Please read the notes associated with each slide for the full context of the presentation

# Who am I?



## Philip Fisher-Ogden

- Director of Engineering @ Netflix

- Playback Services (making "click play" work)

- 6 years @ Netflix, from 10 servers to 10,000s
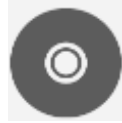
**NETFLIX**

# Story

Netflix streaming – 2007 to present

**NETFLIX**

# Device Growth



| 2007 | 2008 | 2009 | 2010 | 2011+ |
|------|------|------|------|-------|
| 1 device | 10s of devices | 10s of devices | 100s of devices | 1000+ devices |

**NETFLIX**

# Experience Evolution

# Subscribers & Viewing

53M global subscribers

50 countries

>2 billion hours viewed per month

**NETFLIX**

# Virtuous Cycle



Viewing

Improved Personalization

Better Experience

NETFLIX

# Viewing Data

## Who,   What,   When,   Where,   How Long

"city":"PLEASANTON", "region_code":"CA",

10/13/14

Duration

0:15:11

Latest Position

14:41

NETFLIX

# Real time data use cases

## What have I watched?





NETFLIX

# Real time data use cases

## Where was I at?

# Real time data use cases

## What else am I watching?



Too many people are using your account right now.
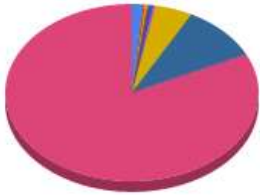
To watch House of Cards, stop playing on this screen:

**iPhone** - Orange Is the New Black (Doppelganger)

Retry

NETFLIX
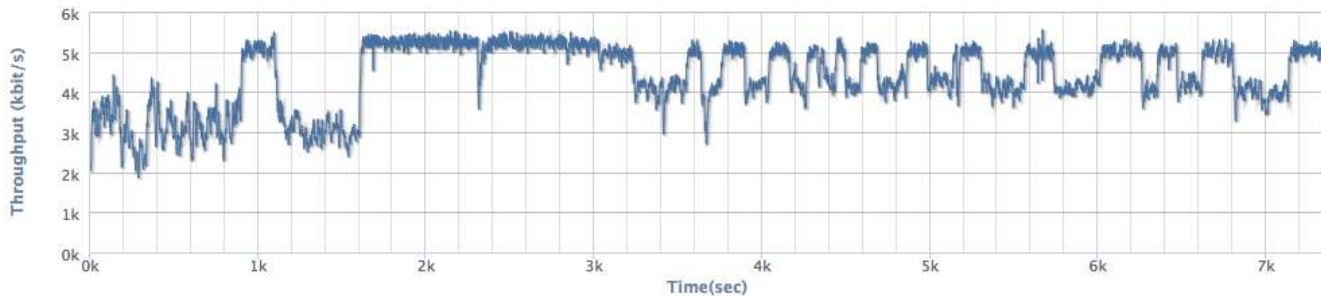
# Session Analytics



Video Bitrate Usage



5%


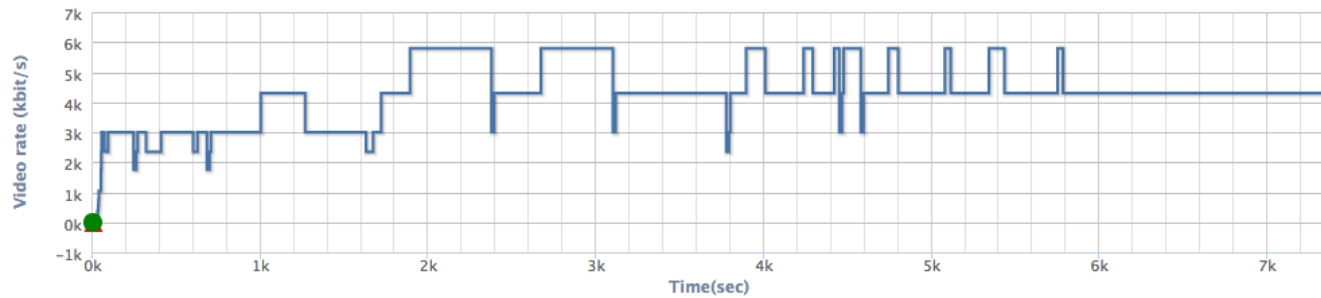
Whoops, something went wrong…

Internet Connection Problem

An Internet or home network connection problem is preventing playback. Please check your Internet connection and try again.
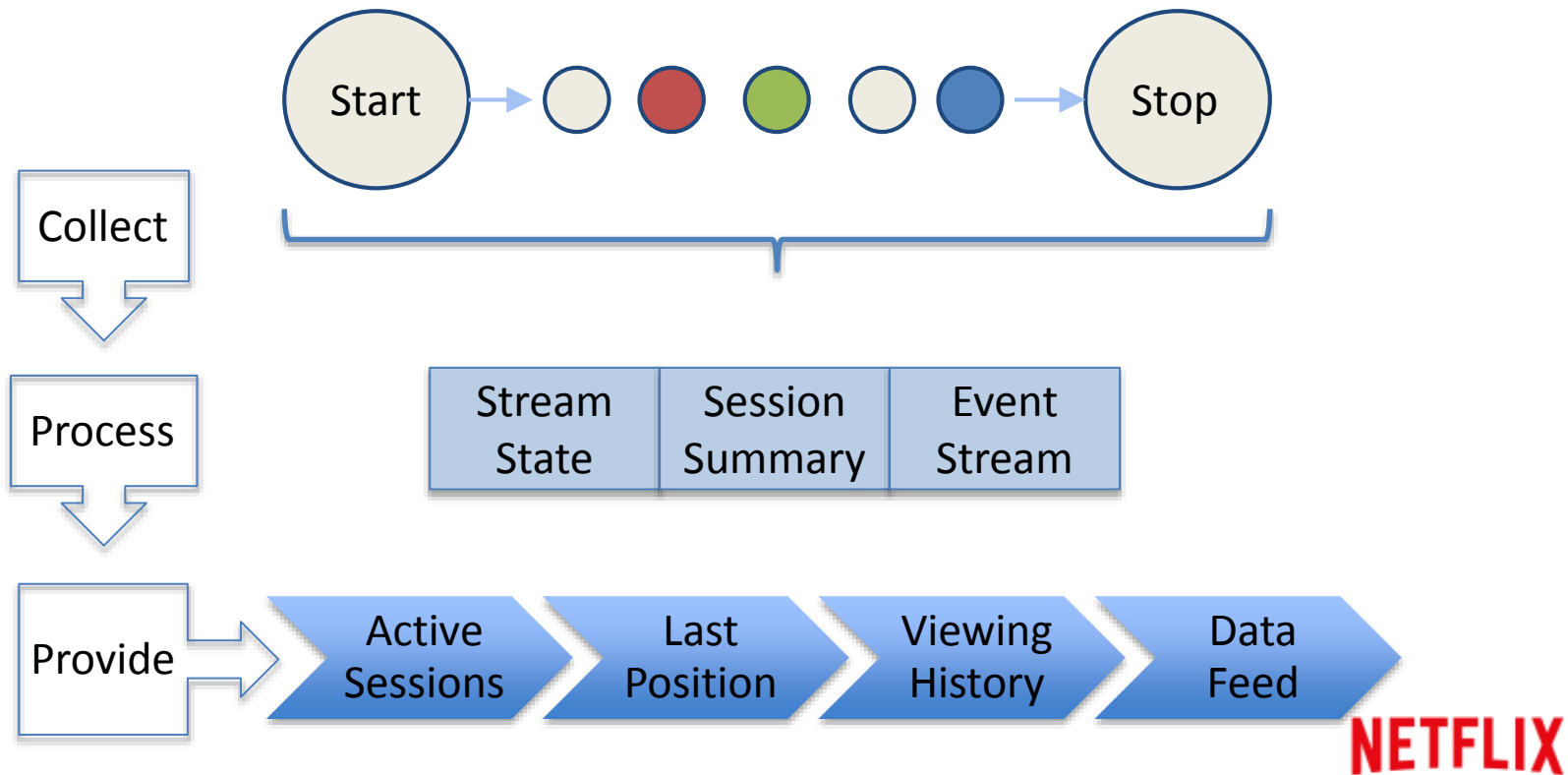
NETFLIX

