QCON
London 2013

# My 'number'

# #6752

Data On Demand

PUSH
TECHNOLOGY

- Distributed Systems / HPC guy.

- **Chief Scientist** :- at **Push Technology**

- Alumnus of :-

  Motorola, IONA, Betfair, JPMC, StreamBase.

- School: Trinity College Dublin.
  - BA (Mod). Comp. Sci.+
  - M.Sc. Networks & Distributed Systems

- Responds to: Guinness, Whisky

# About me?

Darach@PushTechnology.com

Highly Available
Real-Time
Data Distribution
at the network's
Edge

# HA. Black swans & Pigfish!

**PUSH**
TECHNOLOGY

# When this …

Service

*esse hic draconum*

Data On Demand
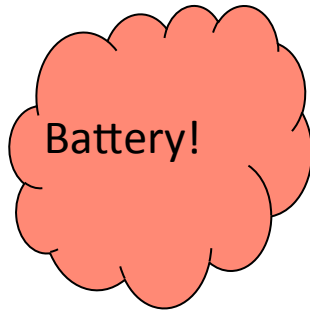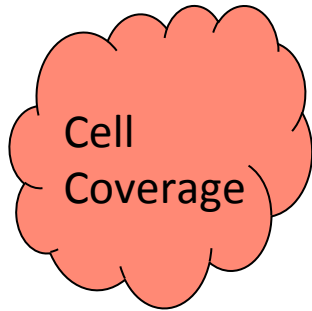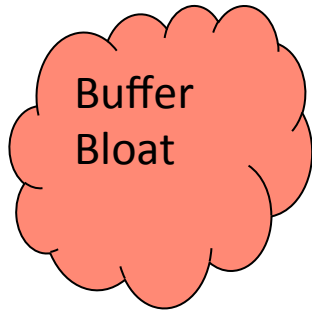
PUSH
TECHNOLOGY

# … inevitably happens
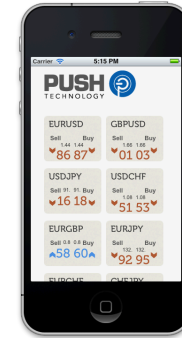
# Over the Edge, where Fallacies roam

- The network is reliable
- Latency is zero
- Bandwidth is infinite
- There is one administrator
- Topology does not change
- The network is secure
- Transport cost is zero
- The network is homogeneous

http://en.wikipedia.org/wiki/Fallacies_of_Distributed_Computing

Data On Demand

PUSH
TECHNOLOGY

# & the last mile. Dragons there be.



Buffer Bloat

Cell Coverage

Battery!

## Download Bandwidth



- Expected (Mbps)
- Download (Mbps)

## Latency

$y = -1.1667x + 45703$
$R^2 = 0.02496$



- Latency (ms)
- Expected (ms)
- Linear (Latency (ms))

Data On Demand

PUSH TECHNOLOGY

# Distributed Systems

"A **distributed system** is one in which the **failure** of a computer **you didn't even know existed** can **render your computer unusable**"

- Leslie Lamport

Data On Demand

**PUSH** ®
TECHNOLOGY

# High Availability

"The price of **reliability** is the pursuit of the **utmost simplicity**."

- Sir Tony Hoare, Turing Award speech, 1980

"This is the unfortunate truth: **Simplicity** is **Hard**."

- Mosely & Marks, Out of the Tar Pit, 2006

Data On Demand

PUSH
TECHNOLOGY

# High **Experiences ~ Man vs Machine**

Image: http://www.phillymarketinglabs.com/man-vs-machine/

Data On Demand

PUSH
TECHNOLOGY

# Machines are 'Easy'

A more unfortunate truth is:

**Humans\*** are **Harder.**

**Allegedly Sergey\* can be *Emasculated by Glass*. But, I bet he owns a god phone.**

*\* Humans are grumpy!*        *\* Colors chosen at random*

# Once Upon a time ...



**SPOF**

There was a very useful little single point of failure.

Data On Demand

**PUSH** ℗ ™
TECHNOLOGY

# It was so successful



**SPOF**

It grew and it grew! So we put it in the cloud.

Data On Demand

PUSH
TECHNOLOGY

# It grew, and it grew …



SPOF

We touched it with all of our devices.

PUSH
TECHNOLOGY

# See SPOF spew ☹



**Wat?**

Is there nothing we can do? Esse hic draconum, basically. But, yes we can too …

Data On Demand

PUSH
TECHNOLOGY

# Bad things happen ...

RapGenius
vs
Heroku.

