# Apache Geode Summit



[https://www.youtube.com/playlist?list=PLgCUiLdOC2pi9SYptnVamh6xxoeRVZx8s](https://www.youtube.com/playlist?list=PLgCUiLdOC2pi9SYptnVamh6xxoeRVZx8s)

Pivotal

SpringOne Platform by Pivotal **re:Cap Seoul 2019**
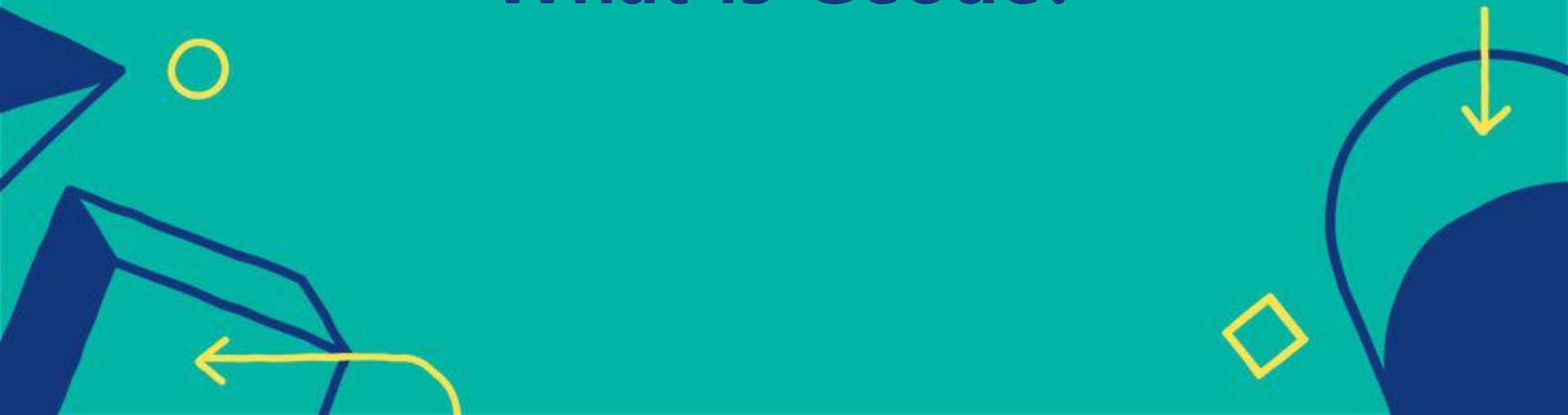
# Geode Sessions

- **Breaking Open Apache Geode: How It Works and Why**
- Introducing the Geode Native Client
- **Performance in Geode: How Fast Is It, How Is It Measured, and How Can It Be Improved?**
- Using Apache Geode: Lessons Learned at Southwest Airlines
- A Fireside Chat with Apache Geode Committers
- **Visualize Your Geode Metrics**
- **Reactive Event Processing with Apache Geode**
- Data Serialization and CI/CD Techniques for Apache Geode
- High-Performance Data Processing with Spring Cloud Data Flow and Geode
- Scaling Beyond a Billion Transactions Per Day with Sub-Second Responses
- Scalable, Cloud-Native Data Applications by Example
- Simple Data Movement Patterns: Legacy Application to Cloud-Native Environment and Apache Geode

SpringOne Platform by Pivotal **re:Cap Seoul 2019**

# Agenda

- What is Geode?

- How to Monitor: Visualize Geode Metrics

- Performance in Geode

- Reactive Programming with Geode

# 피보탈이 투자하고 있는 오픈소스 생태계

# Geode / GemFire / Pivotal Cloud Cache (PCC)

**Continuous investment in data R&D**

**What is good for one is good for the other**

Pivotal
**Cloud Cache**

On-platform & clouds

High-performance caching & data
acceleration for microservices

https://pivotal.io/pivotal-cloud-cache

Pivotal
**GemFire**®

Off-platform & native clients

In-memory data grid enabling
event-driven architecture and
fast-data access patterns