# Session Analytics

# Generic Architecture
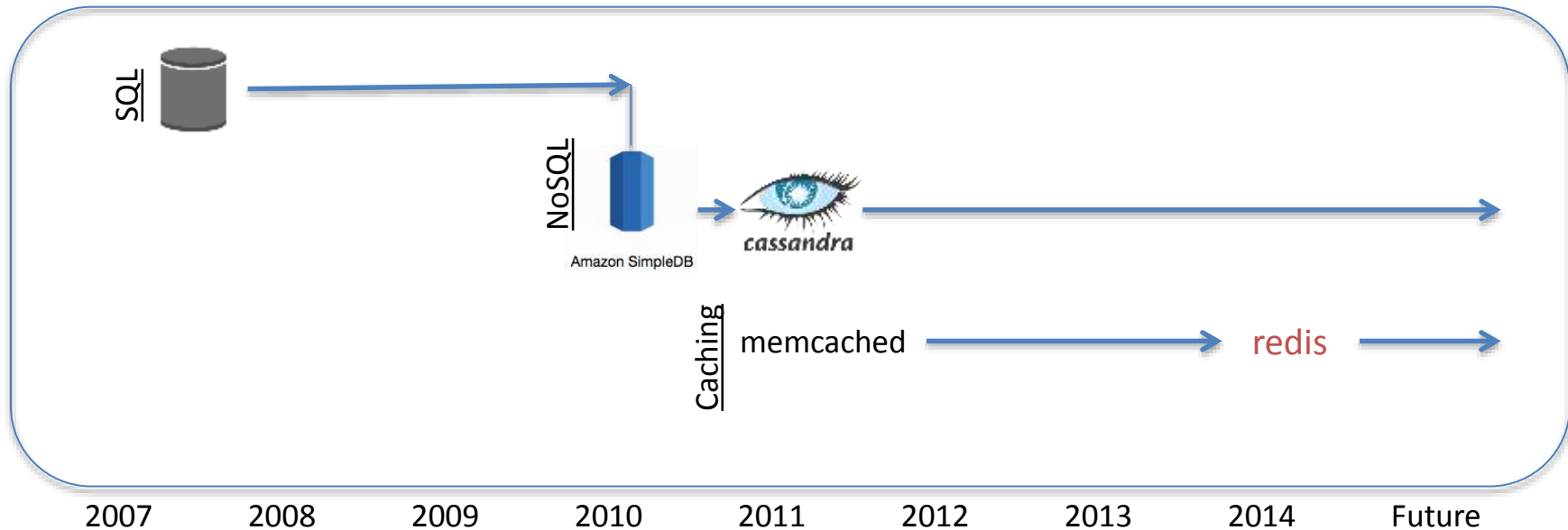
# Architecture Evolution

- Different generations
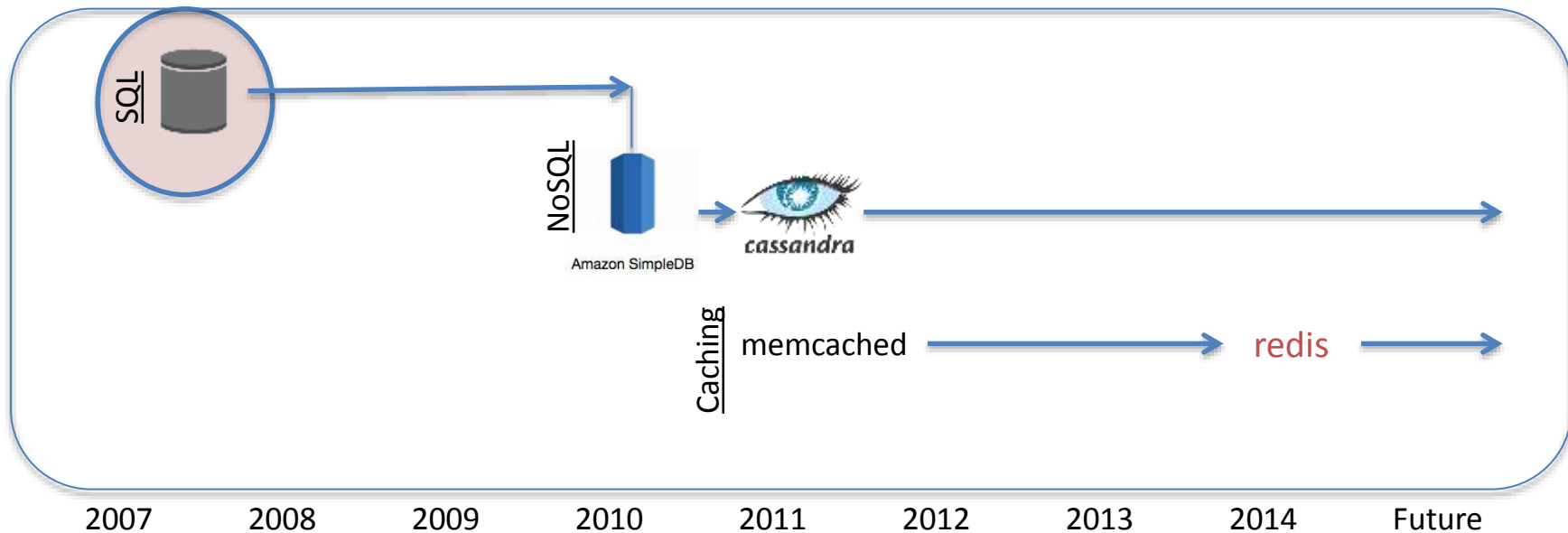
- Pain points & learnings