Anyone?
Ferris?

It's either dyno park this or whale fail in the zoo ...

Data On Demand

PUSH
TECHNOLOGY

# So we spun up another instance

**Master**

**Slave**

And master / slave dabbled with us

# Active/Passive – Maybe?



Active (Hot)

Passive (Cold)

MTTR relatively high. MTBF relatively low.  Ok for 'nice-to-have' or  'too-late-to-cry-now'  facilities.

Data On Demand

PUSH
TECHNOLOGY

# Active/Passive – Hmm?

**Active**
**(Hot)**

**Passive**
**(Warm)**

MTTR slightly better. Can replicate **state** actively for cheaper recovery.

Data On Demand

PUSH
TECHNOLOGY

# Active/Active – Ugh…

**Active**
**(Hot)**

**Active**
**(Hot)**

MTTR negligeable. Use 2x the bandwidth. Get 'half' the value. Dups

# Active/Active – Ugh…

Active
(Hot)

Active
(Warm)

MTTR negligeable. Minify bandwidth usage mostly.  Dups

Data On Demand

PUSH
TECHNOLOGY

# Active/Active – Ugh...

**Active (Hot)**

**Active (Warm)**

MTTR negligeable. Minify bandwidth usage mostly.  Dups

# Spin up another box?

**Active**
**(Hot)**

**Passive**
**(Cold)**

**Standby**
**(Cold)**

Speak no evil

Hear no evil

See no evil

Failback, Failover, Takeover? Load Balance? This is getting costly!

PUSH
TECHNOLOGY

# Can implement N-Plex HA with a 'simple' FSM



Consider roles for Application, Controller, callbacks, …

# Spin up a cluster, ring?



A simple implementation can be easier to build than HA pairs. Your mileage may vary. Can be masterless too.

Where possible, use existing solutions such as Akka, riak_core, doozer, Apache Zookeeper.

You'll still need to roll your own data recovery - mostly. Although CRDT's are changing that in the database community.

# One more thing … Real-Time streams

# Data!



**Meh** ══════════ **Meh**

Real-time data streams are long-lived, conversational, can be stateful.

Data On Demand

PUSH
TECHNOLOGY

# Streaming Data Recovery

| Type | Before | | | After | | | Comments |
|---|---|---|---|---|---|---|---|
| **Precise** | T1 | T2 | T3 | T4 | T5 | T6 | Completely masks failure. |
| **Gap** | T1 | T2 | T3 | ? | T5 | T6 | May result in data loss |
| **Repeating Rollback** | T1 | T2 | T3 | T2 | T3 | T4 | IO preservation. P/S dups are equivalent. |
| **Convergent Rollback** | T1 | T2 | T3 | X2 | X3 | X4 | IO preservation. P/S dups may differ initially but EC. |
| **Divergent Rollback** | T1 | T2 | T3 | X2 | X3 | X4 | IO preservation. P/S dups will differ |

IIBECM

http://nms.lcs.mit.edu/papers/ha-icde05.pdf - HA algos for DSPs

# But, problem: The bird, basically.

**Immediately Inconsistent.**
**But, Eventually Consistent**
… **Maybe.**

**Humans**
**Hate**
**Repeating**
**Rollback**
**Recovery**

**IIBECM**

# Ugh, just one more thing ...

# State?



**Finagle** ──────── **Finagle**

Data rarely passes through a box without being fiddled with.

# Streaming Operations – 4 types

| Type | Comments |
|---|---|
| **Arbitrary** | Completely masks failure. |
| **Deterministic** | Given Same Input, Same State, produce Same Output |
| **Convergent-capable** | Given Same Input, Earlier Start, Empty State, converge to Same Output as Deterministic case |
| **Repeatable** | Given Same Input, Earlier State, Empty State, requires at most one input to product at most one output to be convergent-capable and same as deterministic case |

[http://nms.lcs.mit.edu/papers/ha-icde05.pdf](http://nms.lcs.mit.edu/papers/ha-icde05.pdf) - HA algos for DSPs

PUSH
TECHNOLOGY

# DOH?

# Data Recovery Type x Operation Type = The Hard Part

Data On Demand

PUSH
TECHNOLOGY

# Just one last one more thing …

# Human Expectations

We are fickle, grumpy, changeable, intolerant …

# Let's see. Trading systems?



Market Data → **Prices** → Broker Dealer Service → **Prices** → Fat Cat Banker

Market Execution ↔ "Flow" ↔ Broker Dealer Service

Broker Dealer Service ← **Orders** ← Fat Cat Banker

Broker Dealer Service → **Trades** → Fat Cat Banker

**It Depends…
It's Complicated!**

**Precise Recovery**

**Gap Recovery**

There are at least 26 race conditions in FIX 4.2 Order Cancel/Replace/Reject

And please. No stale data!