https://pivotal.io/pivotal-gemfire

APACHE
**GEODE**

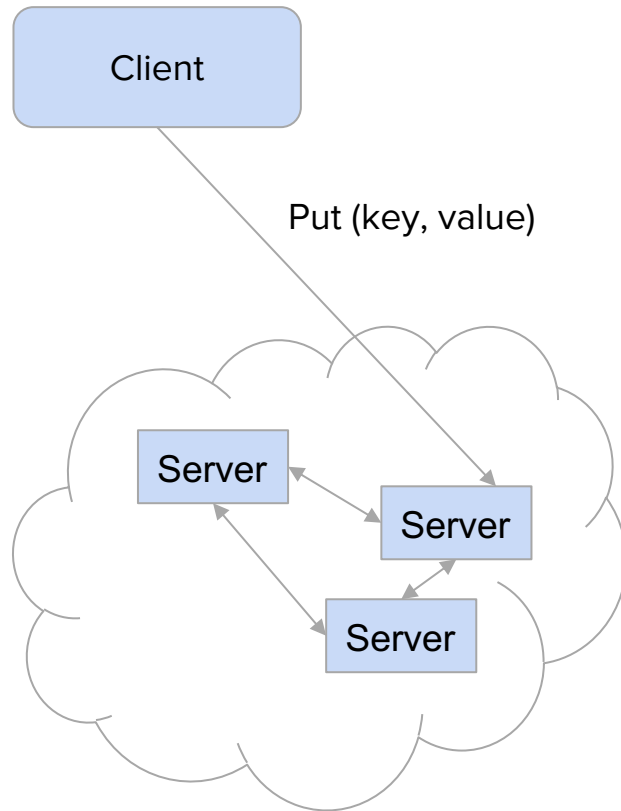Open source - Performance is key. Consistency is a must.

Reliable transaction processing and shared-nothing architecture for very low

latency performance with high concurrency processing

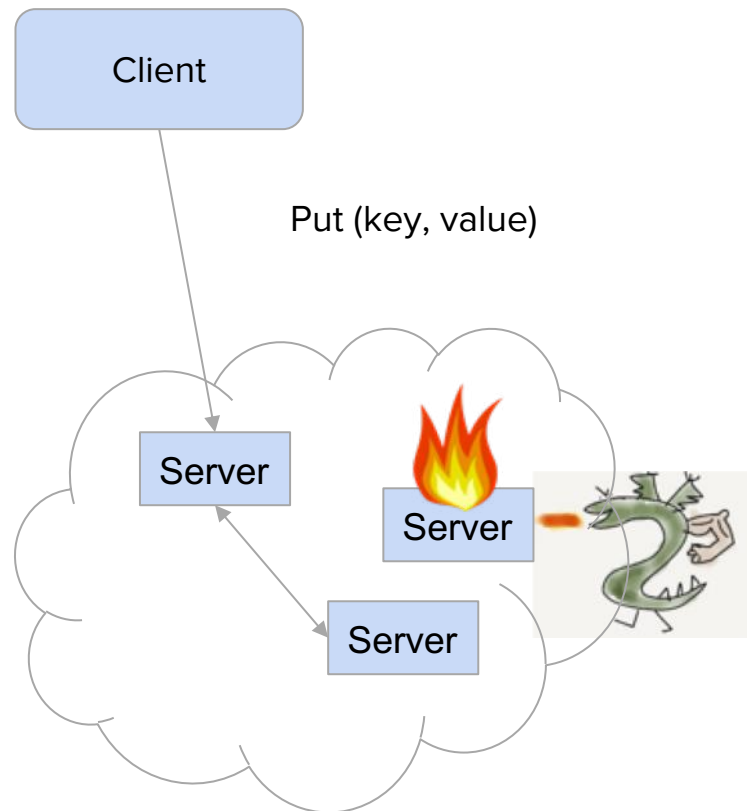https://content.pivotal.io/pivotal-gemfire/scaling-data-services-with-pivotal-gemfire

Pivotal

SpringOne Platform )»Pivotal **re:Cap Seoul 2019**

# What is Geode?

- Distributed key-value store

Client

Put (key, value)