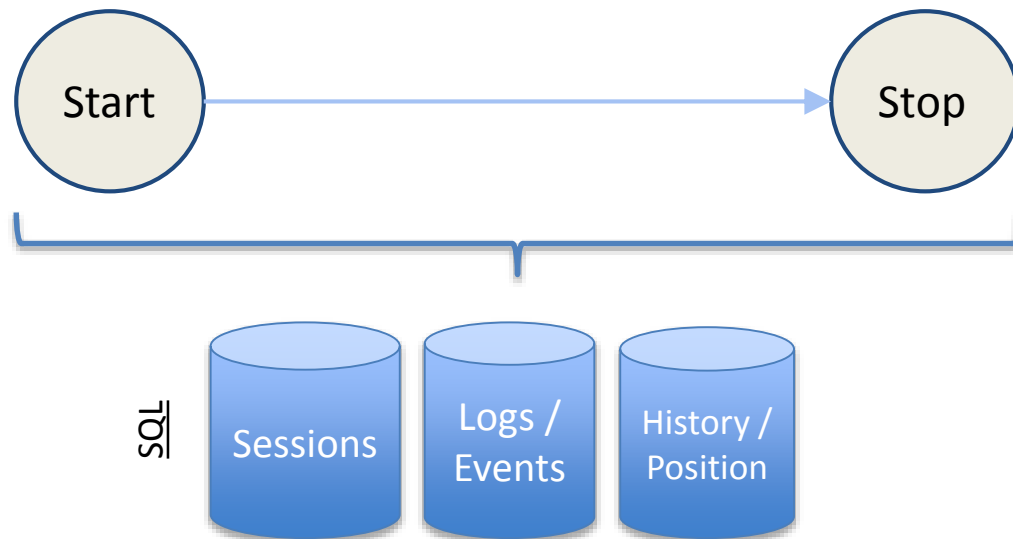- Re-architecture motivations

NETFLIX

# Real Time Data



SQL

NoSQL

Amazon SimpleDB

cassandra

Caching

memcached → redis →

2007  2008  2009  2010  2011  2012  2013  2014  Future

NETFLIX

# Real Time Data – gen 1

SQL

NoSQL

Amazon SimpleDB

cassandra

Caching   memcached   redis

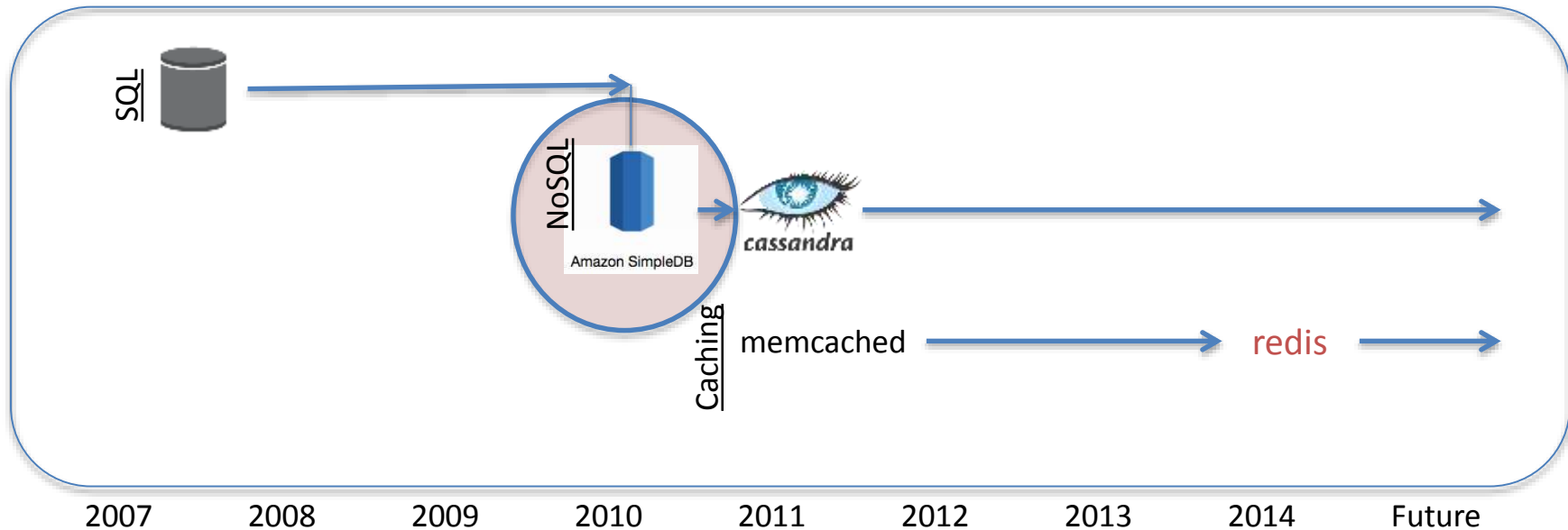2007    2008    2009    2010    2011    2012    2013    2014    Future

NETFLIX

# Real Time Data – gen 1

# Real Time Data – gen 1 pain points

- Scalability
  - DB scaled up not out
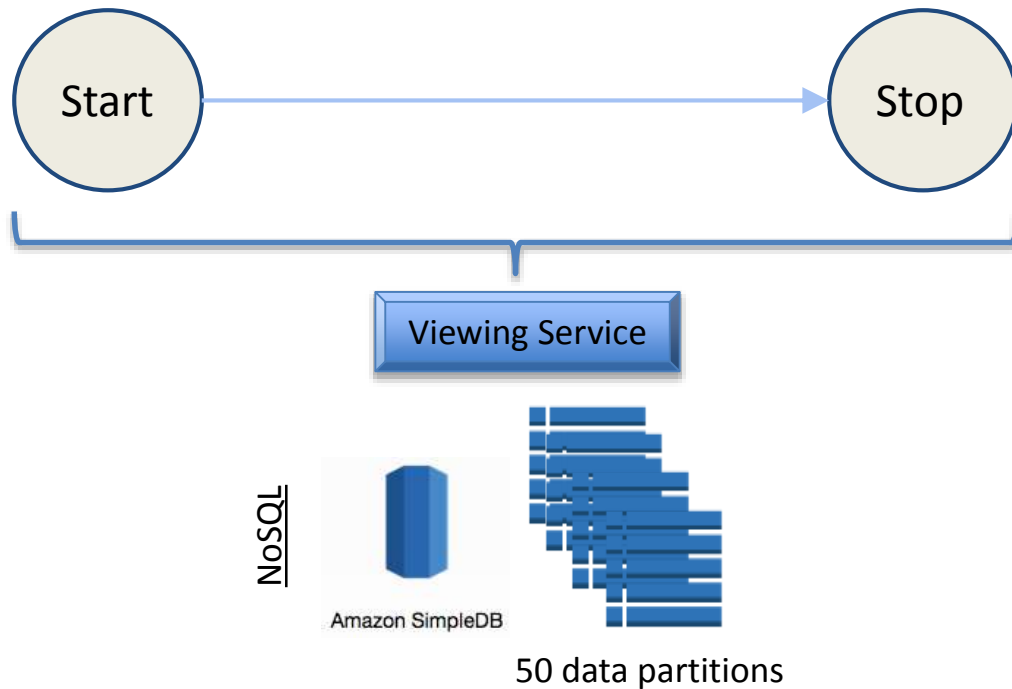- Event Data Analytics
  - ad hoc
- Fixed schema

**NETFLIX**

# Real Time Data – gen 2



SQL

NoSQL

Amazon SimpleDB

cassandra

Caching

memcached    redis

2007    2008    2009    2010    2011    2012    2013    2014    Future

NETFLIX

# Real Time Data – gen 2 motivations

- Scalability
  - Scale out not up
- Flexible schema
  - Key/value attributes
- Service oriented

**NETFLIX**

# Real Time Data – gen 2



Start → Stop

Viewing Service

NoSQL

Amazon SimpleDB

50 data partitions

NETFLIX
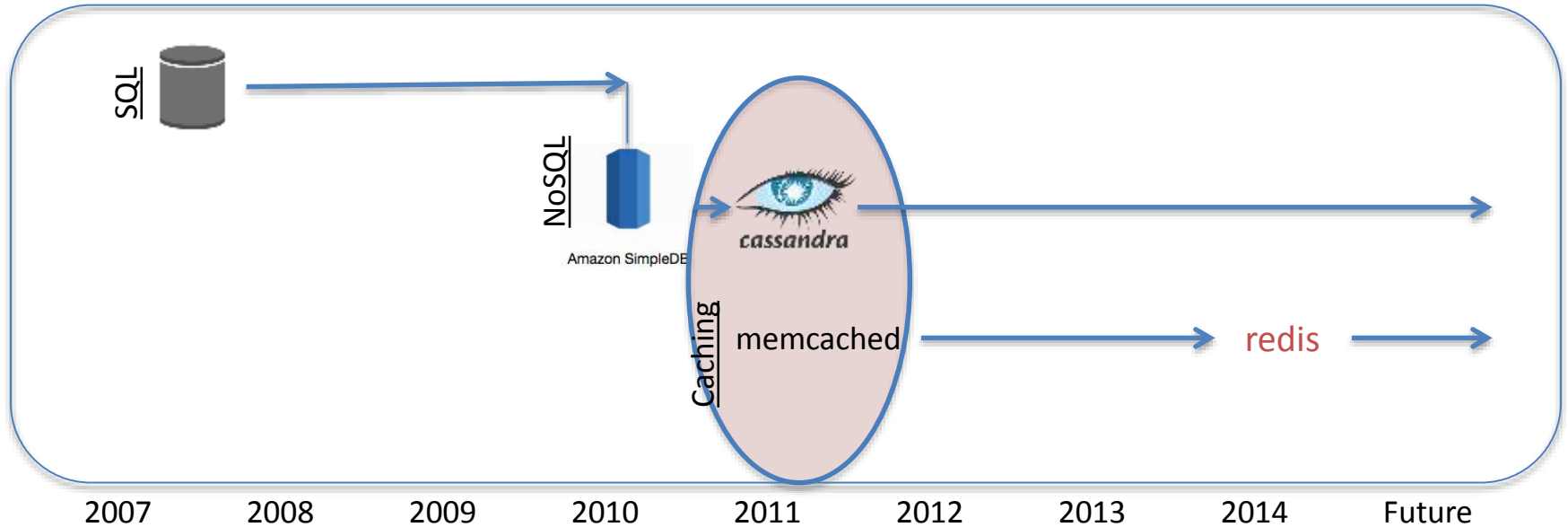
# Real Time Data – gen 2 pain points

