# How In-Memory solutions assist with SaaS integrations

In-Memory Computing | SUMMIT

Craig Gresbrink – Solutions Architect – 24 Hour Fitness

24 FITNESS
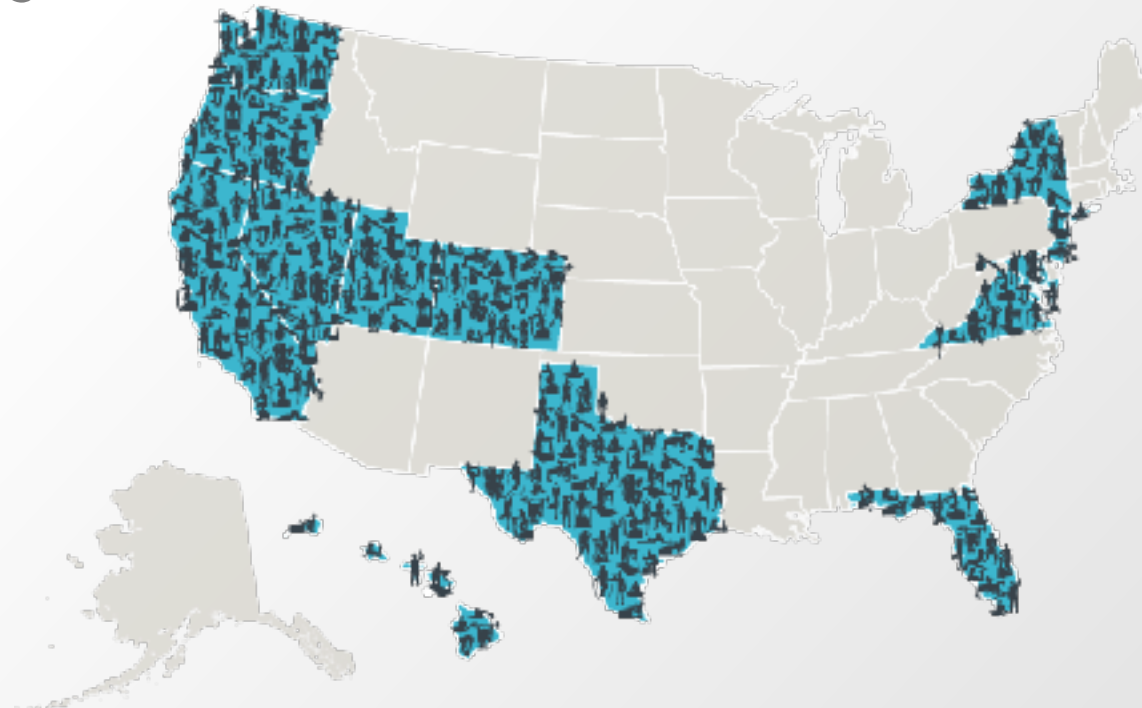
**Craig Gresbrink**
**cgresbrink@24hourfit.com**
**Solutions Architect**

## 24 Hour Fitness – 400+ clubs in 13 States

 - We are a leading fitness industry pioneer with nearly four million members in more than 400 clubs across the U.S. For more than 30 years, we've held fast to our mission of helping people improve their lives through fitness.
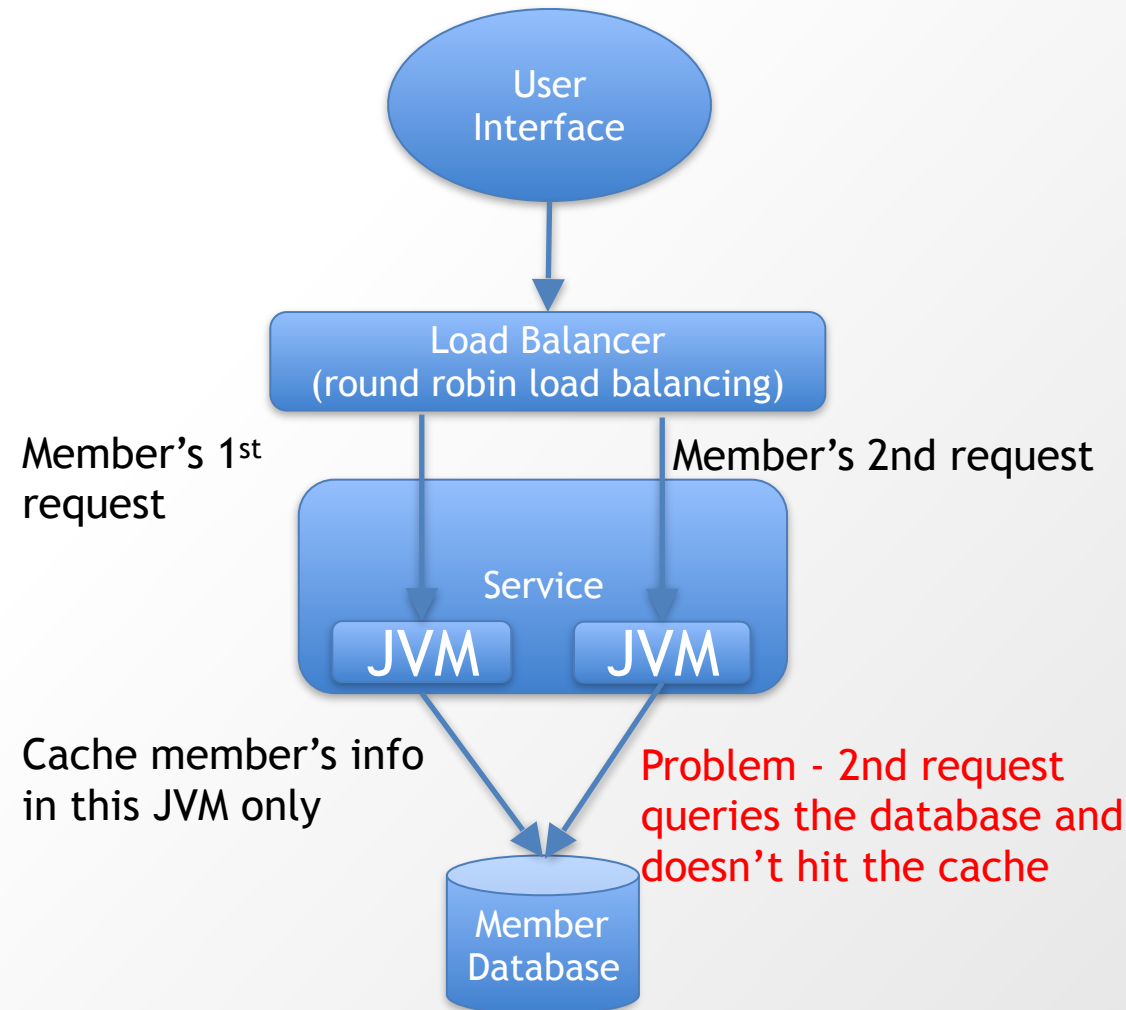
# *What am I going to show you?*

**How In-Memory solutions assist with SaaS integrations**

**Agenda**

- 24 Hour Fitness' historical architecture and some limitations

- How in-memory solutions solved use cases at 24 Hour Fitness

- Issues we ran into

- How in-memory solutions might be leveraged in the future at 24 Hour Fitness

- Q &A

- Reception!