Server

Server

Server

# What is Geode?

- Distributed key-value store
- High available

Client

Put (key, value)

Server

Server

Server

SpringOne Platform by Pivotal re:Cap Seoul 2019

# What is Geode?

- Distributed key-value store
- High available
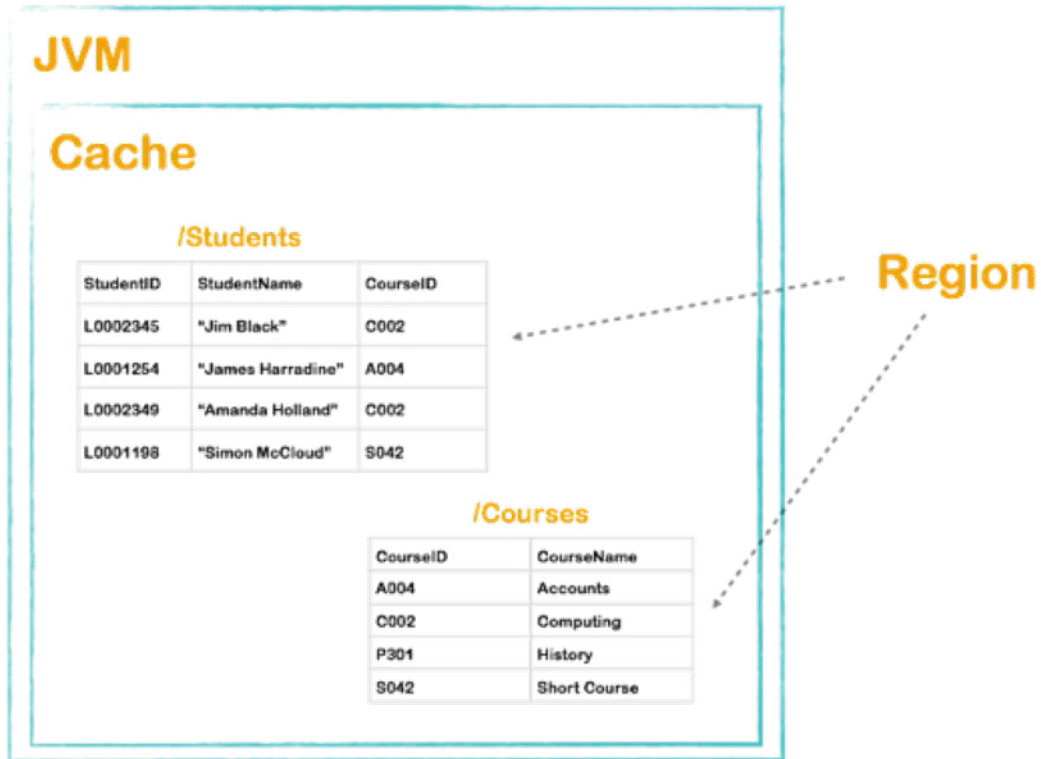- Low Latency

Whoah!

Client

Put (key, value)

< 1ms

Server

Server

# What is Geode?

- Distributed key-value store
- High available
- Low Latency
- Consistent and Partition Tolerant

Client

Put (key, value)
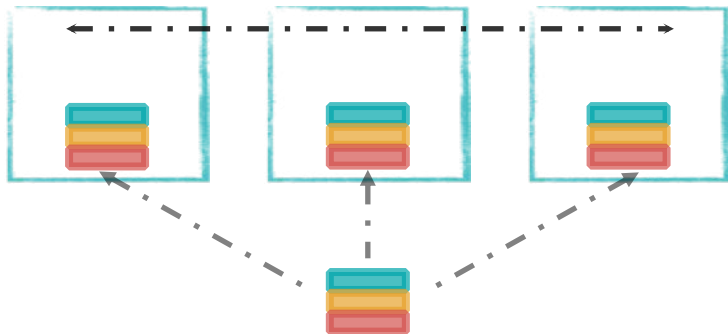
Server

Oh, no! A network partition!