- Scale out
  - Resharding was painful
- Performance
  - Hot spots
- Disaster Recovery
  - SimpleDB had no backups

**NETFLIX**

# Real Time Data – gen 3

SQL

NoSQL

Amazon SimpleDB

Caching

cassandra

memcached → redis →

2007   2008   2009   2010   2011   2012   2013   2014   Future

NETFLIX

# Real Time Data – gen 3 landscape

- Cassandra 0.6
- Before SSDs in AWS
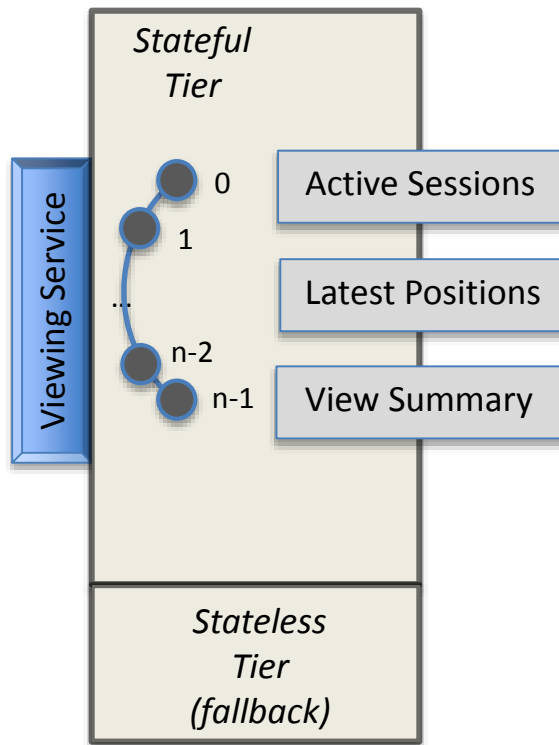- Netflix in 1 AWS region

**NETFLIX**

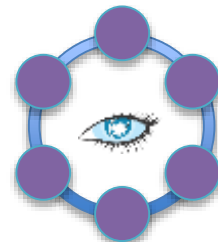# Real Time Data – gen 3 motivations

- Order of magnitude increase in requests

- Scalability
  - Actually scale out rather than up
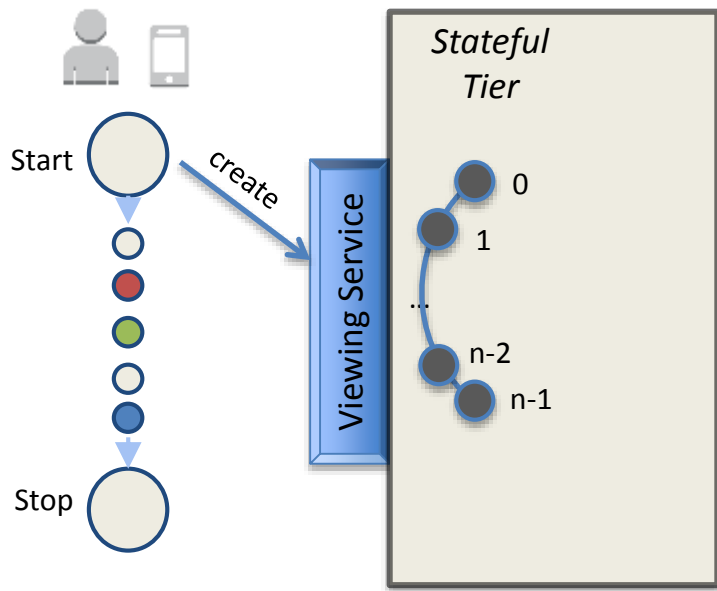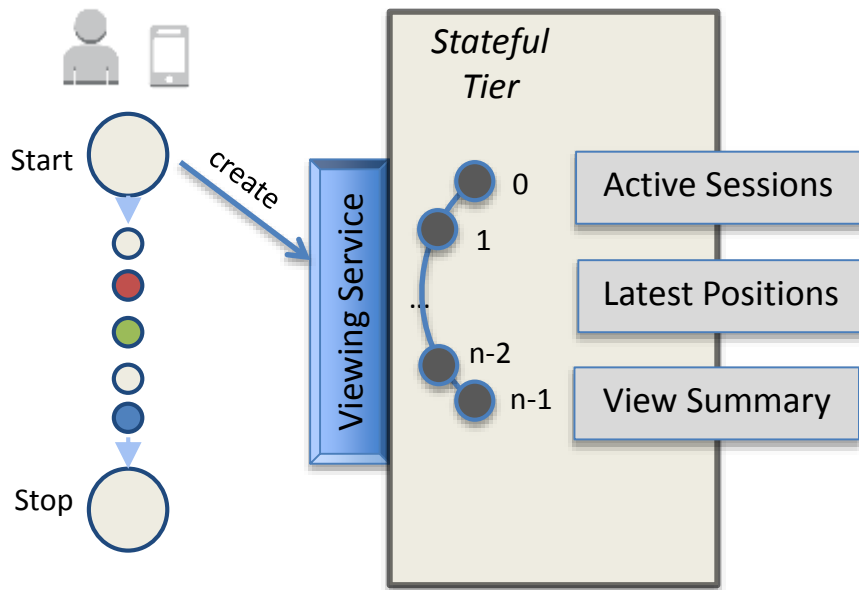
NETFLIX
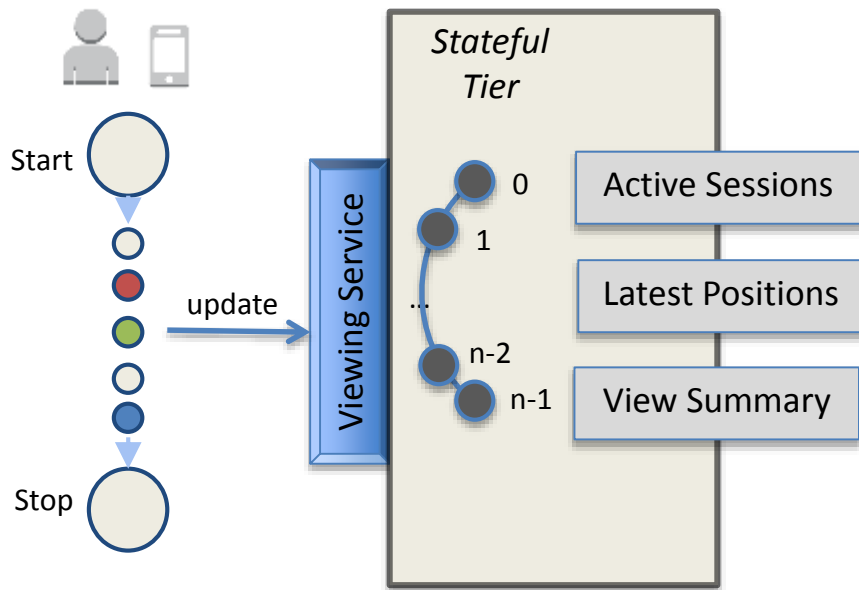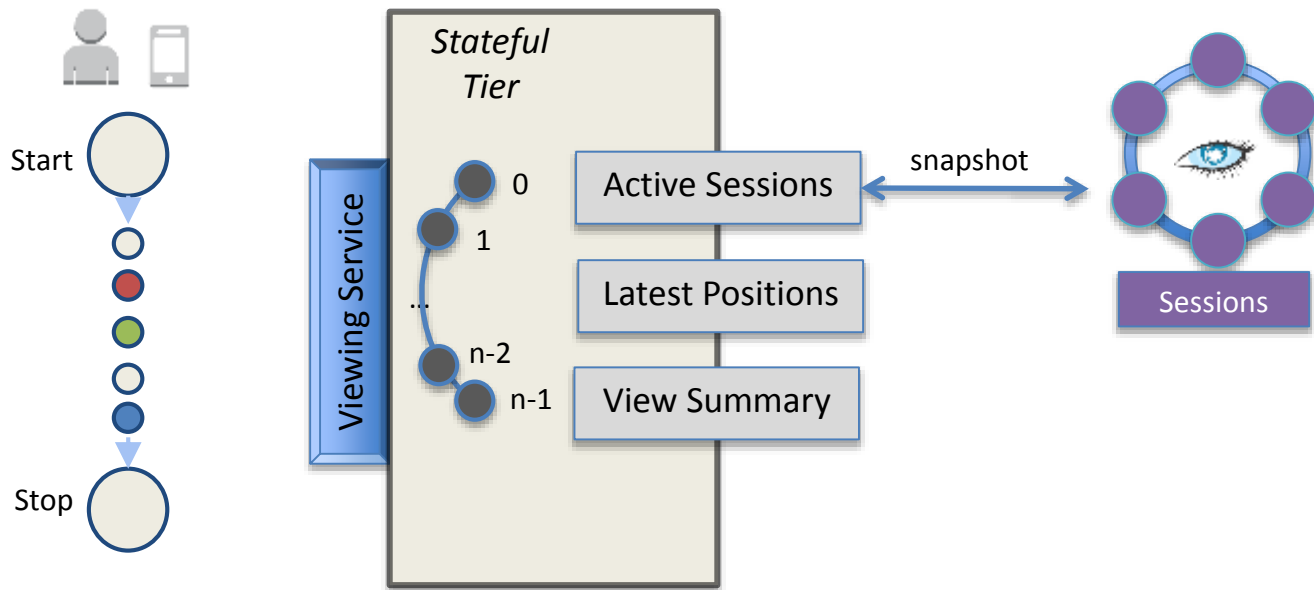
# Real Time Data – gen 3