# A most inconvenient truth

# Simple is

# Hard

So **relax**. Let it **fail**. Let **live** and let's **learn**.

Data On Demand

PUSH
TECHNOLOGY

# Stop.

- The network is **not** reliable
    **nor** is it cost free.
- Latency is **not** zero
    **nor** is it a democracy.
- Bandwidth is **not** infinite
    **nor** predictable especially the last mile!
- There is **not only** one administrator
    trust, relationships **are** key
- Topology **does** change
    It should, however, converge eventually
- The network is **not** secure
    **nor** is the data that flows through it
- Transport cost is **not** zero
    **but** what you **don't do** is **free**
- The network is **not** homogeneous
    **nor** is it **smart**

# Look.

- **High Availability is not enough for Machine to Human**
- **High Experience requires more**

  Data has a time value – out of date data is bad.

  Data has variable relevance, based on location, device …

  Data rates are faster than human perception

  You cannot send 100mb/sec down 3G networks

  You cannot send 100mb/sec down 3G networks

  You cannot … waste is bad, bloats buffers, and 'slows' data.

  Currency & (Immediate) Consistency are important in M2H.

- **M2H 'High Available Experience' might work for M2M …**

**Data On Demand**

**PUSH** TECHNOLOGY ™

# Listen.

– **Every** nuance comes with a set of tradeoffs.

– Choosing the right ones can be hard, but it pays off.

– Context, Environment are critical

– Break all the rules, one benchmark at a time.


## – **Benchmark Driven Development**

Data On Demand

**PUSH** ℗ ™
TECHNOLOGY

# Action: Data Distribution

Messaging remixed around:

**Relevance**        -        Queue depth for conflatable data should be 0 or 1. No more
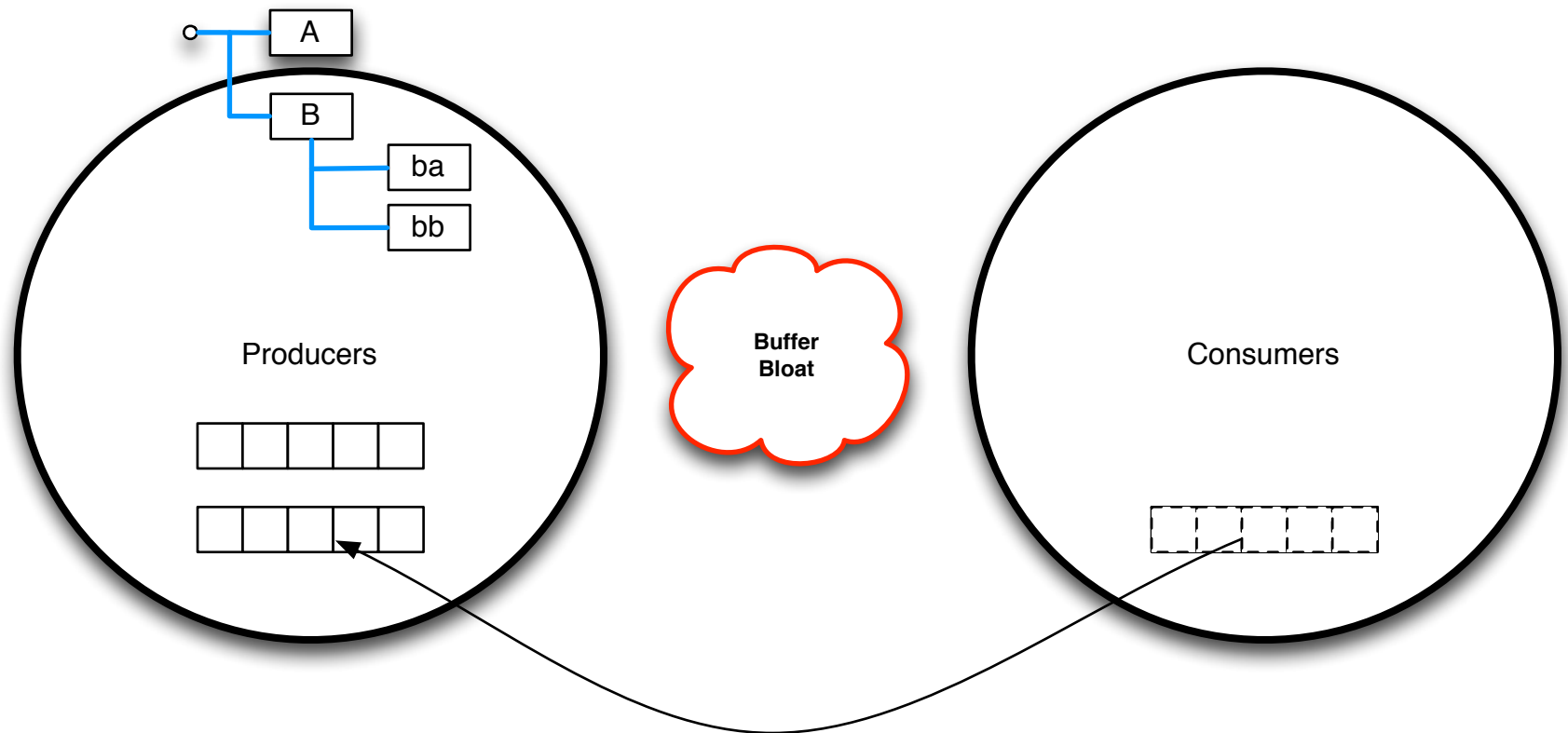
**Responsiveness** - Use HTTP/REST for things. Stream the little things