Server

# Regions

- Synonymous to a Table in NoSQL terminology

- Stores Data in **<Key,Value>** pairs with unique Keys
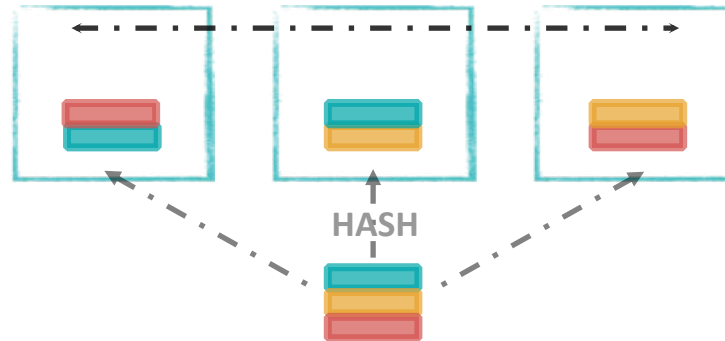
- Divided into buckets across Cache Members

**JVM**

**Cache**

**/Students**

| StudentID | StudentName | CourseID |
|-----------|-------------|----------|
| L0002345 | "Jim Black" | C002 |
| L0001254 | "James Harradine" | A004 |
| L0002349 | "Amanda Holland" | C002 |
| L0001198 | "Simon McCloud" | S042 |

**Region**

**/Courses**

| CourseID | CourseName |
|----------|------------|
| A004 | Accounts |
| C002 | Computing |
| P301 | History |
| S042 | Short Course |

# Regions Data Placement Design



**Replicated**

**Partitioned**

HASH

# Viewing Geode Metrics Before Now

- Internal metrics written to a local file in a proprietary format
    - Viewable with a custom viewing tool
    - Generally viewed after the events happened
- A subset published via JMX

**Main**   **File**   **Chart**   **Template**                                                    **Help**

File: 1. 12_Delta_Solution/stats1007 Windows 7 6.1 x86 blr-bchoudhury GemFire 6.6.2 #build 34761 as of 03/06/2012 17:10:10 PST

| StartTime | File | Samples | Pid | Type | Name |
|---|---|---|---|---|---|
| 06/08 16:09:56 | 1 | 78 | 7468 | VMGCStats | Copy |
| 06/08 16:09:56 | 1 | 78 | 7468 | StatSampler | statSampler |
| 06/08 16:09:56 | 1 | 77 | 7468 | ResourceManagerStats | ResourceManagerStats |
| 06/08 16:09:56 | 1 | 77 | 7468 | PoolStats | client->[anyservers] |
| 06/08 16:09:56 | 1 | 77 | 7468 | ClientStats | ClientStats-client-blr-bchoudhury.vmware.com:40406 |
| 06/08 16:09:56 | 1 | 78 | 7468 | ClientSendStats | ClientSendStats-client-blr-bchoudhury.vmware.com:40406 |
| 06/08 16:09:56 | 1 | 77 | 7468 | CachePerfStats | RegionStats-InventoryItem |
| 06/08 16:09:56 | 1 | 77 | 7468 | CachePerfStats | RegionStats-Customer |
| 06/08 16:09:56 | 1 | 77 | 7468 | Cache | |
| 06/08 16:09:56 | 1 | 77 | 7468 | Cache | |
| 06/08 16:09:56 | 1 | 77 | 7468 | Cache | |
| 06/08 16:09:56 | 1 | 77 | 7468 | Cache | |
| 06/08 16:09:56 | 1 | 77 | 7468 | Cache | |

**cacheListenerCallsCompleted**
cacheListenerCallsInProgress
cacheListenerCallTime
cacheWriterCallsCompleted
cacheWriterCallsInProgress

