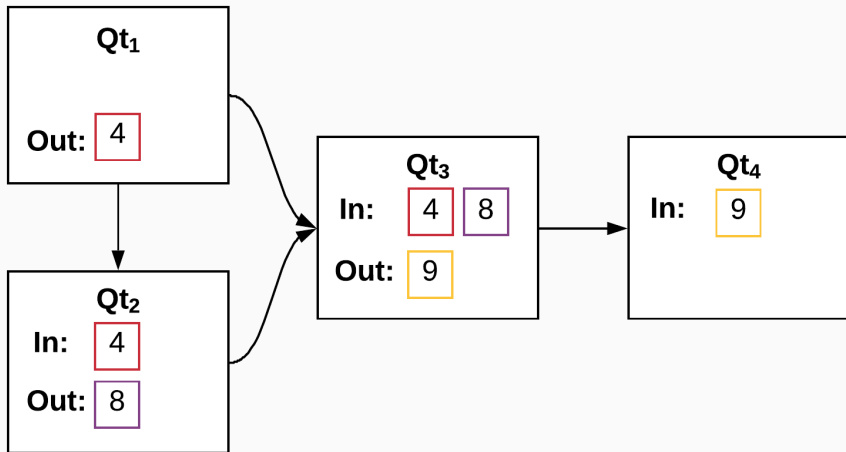
SELECT MAX(O_ID) FROM ORDERS WHERE O_C_ID = **?**

3

SELECT C_ID FROM CUSTOMER WHERE C_UNAME = **?** and
C_PASSWD = **?**

SELECT MAX(O_ID) FROM ORDERS WHERE O_C_ID = **?**

3

SELECT $\boxed{\text{C\_ID}}$ FROM CUSTOMER WHERE C\_UNAME = **?** and
C\_PASSWD = **?**

SELECT MAX(O\_ID) FROM ORDERS WHERE O\_C\_ID = $\boxed{\textbf{?}}$
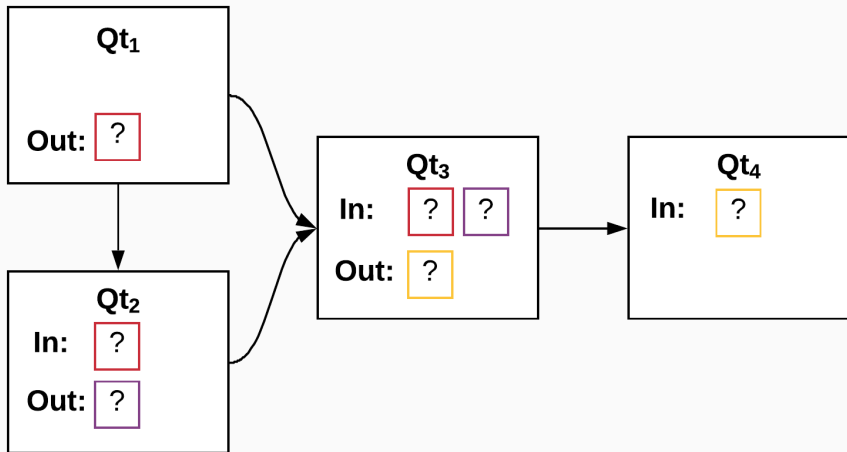
3