**Timeliness**      -        It's relative. M2M != M2H.

**Context**          -        Packed binary, deltas mostly, snapshot on subscribe.
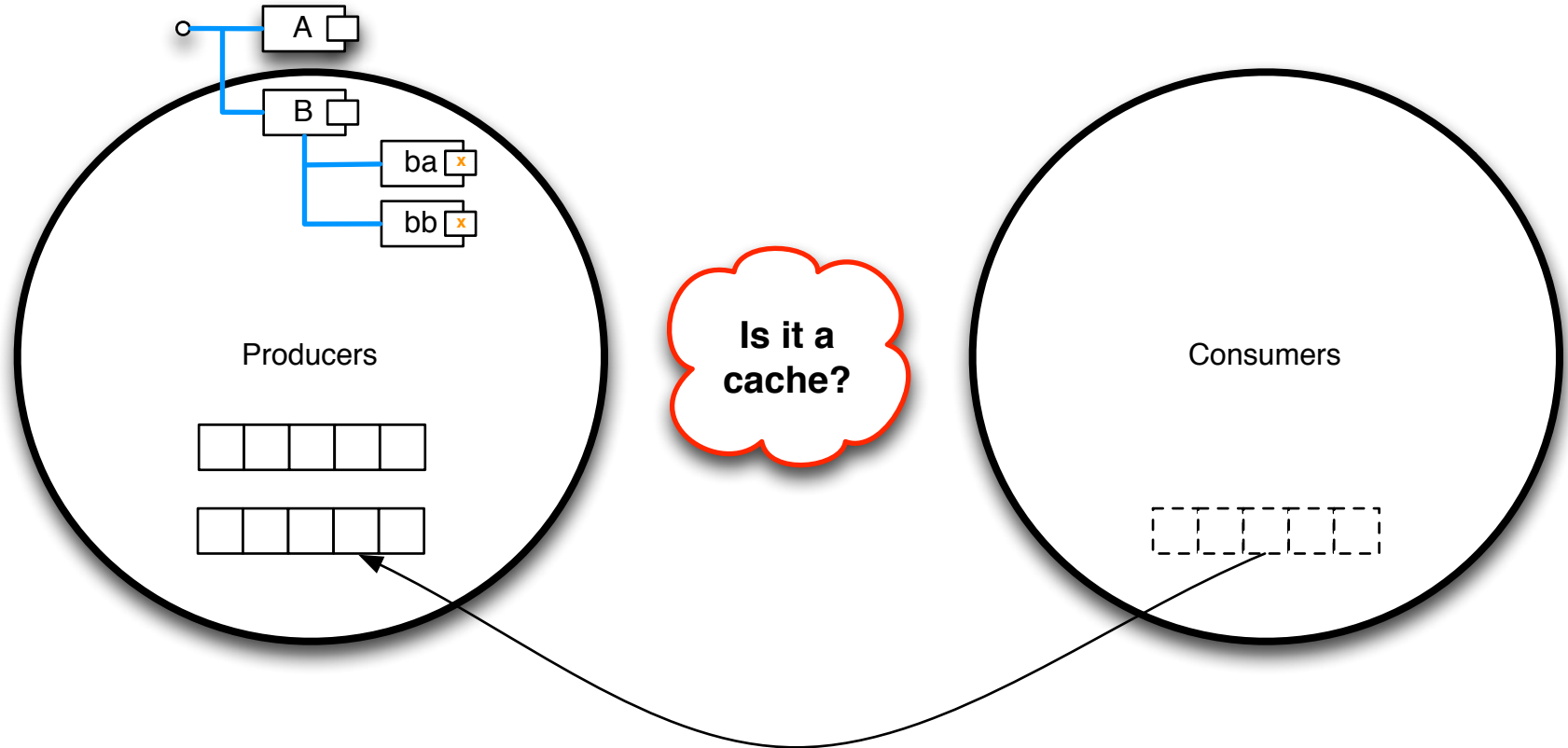
**Environment**-        Don't send 1M 1K events to a mobile phone with 0.5mbps.