Chart: Chart 4

**Add Line**

---

**7% Chart 4**

**Chart** **Line**          X:          Y:        cacheListenerCallsCor ▼    PerSecond ━    S: 1.0

points: 77 min: 0.0 max: 1.0 mean: 0.181818 sd: 0.388224

Operations / second — 1, 0.8, 0.6, 0.4, 0.2, 0

16:10:00      16:10:20      16:10:40      16:11:00

1. 12_Delta_Solution/stats1007 blr-bchoudhury GemFire 6.6.2 #build 34761

1. RegionStats-InventoryItem, cacheListenerCallsCompleted/sec

Piv                                                                                    019

# Micrometer in Geode

- Make key metrics visible in external monitoring systems
- View metrics while the system is running
- Augment metrics with details to aid understanding

# Meter Registries

# Connecting to Monitoring Systems

# Connecting to Monitoring Systems

- Chose a Micrometer meter registry implementation that publishes to your monitoring system
  - For a list of implementations, see: http://micrometer.io/docs
- Create a MetricsPublishingService that adds your registry to Geode

# Live View of Metrics

# How Geode Records Measurements

- Uses Micrometer *meters* to record measurements
  - **eventsReceivedCounter.increment(delta);**
- Collects the meters in its *meter registry*

Pivotal

SpringOne Platform by Pivotal re:Cap Seoul 2019

# Current Geode Meters

- System metrics
  - Processors
  - File descriptors
- JVM metrics
  - Process
  - Memory
  - Garbage collection
  - Threads
- Cache Entries

# jvm.memory.used

# Scenario – Partition Region Redundancy 1

SpringOne Platform by Pivotal re:Cap Seoul 2019

**Geode Summit Demo**

⏱ 2019-09-17 16:21:44 to 2019-09-17 16:33:20

**Cache Entries**

- server1 Current: — server2 Current: 298.3 K — server3 Current: 298.3 K

**Get Max Latency**

- server1 hit — server2 hit — server3 hit

**CPU Usage**

- server1 — server2 — server3 — locator

**Put Max Latency**

- server1 — server2 — server3

**Memory Usage**

- server1 PS Eden Space used — server2 PS Eden Space used — server3 PS Eden Space used — locator PS Eden Space used

# The Future Depends on You

- Send Geode's metrics to your favorite monitoring system
- What works well for you?
- What additional metrics would best help you manage your Geode systems?

# For Further Information

- Micrometer: http://micrometer.io/
- Prometheus: https://prometheus.io
- Grafana: https://grafana.com
- Publishing Geode Metrics to External Monitoring Systems:
- https://cwiki.apache.org/confluence/display/GEODE/Publishing+Geode+Metrics+to+External+Monitoring+Systems
- Example code to publish metrics to Prometheus:
- https://github.com/moleske/geode-registry-example

# Performance in Geode

# Performance of Geode 1.9.0



Geode 1.9.0 Average Throughput by Test

Pivotal

Seoul 2019

# Creating the Geode Benchmark – Features

- Run by anyone interested in Geode
- Have others create benchmarks
- Visualize benchmark results over time
- Increase benchmark coverage of Geode

SpringOne Platform by Pivotal **re:Cap Seoul 2019**

# Creating the Geode Benchmark – Goals

- On demand
- Against any revision of Geode
- On AWS cluster deployment of Geode
- On any dev machine in the office
- From Concourse CI pipeline
- With a profiler attached
- Compare two runs of benchmarks for performance changes

# Tests Currently in the Benchmarks

- **ReplicatedGetBenchmark**
- ReplicatedGetLongBenchmark
- **ReplicatedPutBenchmark**
- ReplicatedPutLongtBenchmark
- ReplicatedPutAllBenchmark
- ReplicatedPutAllLongBenchmark
- ReplicatedFunctionExecutionBenchmark
- ReplicatedFunctionExecutionWithArgumentsBenchmark
- ReplicatedFunctionExecutionWithFiltersBenchmark

- **PartitionedGetBenchmark**
- PartitionedGetLongBenchmark
- **PartitionedPutBenchmark**
- PartitionedPutLongBenchmark
- PartitionedPutAllBenchmark
- PartitionedPutAllLongBenchmark
- PartitionedIndexedQueryBenchmark
- PartitionedFunctionExecutionWithArgumentsBenchmark
- PartitionedFunctionExecutionWithFiltersBenchmark

# Tests Currently in the Benchmarks

## https://github.com/apache/geode-benchmarks

Pivotal.