# Client Query Streams

SELECT C_ID FROM CUSTOMER WHERE C_UNAME = **?** and
C_PASSWD = **?**

SELECT MAX(O_ID) FROM ORDERS WHERE O_C_ID = **?**

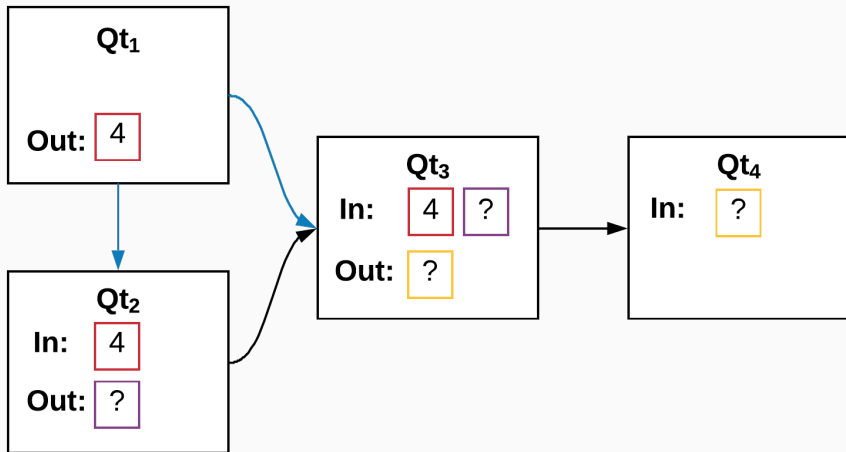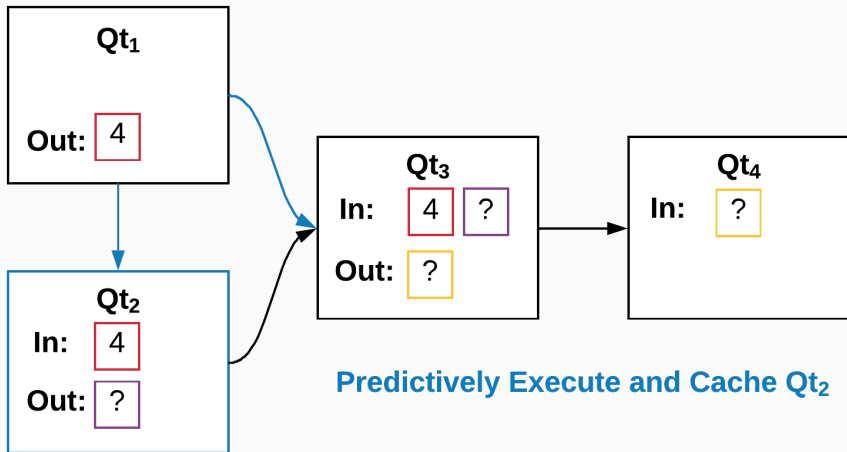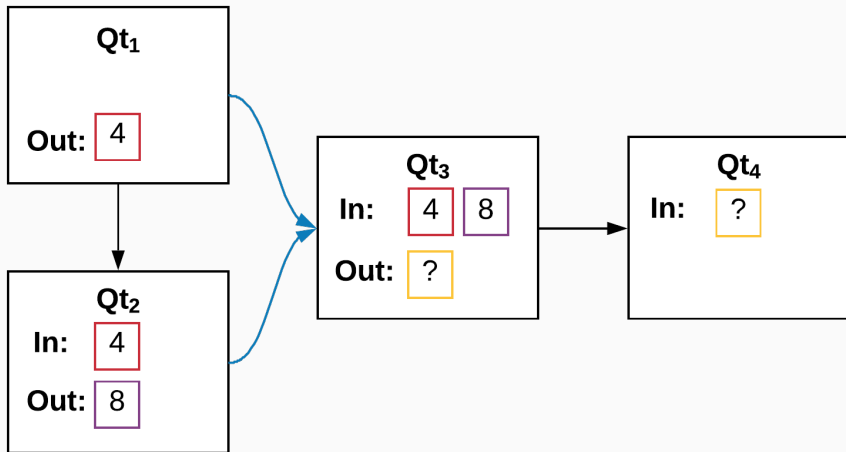# Dependency Graph

**Predictively Execute and Cache $Qt_2$**

# Prediction Routine

**Predictively Execute and Cache $Qt_3$**

**Predictively Execute and Cache Qt$_4$**

# Cache Management

Alice

Data
Cache

Post Product Review

Database

Bob

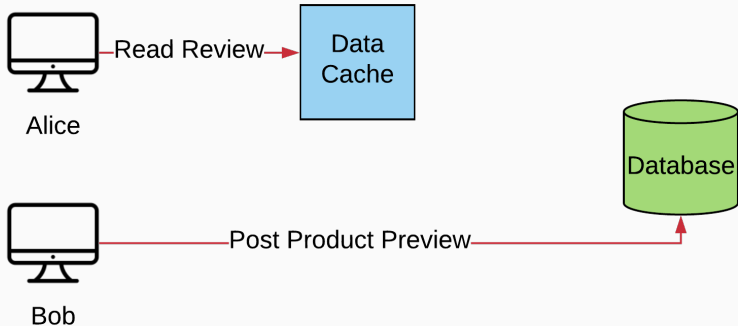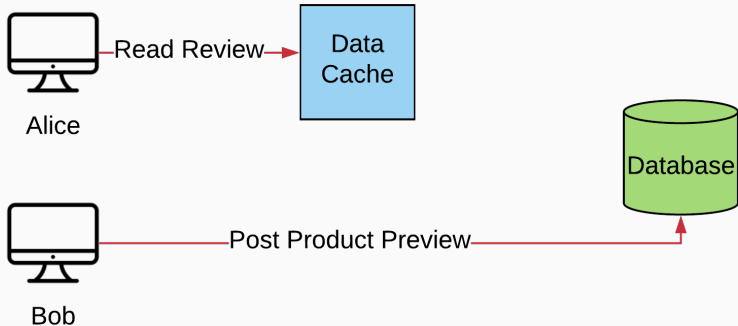# Client-Centric Caching