# Action. Virtualize Client Queues



Nuance: Client telemetry. Tradeoff: Durable subscriptions harder

Data On Demand

PUSH
TECHNOLOGY

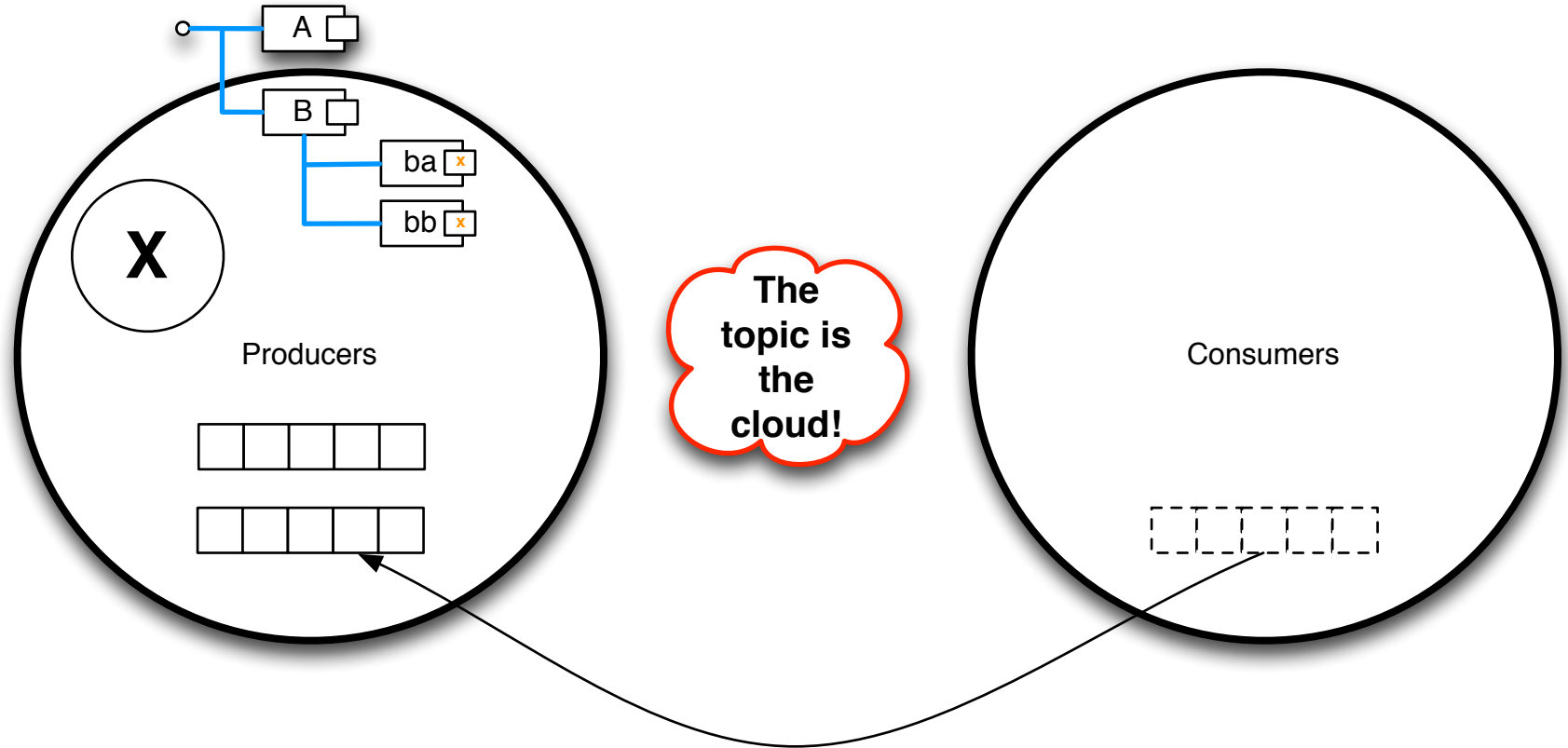# Action. Add data caching



One hop closer to the edge …

# Action. Exploit data structure



Snapshot recovery. Deltas or Changes mostly. Conserves bandwidth

# Action. Behaviors



Extensible. Nuance? Roll your own protocols. Tradeoff? 3$^{rd}$ party code in the engine :/

Data On Demand

PUSH
TECHNOLOGY

# 3$^{rd}$ Party Code in the Engine?

- A bug in 3$^{rd}$ party code can take out the system.
- Is necessary in many environments.
- Can force low density deployments.