# Other Tested Configurations

- With SSL
- With JDKs: 8, 11, 12, 13
- With Security Manager
- With Garbage Collectors:
  - CMS
  - G1
  - Z
  - Shenandoah
- Adjustable max heap size

# Finding Performance Bottleneck

- Monitor locks
- Thread Park/Unpark Reentrant Locks
- Allocations/GC
- Overuse of synchronization
- Getting a system property in a hot path
- Lazy initialization of objects in a hot path
- Synchronization on a container (ex. Hash map)

# Comparing Performance of 1.9.0 & 1.10.0



Geode 1.9.0 versus 1.10.0 Throughput

# Comparing Performance of 1.9.0 & 1.10.0



Geode 1.9.0 versus 1.10.0 Latency

# reactive programming is about...

"
non-blocking, event-driven applications that scale with a small number of threads with backpressure as a key ingredient that aims to ensure producers do not overwhelm consumers"

- Rossen Stoyanchev
spring.io blog July 28, 2016

Pivotal

# Questions

- Geode as reactive consumer (Subscriber)
  - How do active APIs fit with Geode's can Geode's load shedding policy be adapted to reactive back-pressure

- Geode as reactive producer (Publisher)
  - Again, how do reactive APIs fit Geode's
  - Can Geode produce long data streams incrementally

# Conway's Game of Life

```
162    if (wasAlive) {
163      if (liveNeighborsCount < 2) {
164        return Cell.createDead(newCoordinates); // underpopulation
165      } else if (liveNeighborsCount > 3) {
166        return Cell.createDead(newCoordinates); // overpopulation
167      } else {
168        return Cell.createAlive(newCoordinates, isNewborn: false); // survival
169      }
170    } else {
171      if (liveNeighborsCount == 3) {
172        return Cell.createAlive(newCoordinates, isNewborn: true); // reproduction
173      } else {
174        return Cell.createDead(newCoordinates); // status quo
175      }
176    }
```

# A Reactive Toy

Web
Browser

rsocket-js

Cells

RSocket over WebSocket

spring-boot-rsocket

rsocket-core

rsocket-transport-netty

reactor-core

reactor-tools

reactor-test

Game
Server

Pivotal
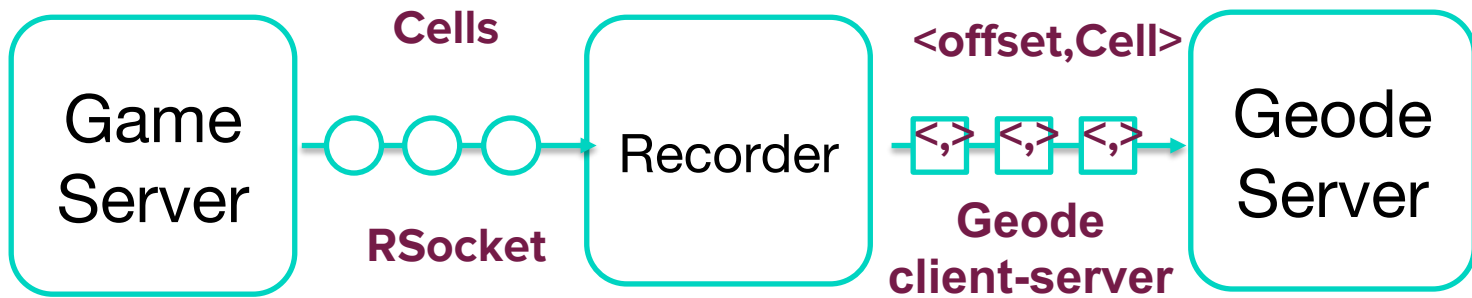
SpringOne Platform by Pivotal **re:Cap Seoul 2019**

```
± |master {1} S:2 U:11 ?:2 ×| → ./gradlew :gameserver:bootRun
```