User
Interface

Load Balancer
(round robin load balancing)

Member's 1st
request

Member's 2nd request

Service

JVM          JVM

Cache member's info
in this JVM only

Problem - 2nd request
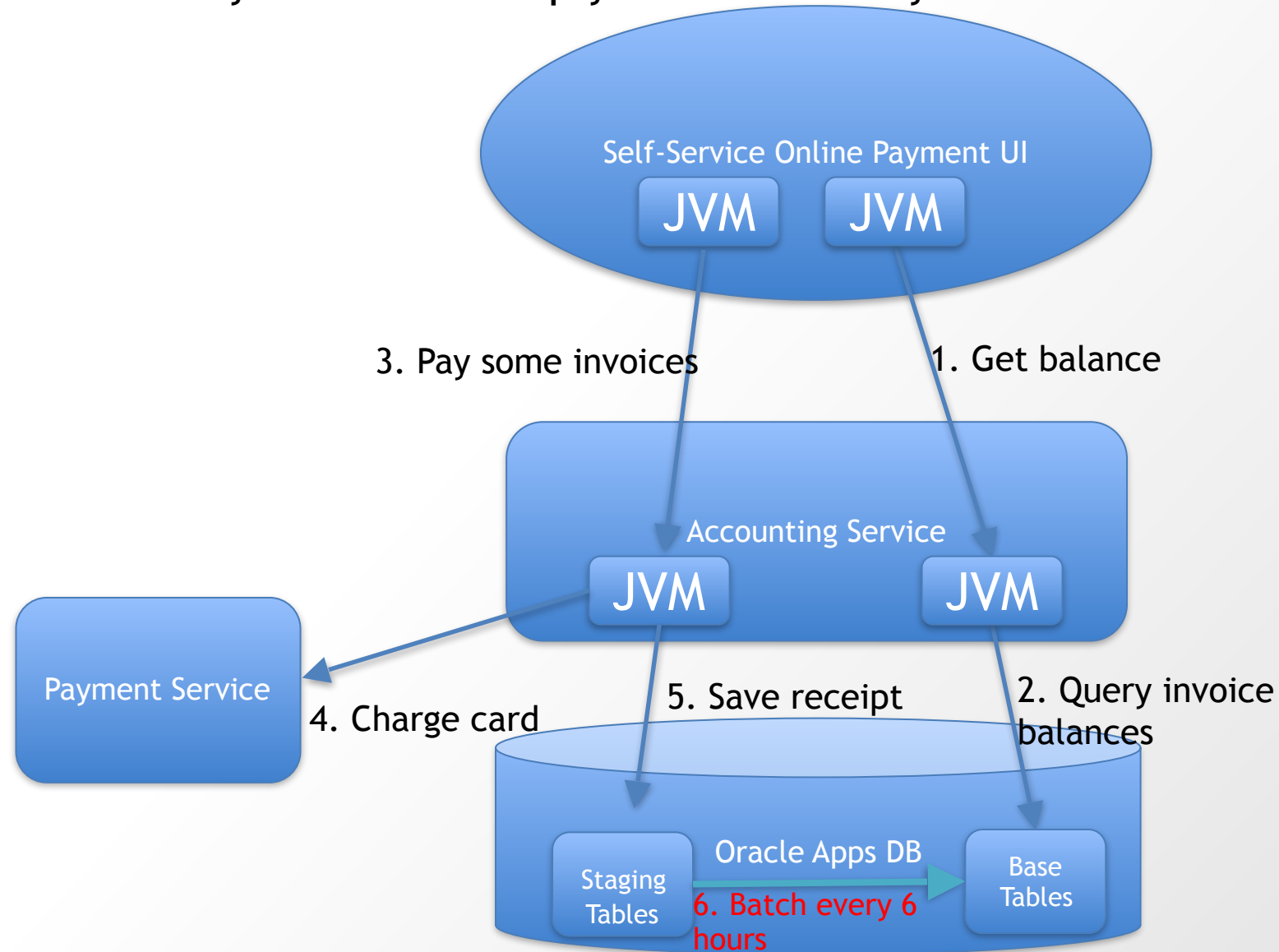queries the database and
doesn't hit the cache

Member
Database

We cache in each JVM only. In our setup there is no guarantee which JVM a user will hit due to round robin load balancing (no sticky session).

More cache misses as we scale (JVMs) horizontally.