# Real Time Data – gen 3 writes

# Real Time Data – gen 3 writes

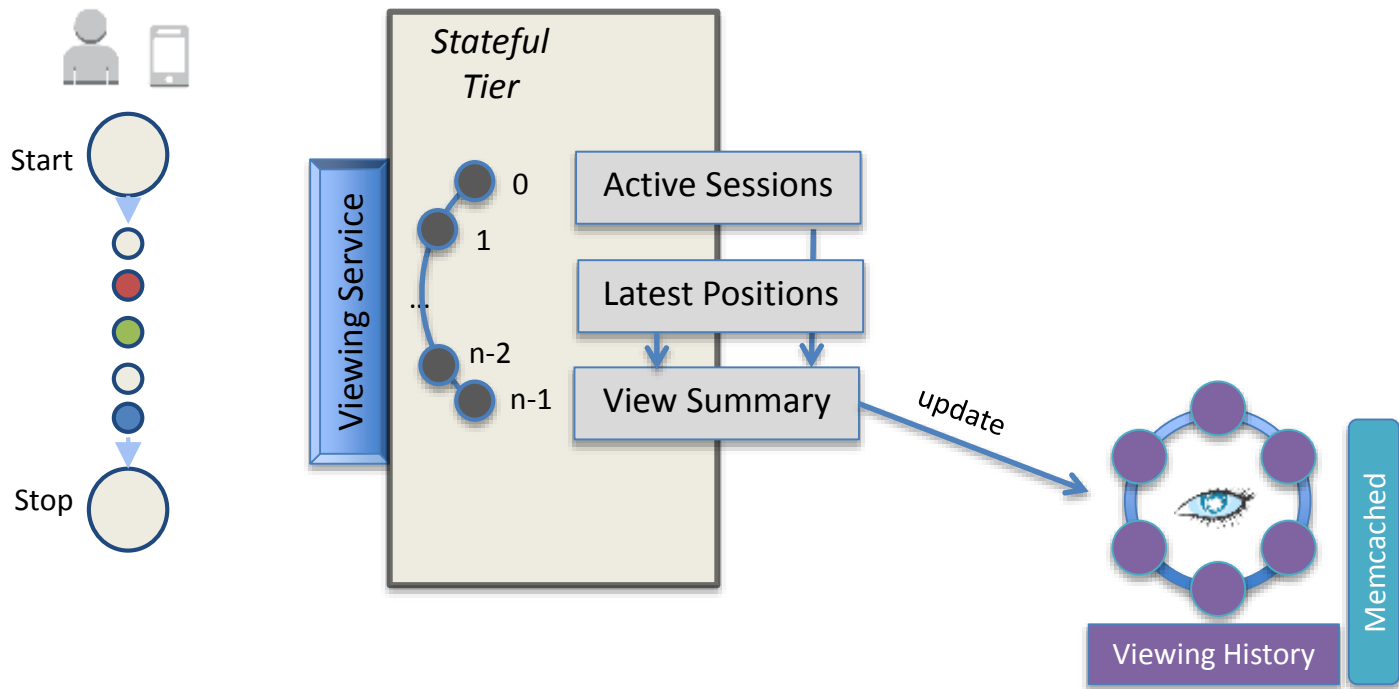# Real Time Data – gen 3 writes

# Real Time Data – gen 3 writes

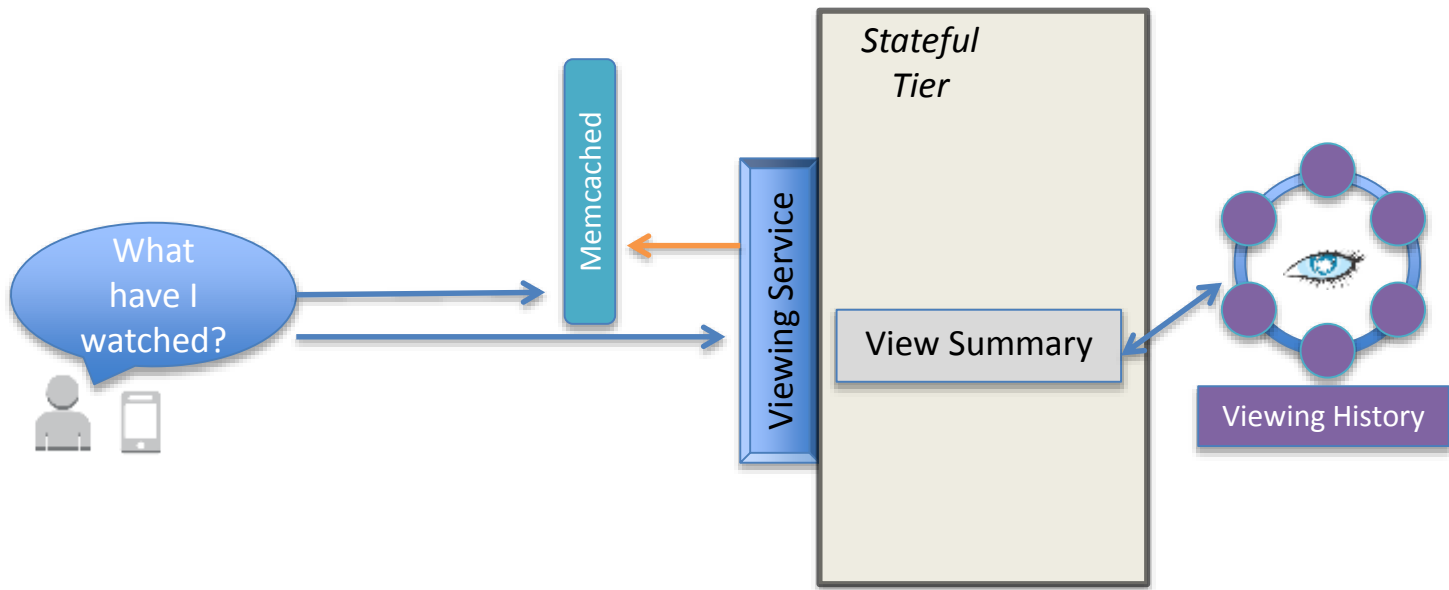# Real Time Data – gen 3 writes
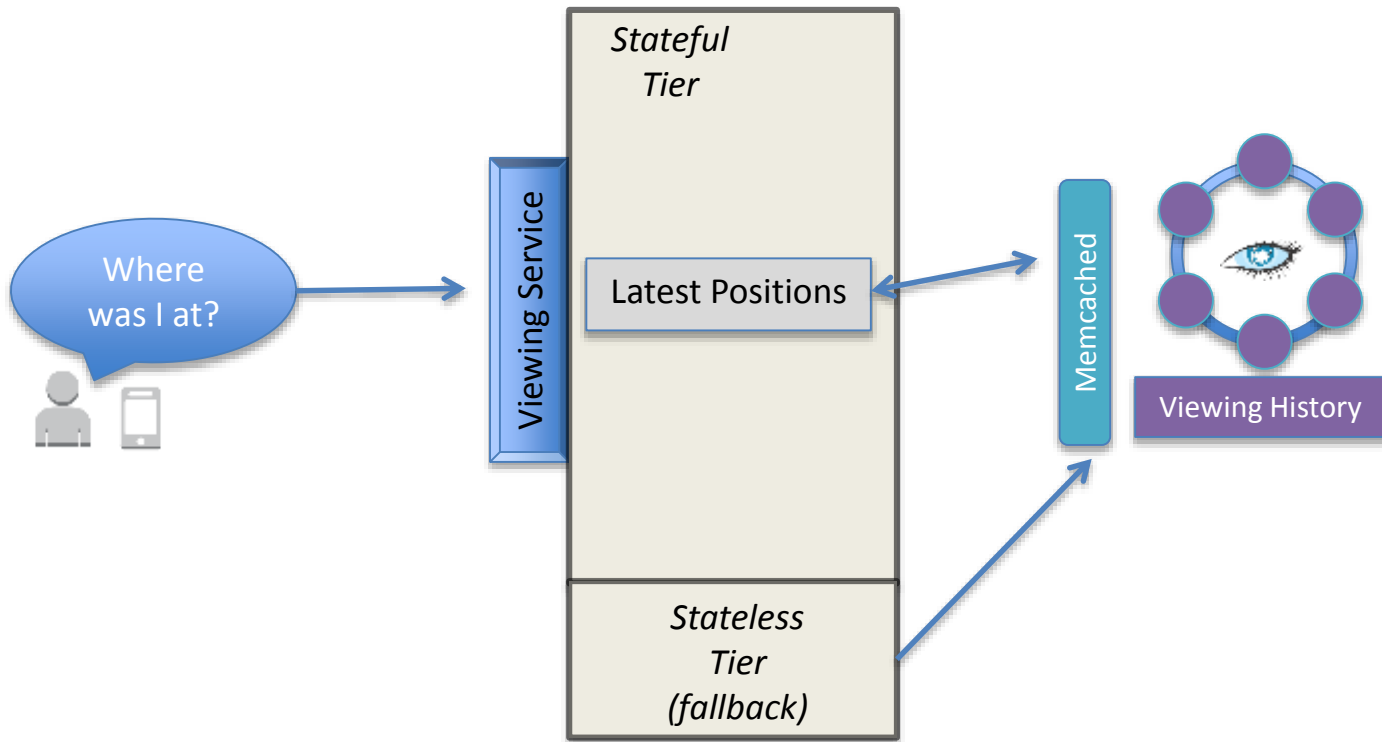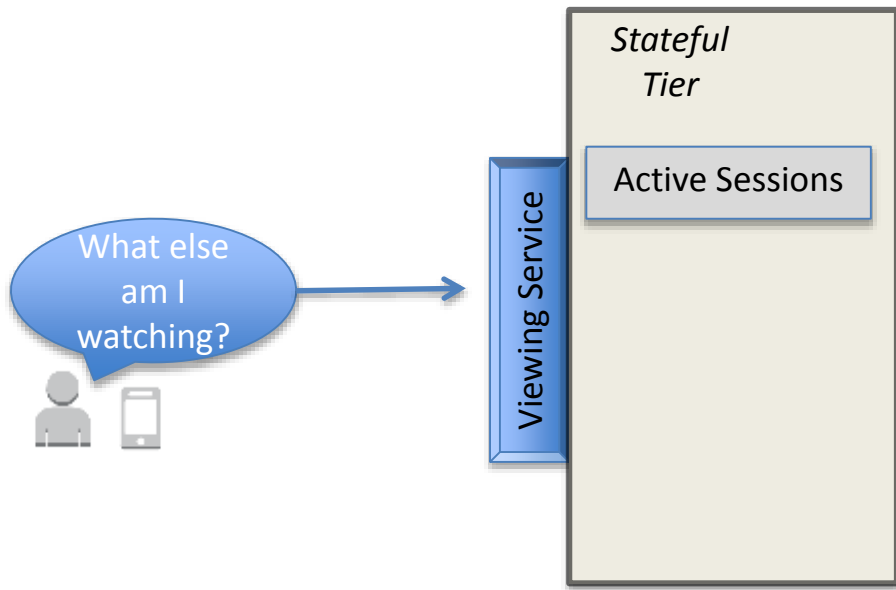
# Real Time Data – gen 3 reads

# Real Time Data – gen 3 reads

# Real Time Data – gen 3 reads

# gen 3 - Requests Scale

| Operation | Scale |
|-----------|-------|
| Create (start streaming) | 1,000s per second |
| Update (heartbeat, close) | 100,000s per second |
| Append (session events/logs) | 10,000s per second |
| Read viewing history | 10,000s per second |
| Read latest position | 100,000s per second |

NETFLIX

# gen 3 – Cluster Scale

| Cluster | Scale |
|---|---|
| Cassandra Viewing History | ~100 hi1.4xl nodes<br>~48 TB total space used |
| Viewing Service Stateful Tier | ~1700 r3.2xl nodes<br>50GB heap memory per node |
| Memcached | ~450 r3.2xl/xl nodes<br>~8TB memory used |