# Record



```
spring-boot-rsocket          … +
rsocket-core
rsocket-transport-netty      spring-geode-starter
reactor-core                 spring-data-geode-test
reactor-tools                spring-boot-starter-test
reactor-test
```

Pivotal

# Parallel putAll

```
175    return Flux.from(source) Flux<Cell>
176        .limitRequest(generations * coordinateSystem.size()) Flux<Cell>
177        .buffer(coordinateSystem.size()) Flux<List<Cell>>
178        .parallel(parallelism) ParallelFlux<List<Cell>>
179        // NB: gotta runOn() after parallel() to actually schedule work in parallel!
180        .runOn(Schedulers.parallel()) // uncomment to demonstrate BlockHound
181    //    .runOn(Schedulers.elastic()) // uncomment to satisfy BlockHound
182        .doOnNext(
183            bulkCellConsumer) ParallelFlux<List<Cell>>
184        .sequential() Flux<List<Cell>>
185        .doOnTerminate(summarizePerformance(n, starting, firstElementReceived));
```

# {put,putAll} x {serial,parallel}

|          | put | putAll |
|----------|-----|--------|
| serial   | 8   | 23     |
| parallel | 21  | 124    |

thousands of Cells per second
MacBook Pro; 8-way parallelism
de-tuned—for *relative* comparison only

# But remember…

" reactive programming is about…

**non-blocking**…"

- Rossen Stoyanchev
spring.io blog July 28, 2016

SpringOne Platform by Pivotal **re:Cap Seoul 2019**

# Geode ResourceManager



start → (eviction threshold) → evicting → (critical threshold) → critical → (out of memory) → crash

evictions happening

puts + queries serviced

LOAD SHEDDING

# Eviction Config (purposely fragile)

```
68    factory.setEvictionAttributes(
69        /*
70         Region grows until evictionHeapPercentage is reached, then TBD elements are evicted
71         daemon monitors heap memory usage--non-cache actions can result in eviction
72        */
73        EvictionAttributes.createLRUHeapAttributes()
74    );
```

```
22    @CacheServerApplication(name = "AutoConfiguredContinuousQueryIntegrationTests",
23        // FRAGILE
24 //     criticalHeapPercentage = 75f, evictionHeapPercentage = 70f)
25
26 💡   // ROBUST-ISH
27        criticalHeapPercentage = 90f, evictionHeapPercentage = 70f)
```