# *Historical architecture and caching*

**Use Case 1 - Why it won't work for payments to a batch system**



Self-Service Online Payment UI

JVM    JVM

3. Pay some invoices    1. Get balance

Accounting Service

JVM    JVM

Payment Service

4. Charge card    5. Save receipt    2. Query invoice balances

Oracle Apps DB

Staging Tables    6. Batch every 6 hours    Base Tables
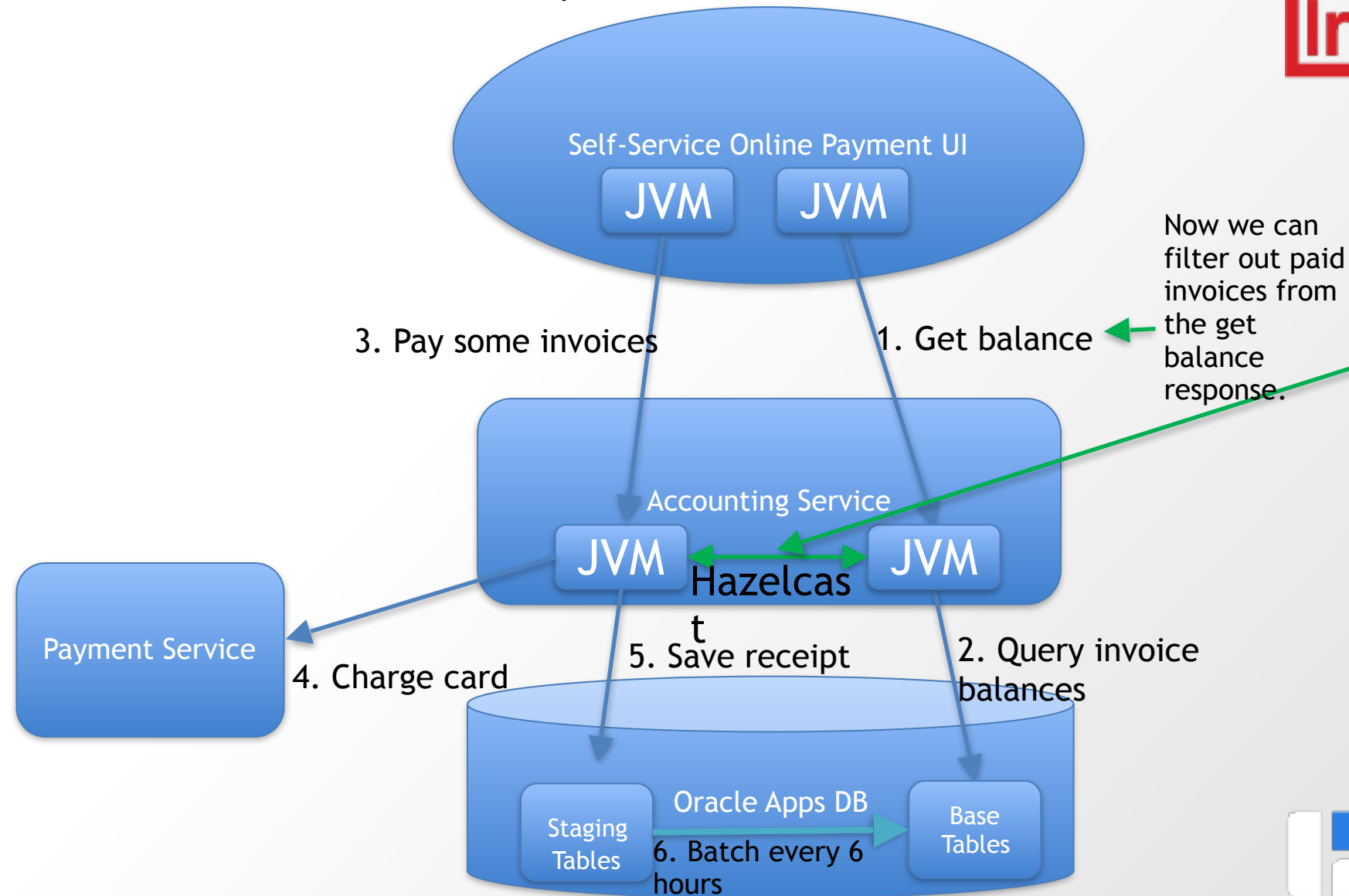
Members, can pay their bill(s) online via our self-service customer portal.

What happens when they come back 2 hours after making their payment?

We could've solved it by querying the staging tables as well as the base tables. It would be slow for all customers. 99% haven't made a previous payment.
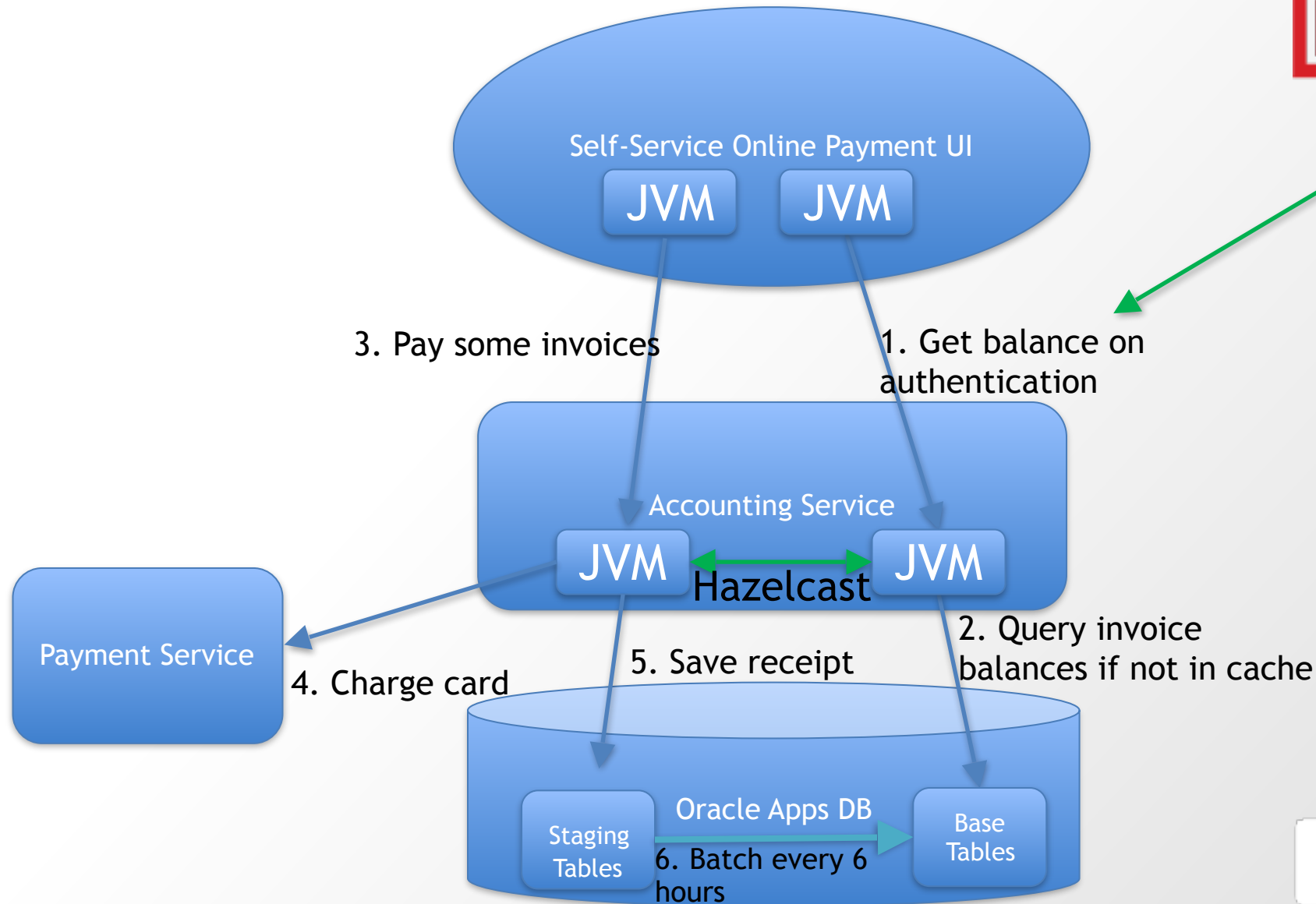
# A distributed cache solves our problem…

**Hazelcast to the rescue for use case 1.**



Self-Service Online Payment UI

JVM    JVM

3. Pay some invoices

1. Get balance

Now we can filter out paid invoices from the get balance response.