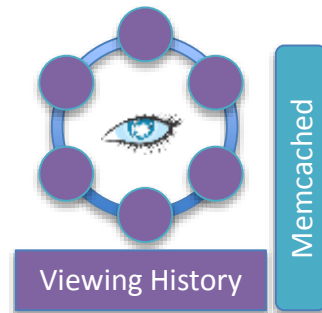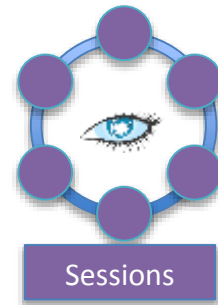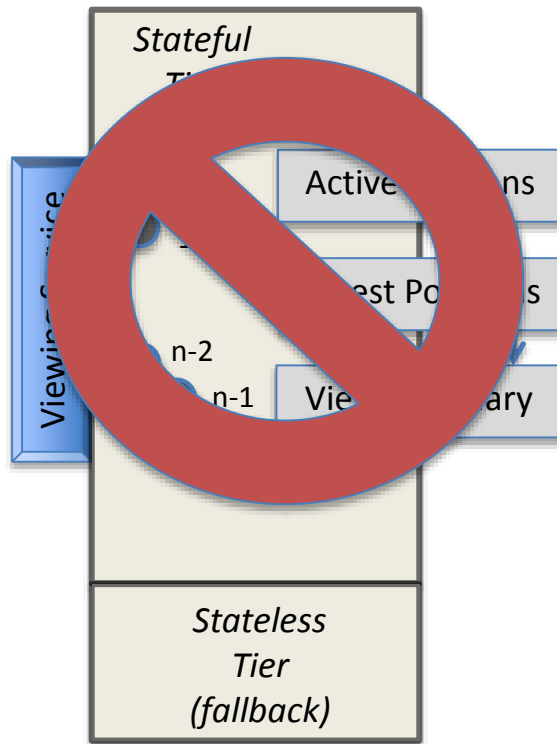NETFLIX

# Real Time Data – gen 3 pain points

- Stateful tier
  - Hot spots
  - Multi-region complexity
- Monolithic service
- read-modify-write poorly suited for memcached

**NETFLIX**

# Real Time Data – gen 3 learnings

- Distributed stateful systems are hard
  - Go stateless, use C*/memcached/redis…
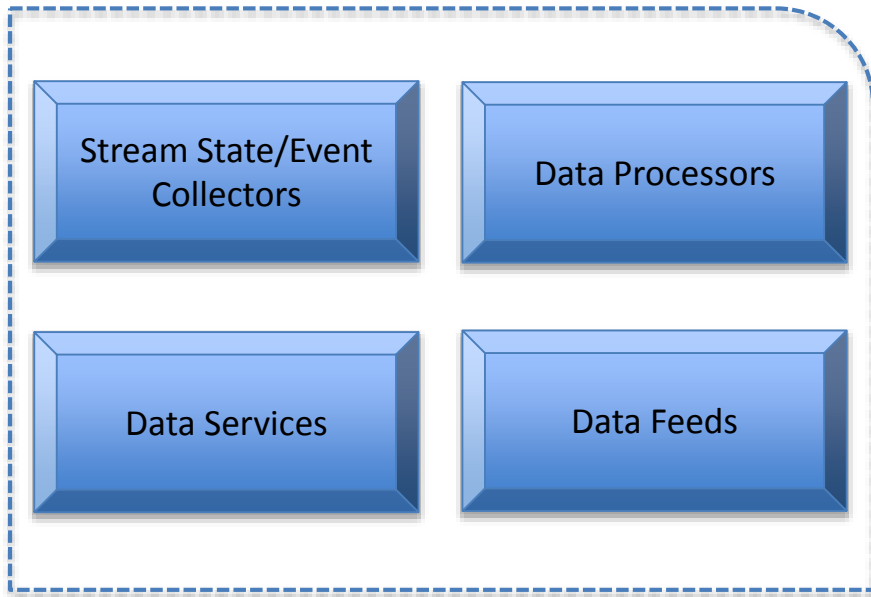- Decompose into microservices