- Solutions are appearing, such as Waratek for the JVM.
  - You can hive out 3$^{rd}$ party code into a lightweight fully isolated container.
  - You can manage the container lifecycle.
  - You can reallocate memory, CPU, network.
  - It's language level virtualization.

# Action. Structural Conflation



Ensures only current + consistent data is distributed. Actively soaks up bloat. Extensible!

# Action. Structural Conflation [EEP]



| C2 | A2 | C1 | B1 | A1 | → | 'op' | → | Cop | Bop | Aop | ... | ... |

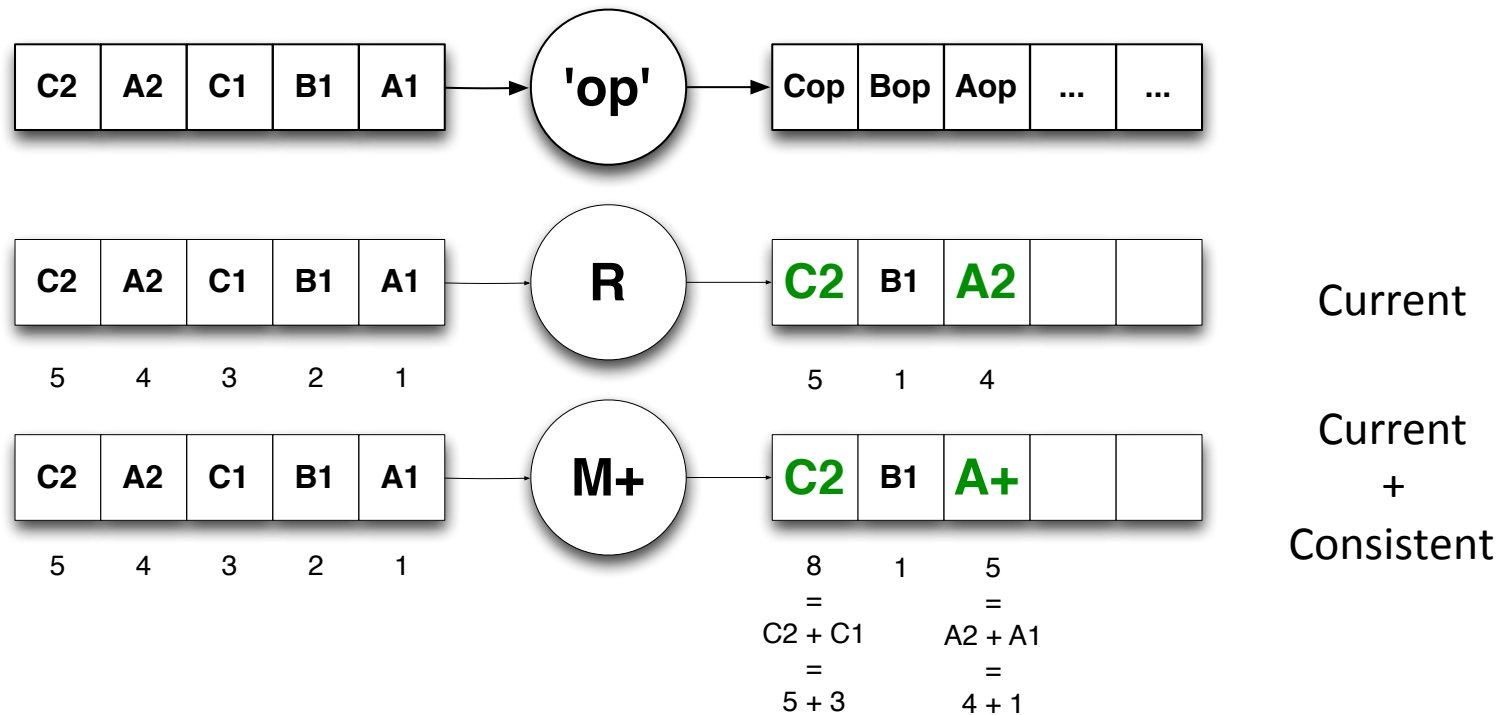| C2 | A2 | C1 | B1 | A1 | → | R | → | **C2** | B1 | **A2** | | |