Accounting Service

JVM    JVM

Hazelcast

Payment Service

4. Charge card

5. Save receipt

2. Query invoice balances

Oracle Apps DB

Staging Tables

6. Batch every 6 hours

Base Tables

What if we implement a distributed cache such that cache consistency is retained across the JVMs so we know which invoices have already been paid?
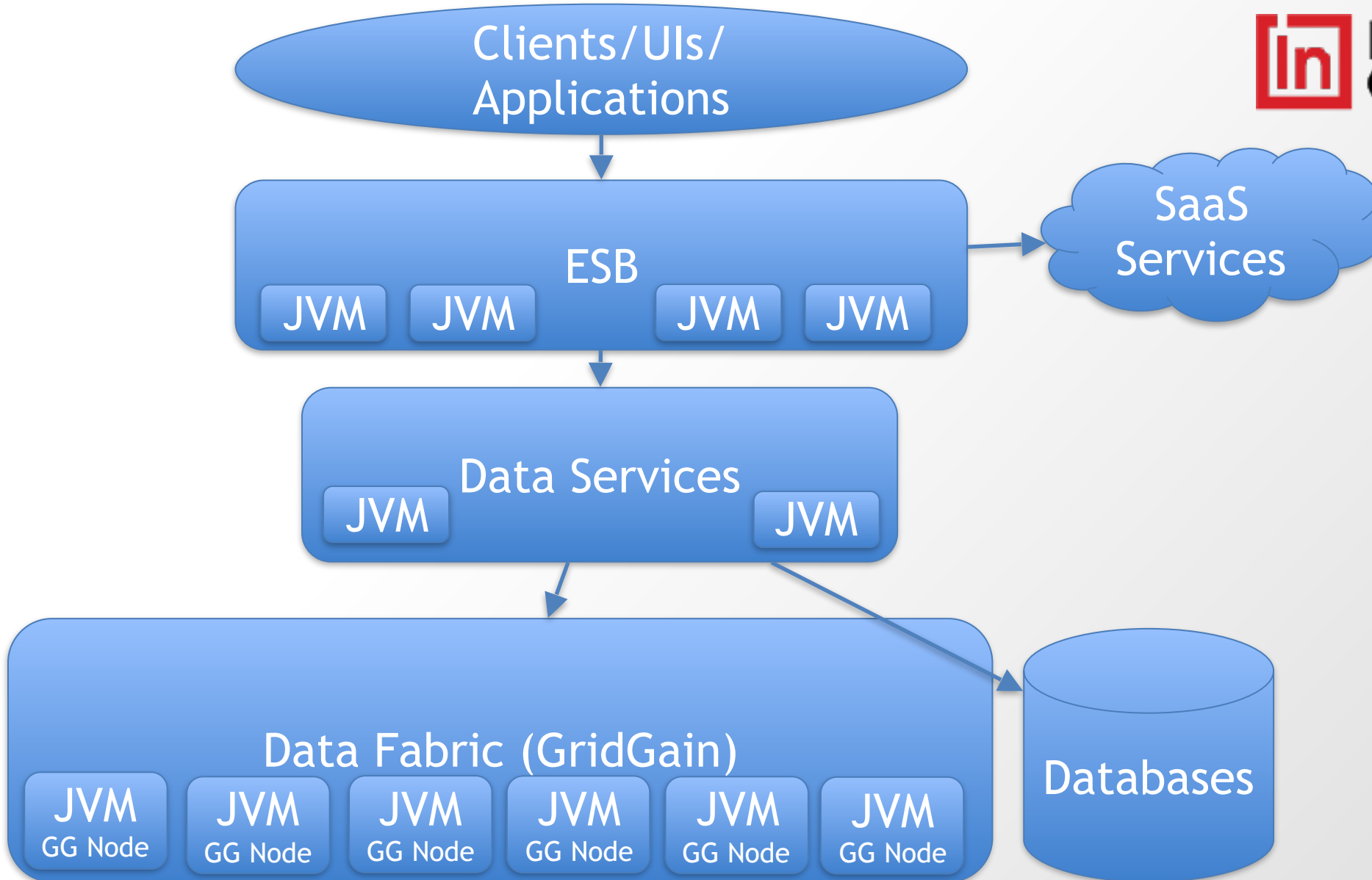
# A distributed cache solves another problem

**Self-Service Online Payment UI**

**JVM**　**JVM**

3. Pay some invoices

1. Get balance on authentication

**Accounting Service**

**JVM** ←→ **JVM**

**Hazelcast**

**Payment Service**

4. Charge card

5. Save receipt

2. Query invoice balances if not in cache

**Oracle Apps DB**

Staging Tables → Base Tables

6. Batch every 6 hours

Eagerly cache, or pre-cache, invoice balances.

Don't wait 'til they go to the make a payment screen. A better Customer Experience (CX).