# Client-Centric Caching



Alice should see her own order, but does not care about Bob's!

# Client-Centric Caching
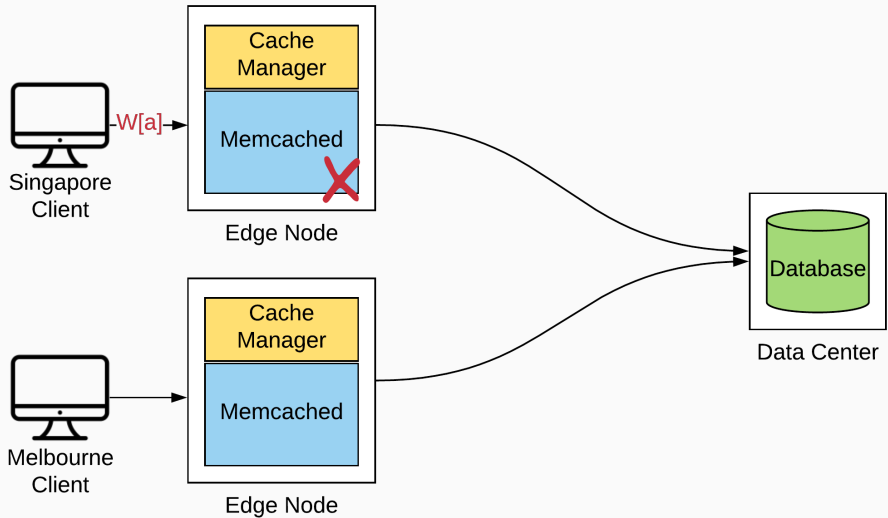


Alice should see her own order, but does not care about Bob's!

- Improves cache performance, client-centric model
- Support for reading latest data if needed

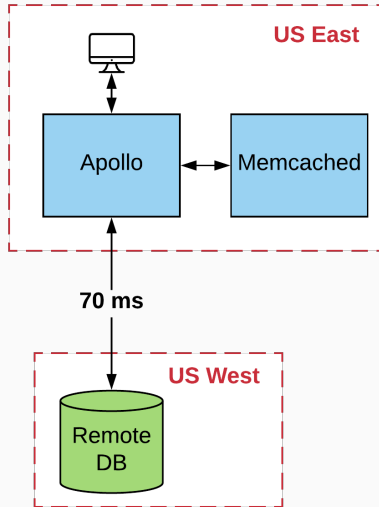# Benefits of Client-Centric Caching

## Benefits of Client-Centric Caching

- Can predict whether a prefetched query result will be used before invalidation!
- Reloading queries upon invalidation

See paper for details!

# Results

# Experiment Configuration

# Experiment Configuration

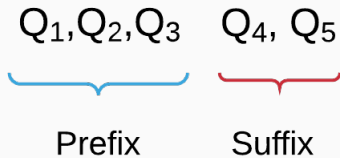Three configurations:

- Apollo configuration: as described in prior sections.
- Memcached configuration: LRU cache — Apollo with predictive features turned off
- Fido configuration: Use Fido predictive engine instead of Apollo's predictive features!