**NETFLIX**

# Real Time Data – gen 4
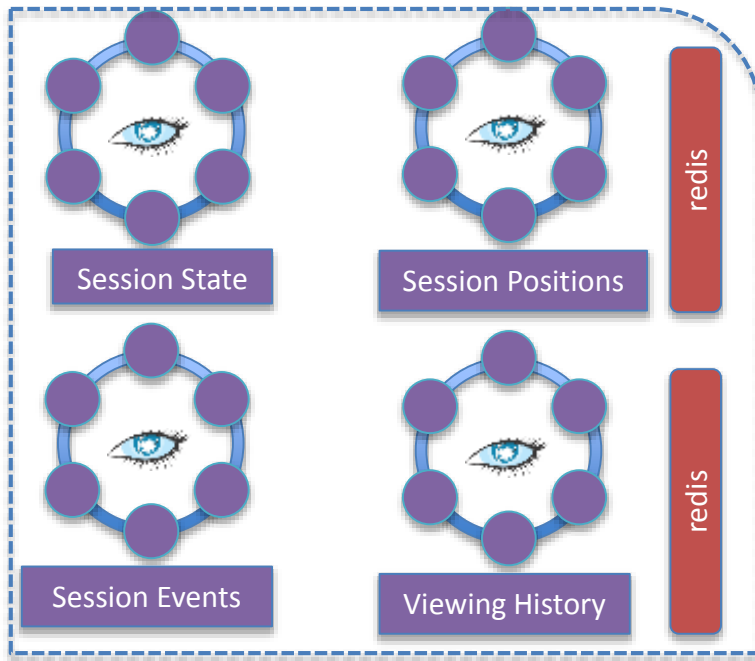
# Real Time Data – gen 4

Stateless Microservices

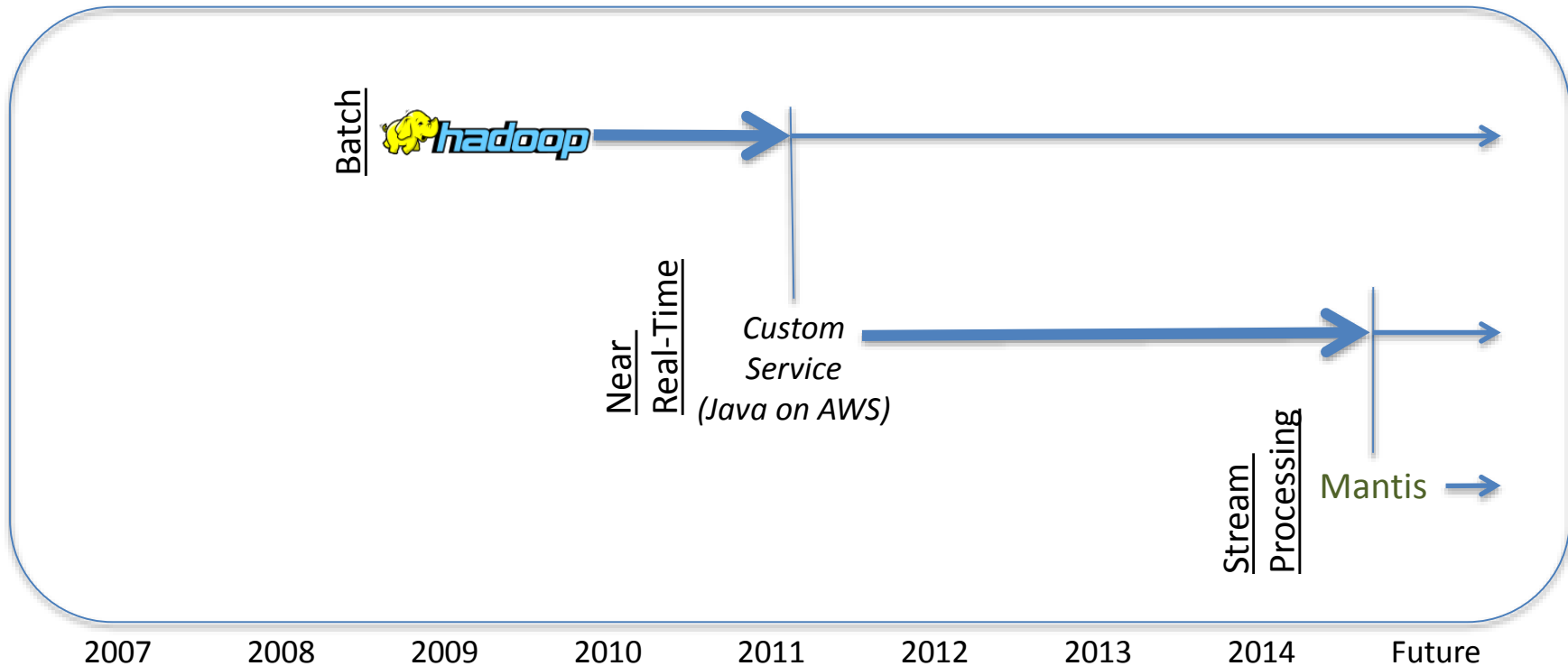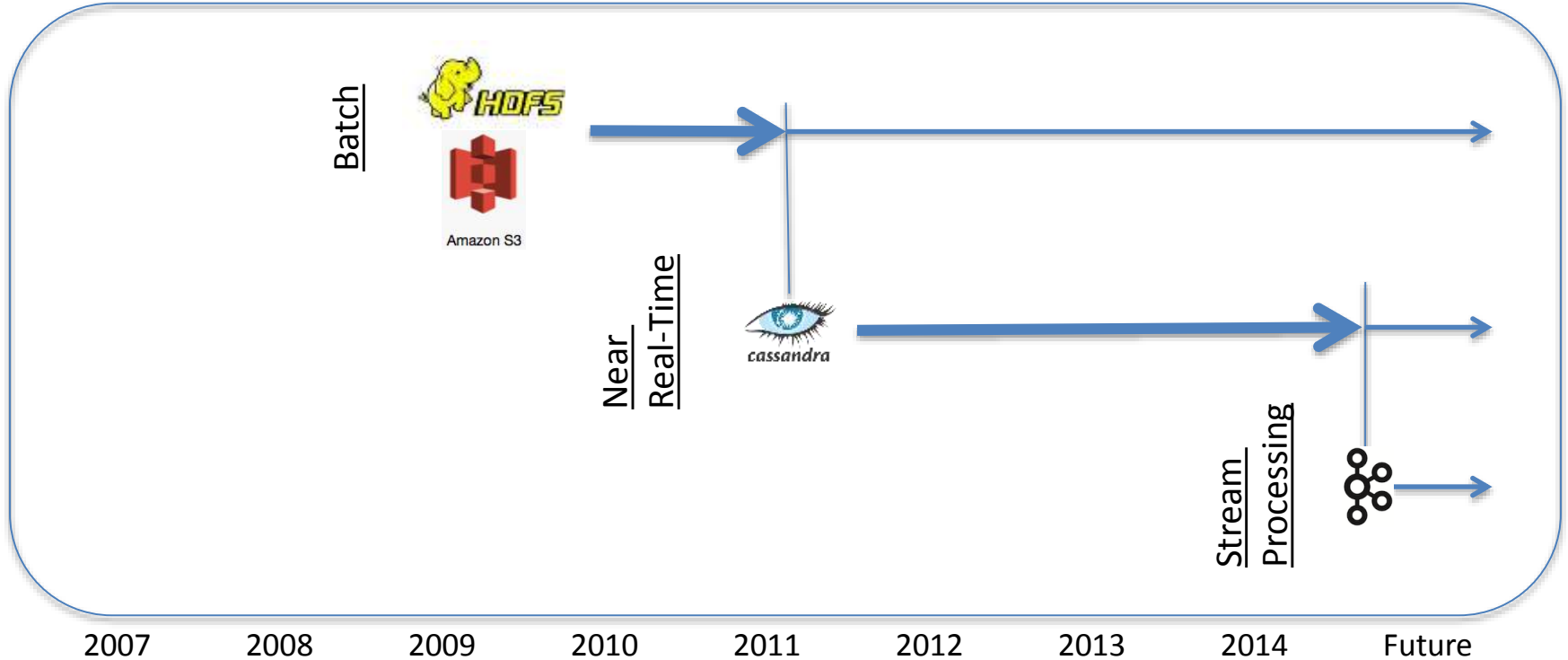| | |
|---|---|
| Stream State/Event Collectors | Data Processors |
| Data Services | Data Feeds |

NETFLIX

# Real Time Data – gen 4

# Session Analytics

- Summarize detailed event data

- Non-real time, but near real time

- Some shared logic with real time

NETFLIX

# Session Analytics - Processing



Batch

Near Real-Time

Custom Service
(Java on AWS)

Stream Processing

Mantis

2007    2008    2009    2010    2011    2012    2013    2014    Future
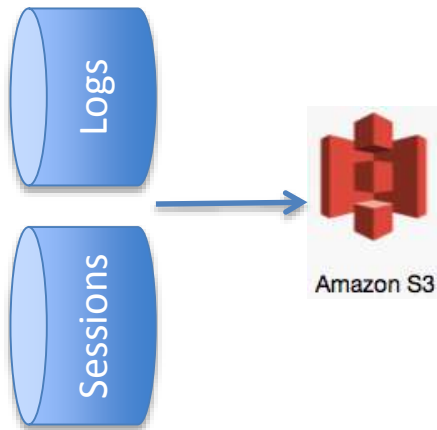
# Session Analytics - Storage

# Session Analytics – gen 1

- Storage
- Processing
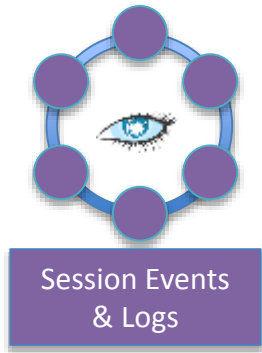


Amazon S3



NETFLIX

# Session Analytics – gen 1 pain points

- MapReduce good for batch
  - Not for near real time
- Complexity
  - Code in 2 systems / frameworks
  - Operational burden of 2 systems

**NETFLIX**

# Session Analytics – gen 2

- Storage

- Processing

Session Events & Logs

AWS

Java

NETFLIX

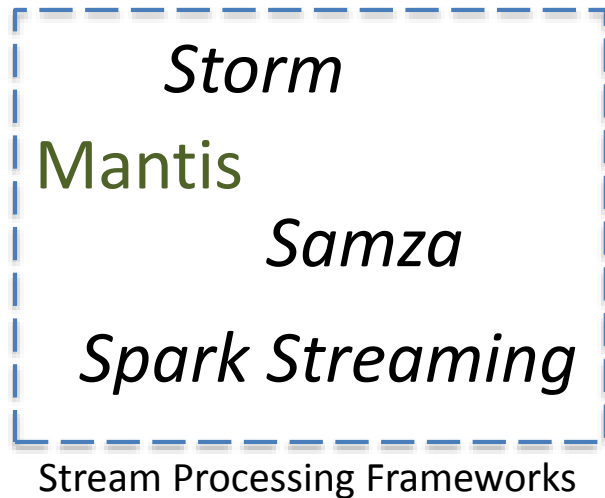# Session Analytics – gen 2 learnings

- Reduced complexity
  - shared code and ops
- Batch still available
- New bottleneck
  - harder to extend logic

NETFLIX

# Session Analytics – gen 3 (*)

- Storage

- Processing



*Storm*

Mantis

*Samza*

*Spark Streaming*

**Stream Processing Frameworks**

**NETFLIX**

# Takeaways

- Polyglot Persistence
  - One size fits all *doesn't* fit all
- Strong opinions, loosely held
  - Design for long term, but be open to redesigns

NETFLIX

# Thanks!

@philip_pfo

**NETFLIX**