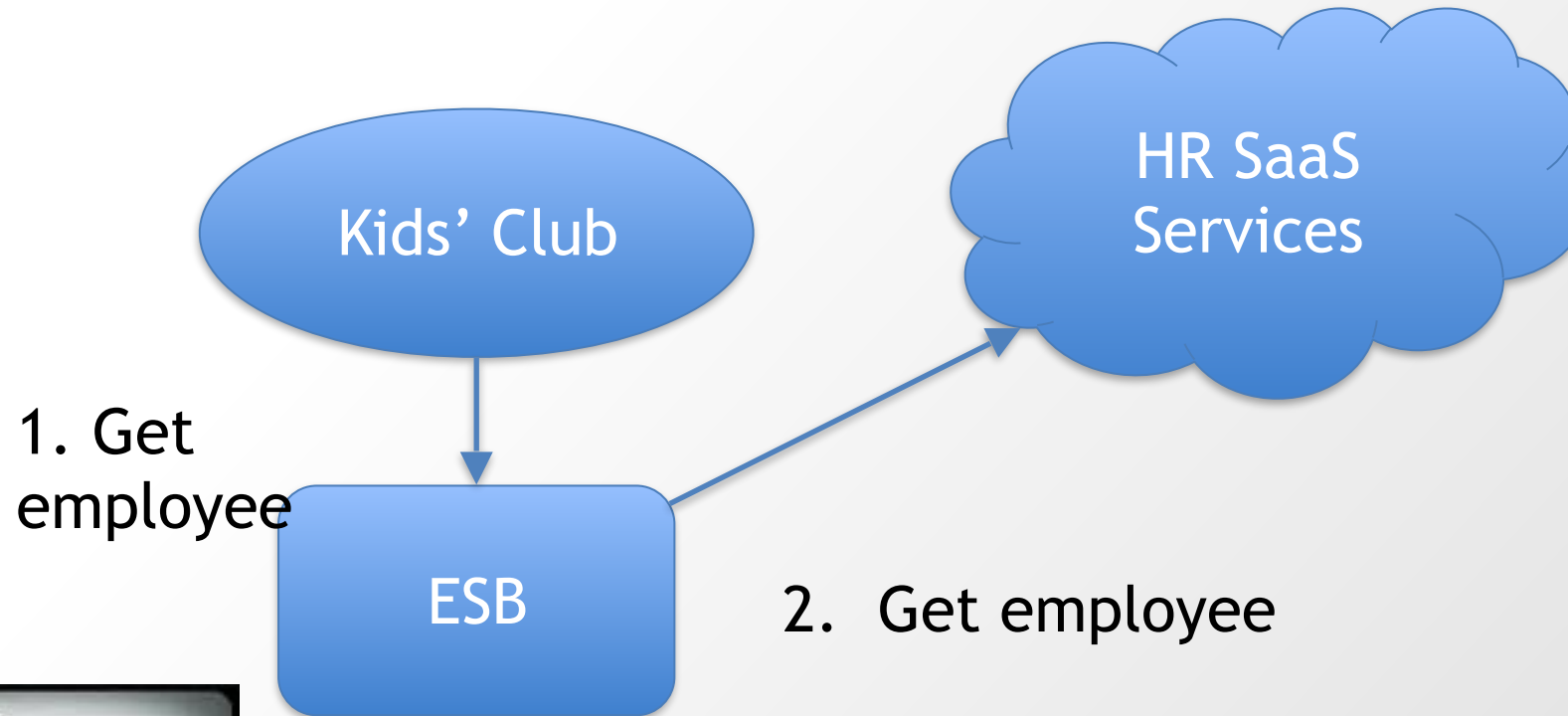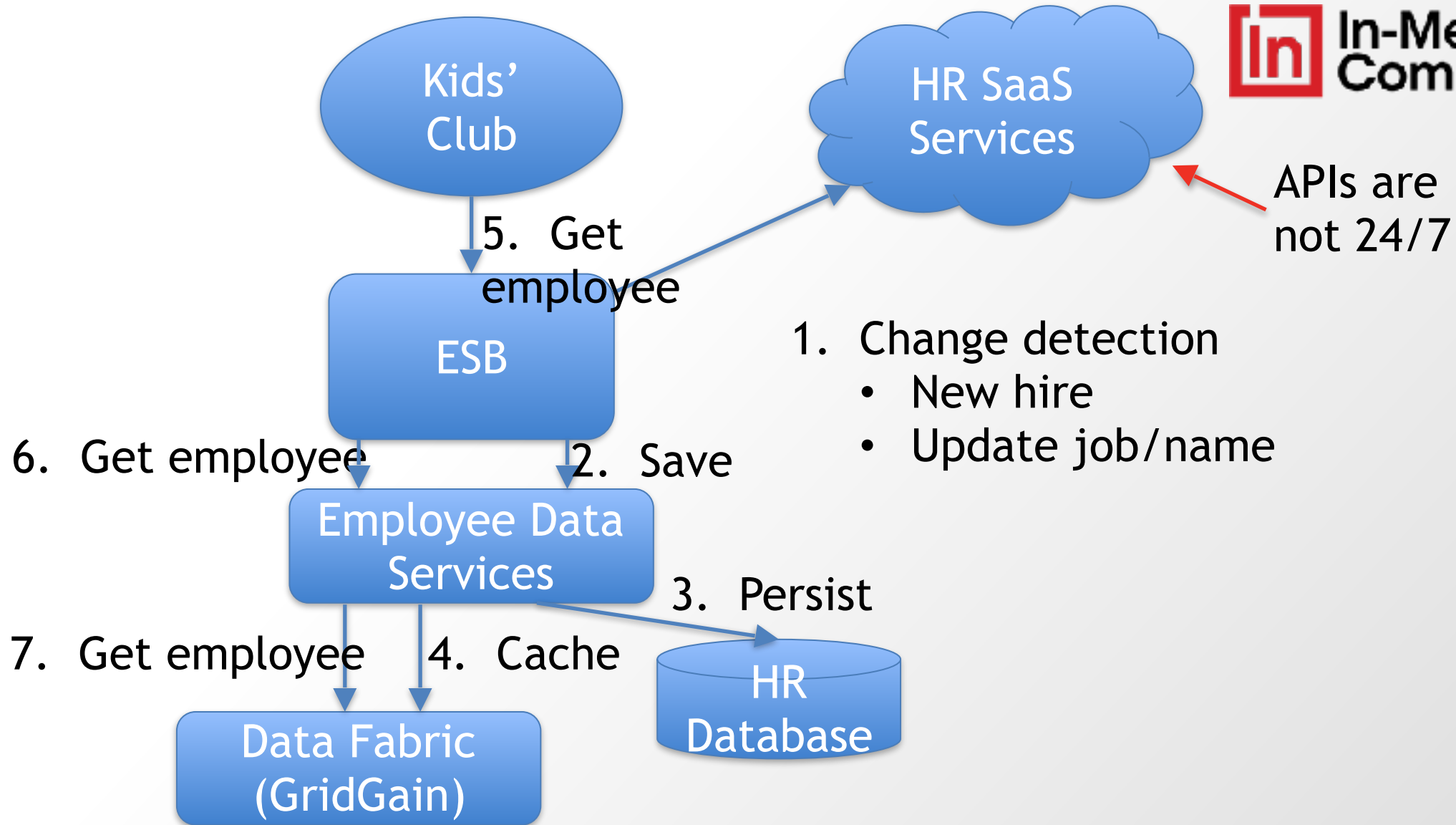But there is another problem, for some customers this query takes 40 seconds.

**Life is perfect!**

Kids' Club

HR SaaS Services

1. Get employee

ESB

2. Get employee

Kids' Club

HR SaaS Services

APIs are not 24/7

5. Get employee

ESB

1. Change detection
   - New hire
   - Update job/name

6. Get employee

2. Save

Employee Data Services

3. Persist

7. Get employee

4. Cache

HR Database

Data Fabric (GridGain)

An improved solution?

Kids' Club UI

ESB

HR SaaS Services

5. Get employee

6. Get employee

2. Save

1. Change detection
- New hire
- Update job/name

Employee Data Services

7. Get employee

3. Cache

Data Fabric (GridGain)

4. Persist – write through

HR Database

An improved, improved solution?

**Kids' Club UI**

**ESB**

**HR SaaS Services**

5. Get employee

6. Get employee

2. Save

**Employee Data Services**

7. Get employee

3. Cache

**Data Fabric (GridGain)**

4. Native persistence

1. Change detection
  - New hire
  - Update job/name

13

Nirvana, finally!

Membership Sales Online

Cloud Lead Management Services

1. Purchase

ESB

2. Create lead

**It's not that easy! We have to make multiple calls to LMS to create a lead...**

Membership Sales Online

Cloud Lead Management Services

1. Purchase

ESB

2. Get lead origin GUID
3. Get state GUID
4. Get club GUID
5. Create lead with GUIDs

Disclaimer:
Chatty APIs have been rectified in future versions of the product.