| 5 | 4 | 3 | 2 | 1 |

| **C2** | B1 | **A2** | | |

| 5 | 1 | 4 |

Current

| C2 | A2 | C1 | B1 | A1 | → | M+ | → | **C2** | B1 | **A+** | | |

| 5 | 4 | 3 | 2 | 1 |

| **C2** | B1 | **A+** | | |

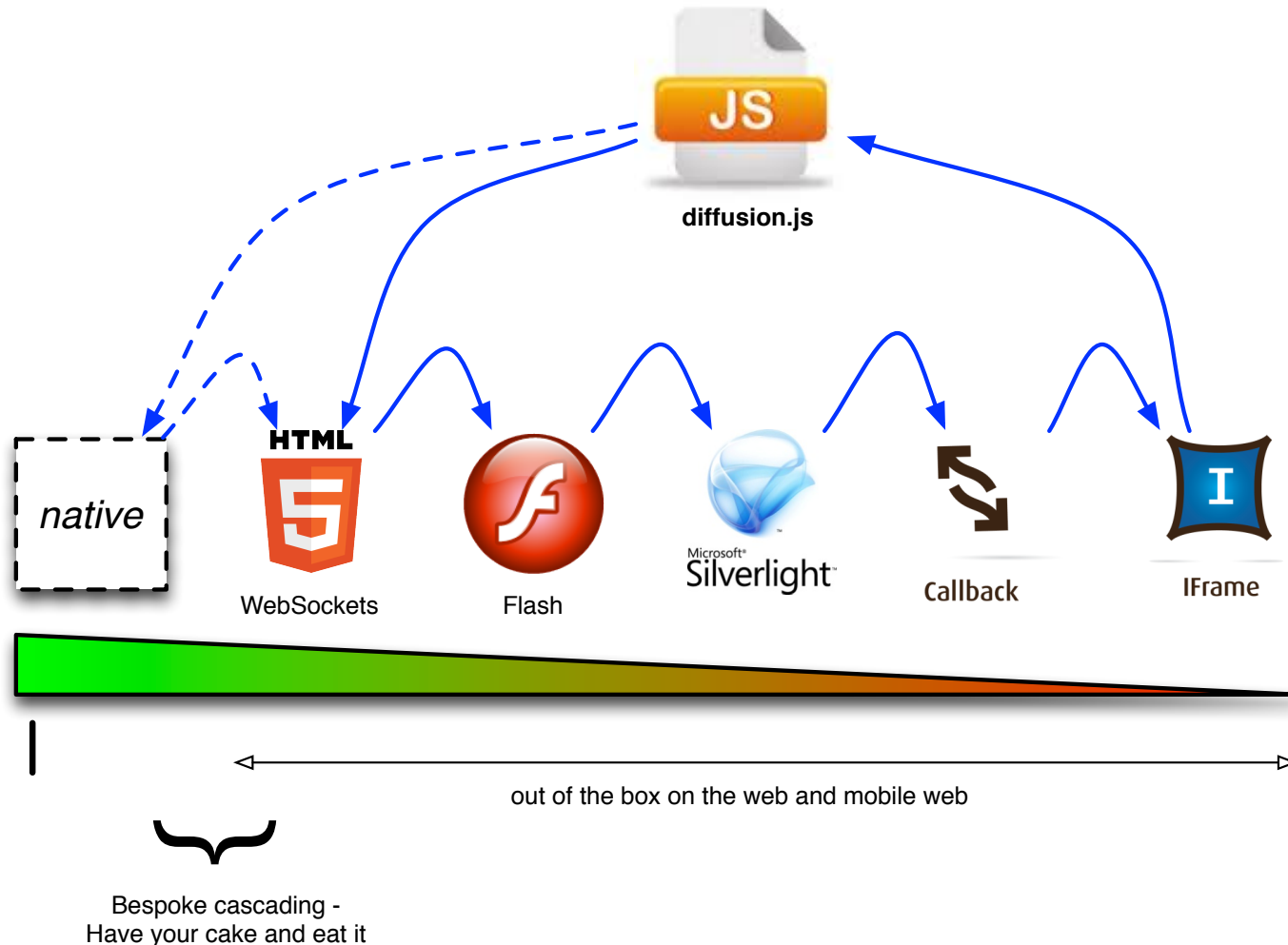| 8 | 1 | 5 |
| = | | = |
| C2 + C1 | | A2 + A1 |
| = | | = |
| 5 + 3 | | 4 + 1 |