# Fido

$$Q_1, Q_2, Q_3$$

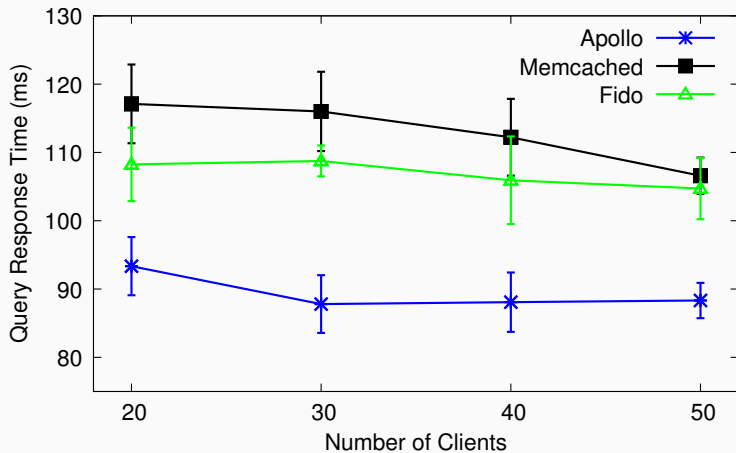Prefix

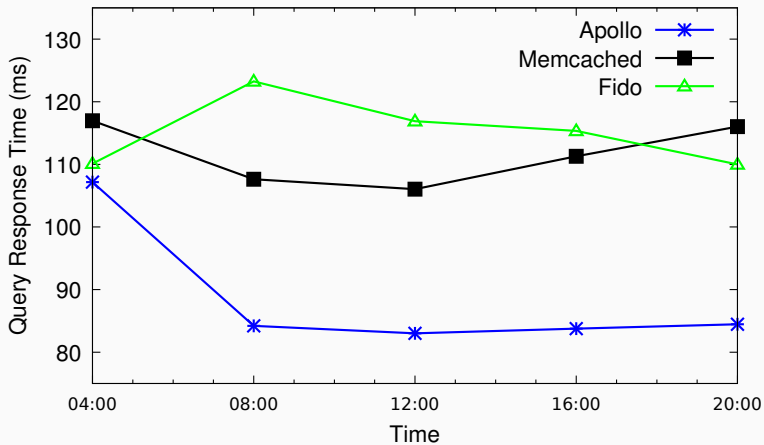$$Q_1, Q_2, Q_3 \quad Q_4, Q_5$$

Prefix    Suffix

# Fido

- Query instance based predictions, instead of query templates.
- Prefix length: 3, Suffix Length: 2
- Requires offline training (Supplied 40 minutes of data).
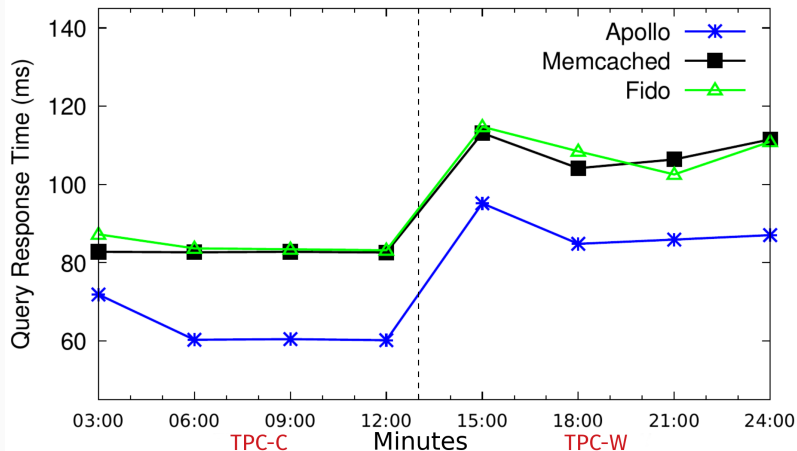
# TPC-W Results

# Learning Over Time

- Online learning enables workload pattern discovery and adapts to client behaviour changes
- Predictive caching is an effective tool to reduce application latency

# Multiple Apollo Instances

## Parameter Settings (TPC-W)

| Parameter | Setting |
| --- | --- |
| Window Width | 15s |
| Correlation Threshold | 0.99 |
| Reload Threshold | 0 |
| Cache Size | 5% of DB |

## Related Work (Systems)

Similarities and Differences Among Related Work

| System | Similarities | Differences |
|--------|-------------|-------------|
| **Scalpel** | • Prefetching via templates | • Offline training<br>• Write Handling<br>• Client-side<br>• Query Rewriting |
| **Fido** | • Prefetching<br>• Server-side/middleware | • Offline training<br>• Query Instances<br>• Write Handling |

📄 K. Amiri, S. Park, R. Tewari, and S. Padmanabhan.
**Dbproxy: a dynamic data cache for web applications.**
In *Proceedings 19th International Conference on Data Engineering (Cat. No.03CH37405)*, pages 821–831, March 2003.

📄 E. Schurman and J. Brutlag.
**Performance related changes and their user impact.**
*Velocity*, 2009.