**Data flow with GridGain**



Membership Sales Online

1. Purchase

ESB

Cloud Lead Management Services

5. Create lead with GUIDs

2. Get lead origin GUID     3. CA GUID   4. Get club team GUID

CRM Entity Data Services

Get GUID

Data Fabric (GridGain)

Eagerly cache LMS GUIDs on startup

**Lead Management System and GUID translations in Camunda BPMN**
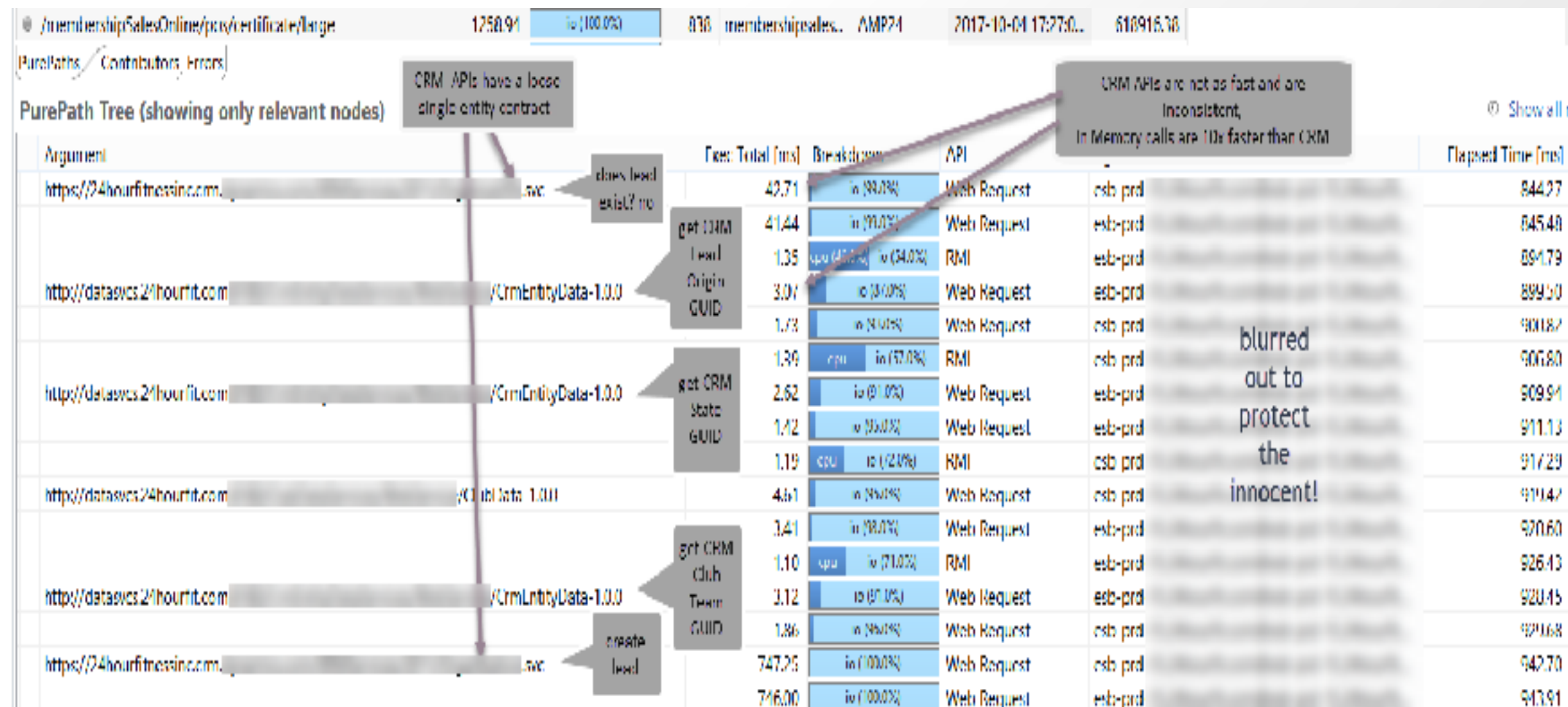
# Use case 3: A chatty real-time integration with CRUD-based APIs is too slow

**Lead Management System and GUIDs translations – How fast is the cache?**

*The road wasn't all smooth!*

Challenges...

1. **Marshalling and class loading**
- **What** – Historically,  the free version, Apache Ignite, didn't have the GridGain Binary Marshaller.
- **Issue** –  Without the GridGain Binary Marshaller, our domain model classes (member, agreement etc..) would need to be loaded into each node in the cluster, and this didn't work when adding a node the cluster due to an issue in the version we were using.
- **Solution**  - We created our own 24 Hour Custom Marshaller to serialize to standard java objects such that Ignite did not need our domain model classes.
- **Future** –  We want to move away from our 24 Hour Custom Marshaller and instead use the Binary Marshaller that is now part of Apache Ignite.

2. **Intermittent Node Connectivity**
- **What** – We were setup with multicast and eagerly loading our employee data.
- **Issue** – We experienced intermittent issues with nodes not returning employee data.
- **Why** - Due to our network issues, nodes could not contact each other and thus nodes redistributed data. There was some data loss depending which nodes disconnect from each other and which node a particular call went to.
- **Solution**  - We switched to TCP/IP where each server knows about all the IP addresses of other servers in the cluster.

3. **Involve other groups/teams even though it is a software solution**
- **Networking /Infrastructure** - Alluded to above, but additionally, we had timeouts that had the same affect of nodes exiting and rejoining the cluster.
- **Operations** – Instrumentation (for us Dynatrace),  cache validation(Visor UI or scripts and REST APIs), reloading the cache (scripts and REST APIs).
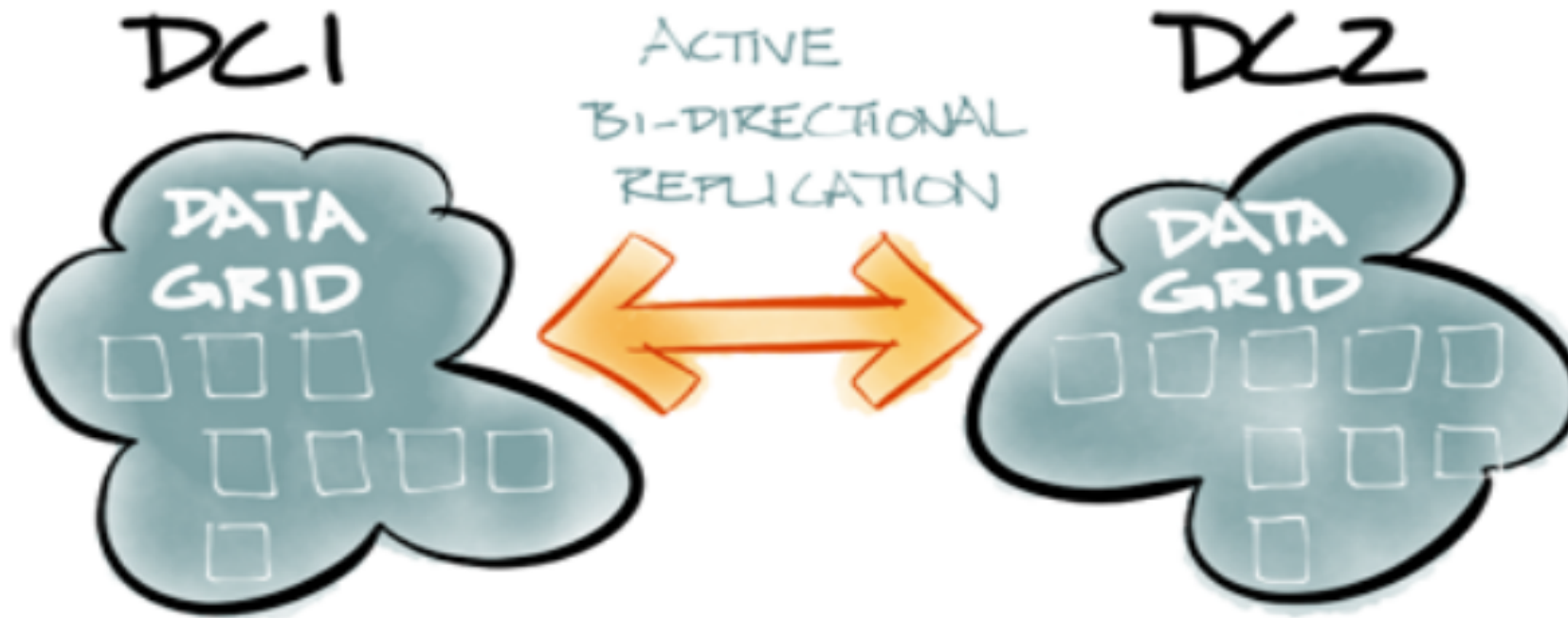
# *Future possibilities*