# Eviction Config (purposely fragile)

```
37    @BeforeClass
38    public static void startGeodeServer() throws IOException {
39
40        startGemFireServer(
41
42            // FRAGILE WITHOUT RETRY, ROBUST WITH RETRY
43            GeodeTestServerConfigurationAutoEvictionNoCQ.class,
44            ...arguments: "-Xmx70m", "-Xms70m",
45
46            // ROBUST (NO RETRY NEEDED)
47 //         GeodeServerConfigurationAutoEvictionAndCQ.class,
48 //         "-Xmx200m", "-Xms200m",
49
50            // While OpenJDK 12 defaults to G1GC now, Geode 1.9 doc says use CMS
51            "-XX:+UseConcMarkSweepGC", "-XX:CMSInitiatingOccupancyFraction=60");
```
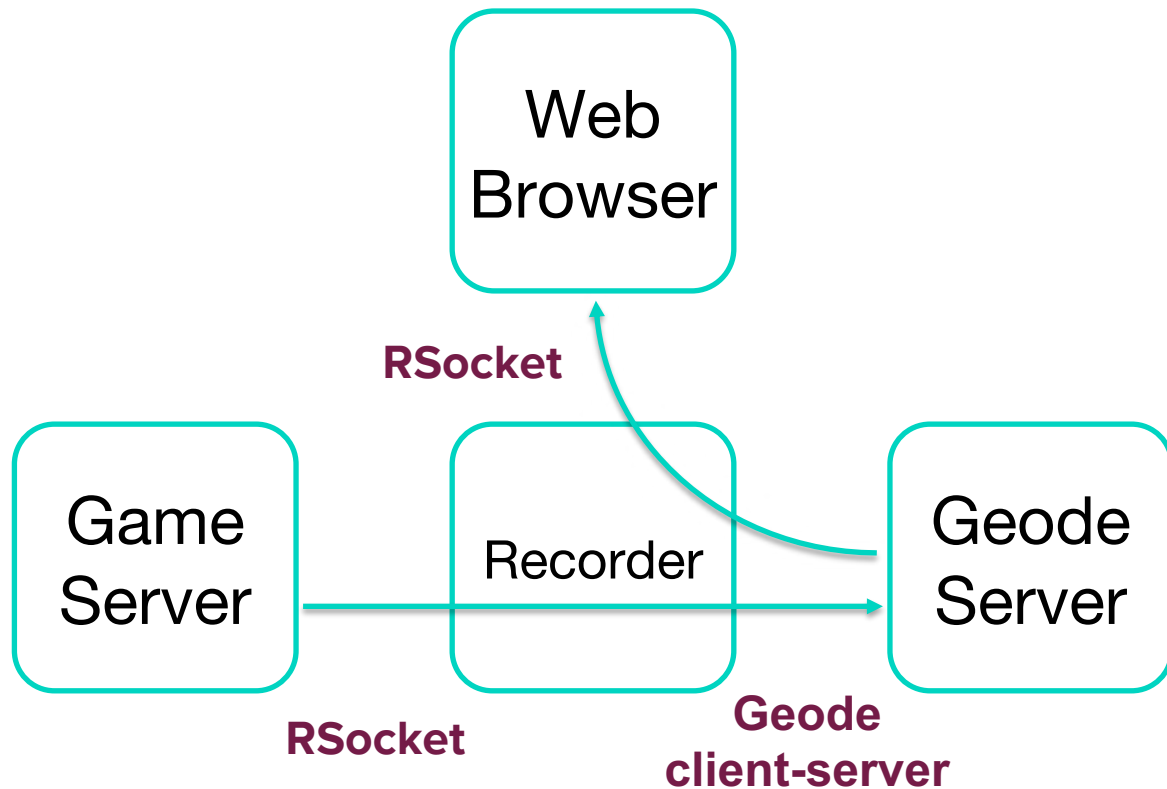
# Playback

# Hot Flux

```java
82    private Publisher<Cell> getCellPublisherNewCellsOnly() {
83
84      return Flux.push(sink ->
85          {
86            try {
87              gemFireCache.getQueryService().newCq( queryString: "SELECT * FROM /Cells",
88                  createCqAttributes(sink)).execute();
89            } catch (final CqException | RegionNotFoundException e) {
90              sink.error(e);
91            }
92          },
93          // BUFFERING == DANGER
94          FluxSink.OverflowStrategy.BUFFER);
95    }

191 @  private static CqListenerAdapter createCqListenerAdapter(final FluxSink<Cell> sink) {
192
193      return (cqEvent) → {
196          if (cqEvent.getBaseOperation() == CREATE) {
197            sink.next((Cell) cqEvent.getNewValue());
198          }
199      };
```

# Endpoint

```java
55   @MessageMapping("/rsocket/all-generations")
56   public Publisher<Cell> allGenerations(final Coordinates _ignored) {
57       return getCellPublisherNewCellsOnly();
58   }
```

Pivotal

SpringOne Platform by Pivotal re:Cap Seoul 2019

# Results I

- Geode as reactive consumer (Subscriber)
  - How do active APIs fit with Geode's can Geode's load shedding policy be adapted to reactive back-pressure
    - Saw serial, parallel, `put`, `putAll`
    - BlockHound + `Schedulers.boundedElastic()`
  - Can Geode's load shedding policy be adapted to reactive back0pressure
    - App design, capacity planning, testing as usual
    - Back off on `LowMemoryException`

# Results II

- Geode as reactive producer (Publisher)
  - Again, how do reactive APIs *fit* with Geode's
    - Demonstrated a simple CQ -> hot Flux
  - Can Geode produce long data streams incrementally
    - Query results chunked cold/cold-ish Fluxes present some challenges
    - CQ ordering of initial results would be nice