**Previously shown future possibilities**

– **Write Through**
   - HR Employee data with database to support legacy systems

– **Native Persistence**
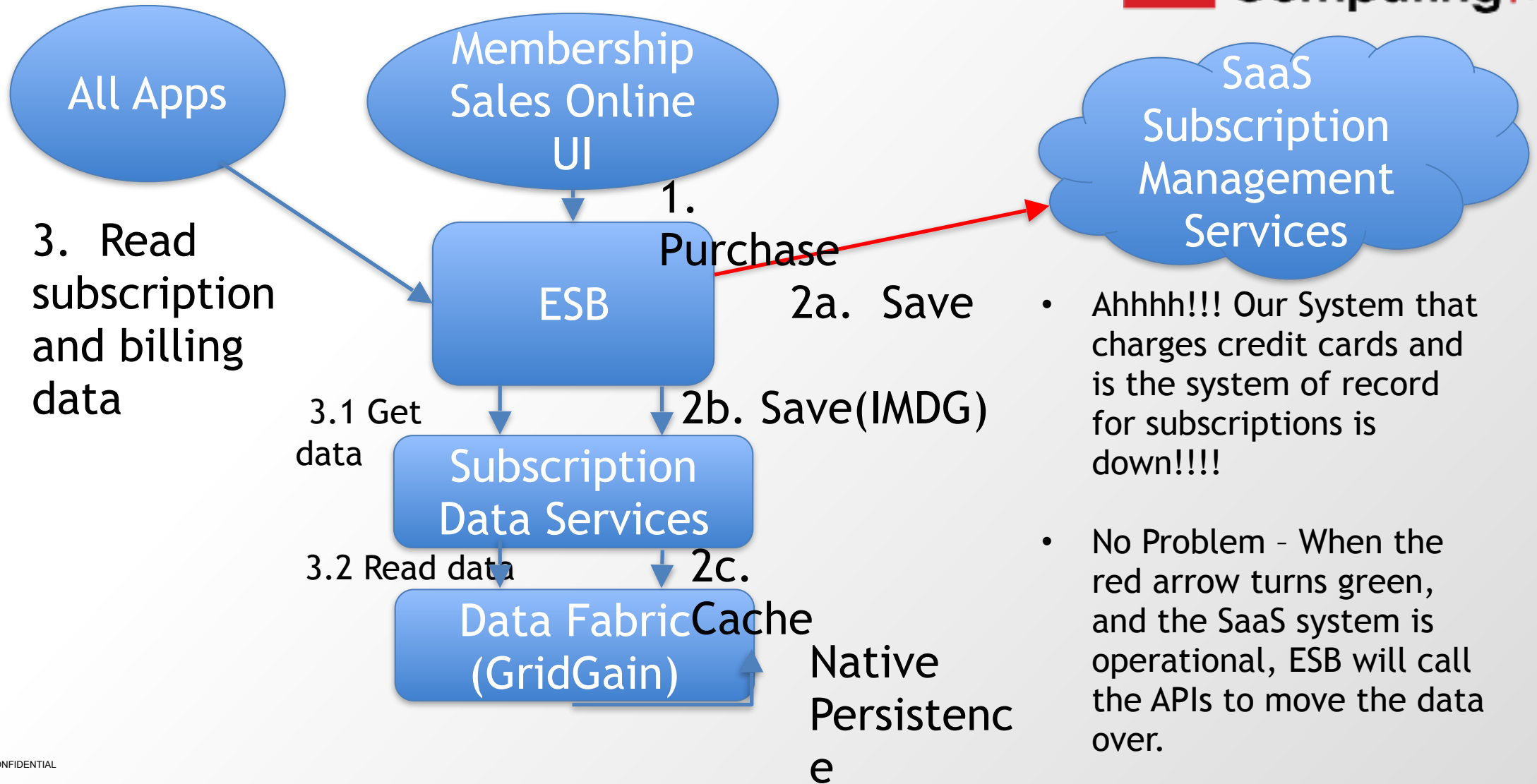   - HR Employee data w/o database

**Other ideas**

– **ID Generation vs. MySQL sequence**
   - We have a vendor that has APIs where we are expected to provide their primary key IDs
   - It is a nightly batch with many rows/records so it would be nice/faster to use Ignite's in-memory ID generation

– **Utilize the distributed nature of the Data Fabric to support CSA (Continuous Service Availability)**
   - We currently have two data centers, and we are implementing an Active/Active solution across both data centers
      – When one data center is down for maintenance we'll utilize the other data center

The picture below illustrates a simple view of Data Center Replication.

Vendor has concurrency API rate limits, so use GridGain to lessen the load

24 FITNESS

In-Memory Computing | SUMMIT

**All Apps**

**Membership Sales Online UI**

**SaaS Subscription Management Services**

3.  Read subscription and billing data

1. Purchase

**ESB**

2a.  Save

2b. Save(IMDG)

3.1 Get data

**Subscription Data Services**

3.2 Read data

2c. Cache

**Data Fabric (GridGain)**

Native Persistence

- Ahhhh!!! Our System that charges credit cards and is the system of record for subscriptions is down!!!!

- No Problem – When the red arrow turns green, and the SaaS system is operational, ESB will call the APIs to move the data over.

22

1. **Not 24/7 (possibly specified in the contract allowing them to be down x hours per day/week/month, not that they actually are)**
2. **Not Performant (calls too slow, or chattiness makes business transaction too slow)**
3. **API Rate Limits (calls per minute or concurrency – shown below)**

This policy prevents tenants from monopolizing SaaS system resources.

# Concurrent Request Limits

Each tenant has the following default concurrent request limits:

| Type | Description | Default Limit for Concurrent, Uncompleted Requests | Retry After |
|---|---|---|---|
| **Total Request** | Refers to UI, REST, and SOAP API requests processed under a tenant.<br>· API calls (REST and SOAP)<br>   ○ Create<br>   ○ Update<br>   ○ Get<br>   ○ Delete<br>   ○ Generate<br>   ○ Query<br>· UI requests (all)<br>Excluded from this policy<br>· Login API calls (REST and SOAP)<br>· UI login requests | 40 | 120 seconds |

## *Summary*

**What have we seen?**

- 24 Hour Fitness' historical architecture and some limitations

- How in-memory solutions solved use cases at 24 Hour Fitness

- Issues we ran into, marshalling and intermittent node connectivity

- How in-memory solutions might be leveraged in the future at 24 Hour Fitness

**Tidbits**

- A thank you to the Target IMC Summit 2016 presenters
  - Maybe you'll be up here next year, and I'll be watching you present!

- Careers at 24 Hour Fitness IT
  - North County San Diego (Carlsbad)
  - http://careers.24hourfitness.com

- A thank you to all my colleagues at 24 Hour Fitness that helped me with details in this presentation, attended my practice runs, and to those who implemented these in-memory solutions

# Questions…..

**Craig Gresbrink**

[cgresbrink@24hourfit.com](mailto:cgresbrink@24hourfit.com)
– Subject – IMC (start with IMC)
• I will do my best to respond within 2 or 3 days, for at least the next month.

# Thank You!

**Don't forget to check out the next slide: References and further reading…**

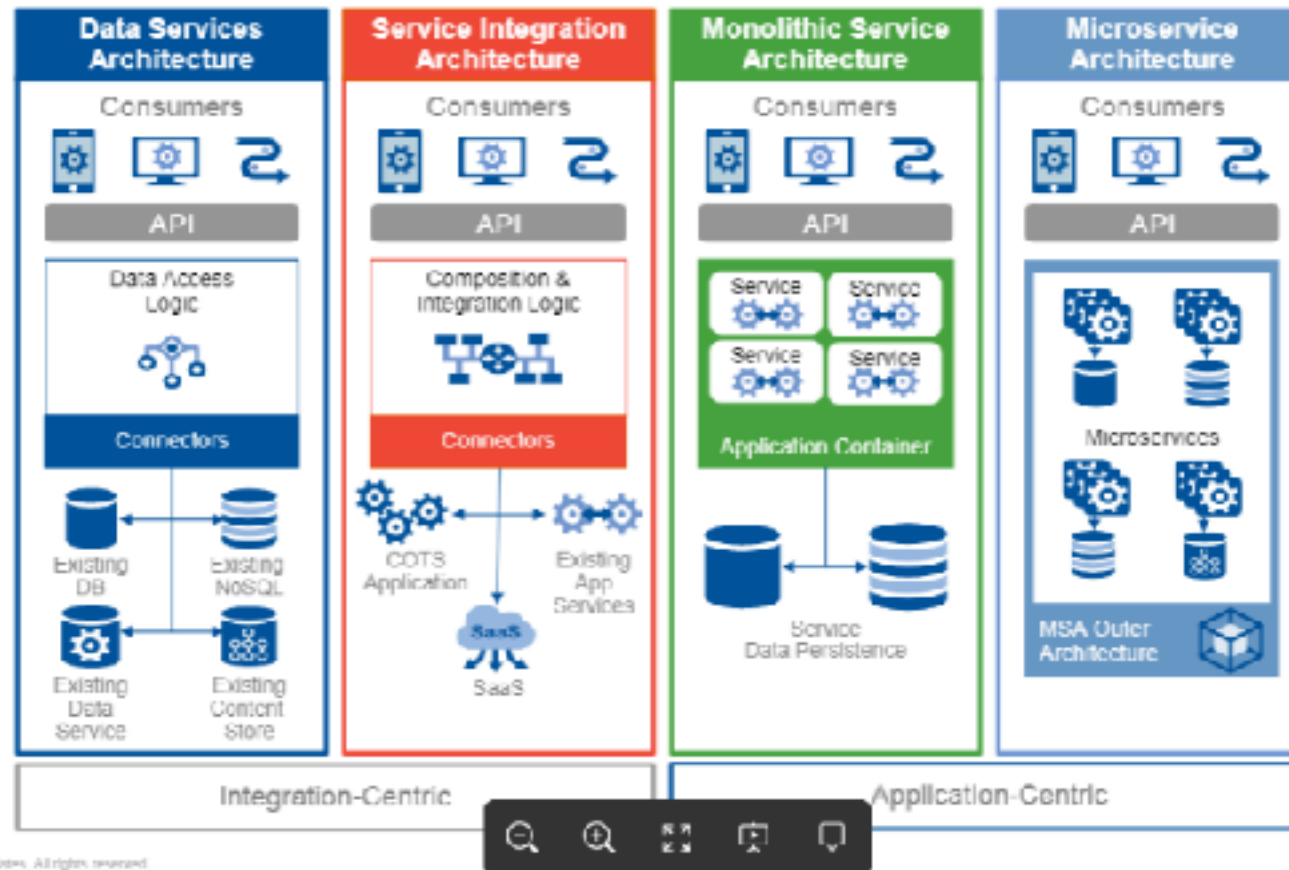# *References and further reading*

## Give credit, where credit is due...

- Hazelcast
  - https://hazelcast.com/

- Cache-Aside Pattern
  - https://blog.cdemi.io/design-patterns-cache-aside-pattern/

- You Are Not Netflix: How and When to Use Microservices in the Enterprise
  - https://www.gartner.com/webinar/3437517

- Microservices vs. Service Oriented Architecture
  - http://www.oreilly.com/programming/free/microservices-vs-service-oriented-architecture.csp

- Apache Ignite
  - https://ignite.apache.org/

- Apache Ignite's Read Through, Write Through, and Write Behind
  - https://apacheignite.readme.io/v1.1/docs/persistent-store

- Apache Ignite Native Persistence
  - https://apacheignite.readme.io/docs/distributed-persistent-store

- Camunda - an open source platform for workflow and business process management
  - https://camunda.org/

- Dynatrace - an Automated Performance Management product (APM)
  - https://www.dynatrace.com/

- Apache Ignite's Binary Marshaller
  - https://apacheignite.readme.io/docs/binary-marshaller

- Apache Ignite's ID Generator
  - https://apacheignite.readme.io/docs/id-generator

- GridGain's Data Center Replication
  - https://docs.gridgain.com/v8.1/docs/data-center-replication

- IMC Summit 2016 - Target's Presentation which was my inspiration (Thank You)
  - https://www.imcsummit.org/2016/videos-and-slides/targets-first-foray-into-an-in-memory-data-grid-and-the-trips-stumbles-and-falls-that-came-with/

- High performance in-memory computing with Apache Ignite (I haven't read this but it looked good)
  - https://www.amazon.com/Performance-memory-computing-Apache-Ignite/dp/1365732355/ref=sr_1_cc_1?s=aps&ie=UTF8&qid=1507758595&sr=1-1-catcorr

# »B-SIDES

# »OUTTAKES

The following slides were contemplated but didn't make the final cut....

*What is this, 2005, ESB? No Microservices?*

**What does lead creation look like in a visual Business Process Management System?**