Current
+
Consistent

## Replace

- Replace/Overwrite 'A1' with 'A2'
- 'Current data right now'
- Fast

## Merge

- Merge A1, A2 under some operation.
  - Operation is pluggable
  - Tunable consistency
  - Performance f(operation)

Data On Demand

PUSH
TECHNOLOGY

# Client? Connection & Transport Cascading



**diffusion.js**

native

**HTML** WebSockets

Flash

Microsoft® Silverlight™

Callback

**I** IFrame

out of the box on the web and mobile web

Bespoke cascading -
Have your cake and eat it

Data On Demand

**PUSH** TECHNOLOGY
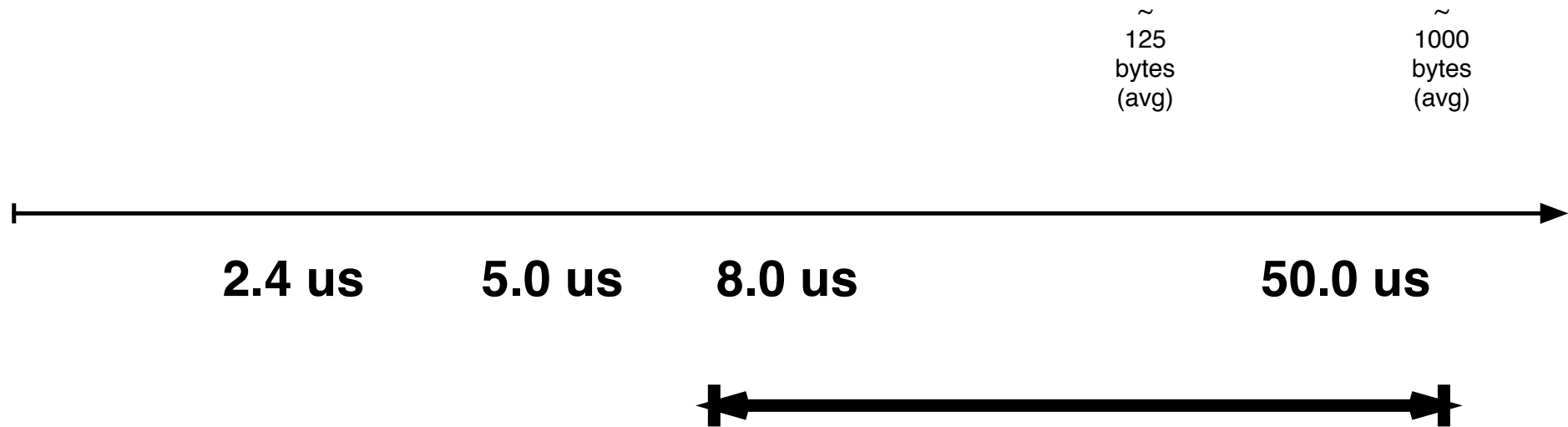
# An extreme example.



Box A

Box B

**Throughput (Worst case)**

- Ramp clients continuously
- 100 messages per second per client
- Payload: 125 .. 2000 bytes
- Message style vs with conflation

**Latency (Best case)**

- Really simple. 1 client
- Ping – Pong. Measure RTT time.
- Payload: 125 .. 2000 bytes

PUSH
TECHNOLOGY

# Latency in perspective

~
125
bytes
(avg)

~
1000
bytes
(avg)

**2.4 us**     **5.0 us**     **8.0 us**          **50.0 us**

- 2.4 us. A tuned benchmark in C with low latency 10Ge NIC, with kernel bypass, with FPGA acceleration
- 5.0 us. A basic java benchmark – as good as it gets in java
- Diffusion is measurably 'moving to the left' release on release
- We've been actively tracking and continuously improving since 4.0

# Throughput. Server view. Non-conflated data



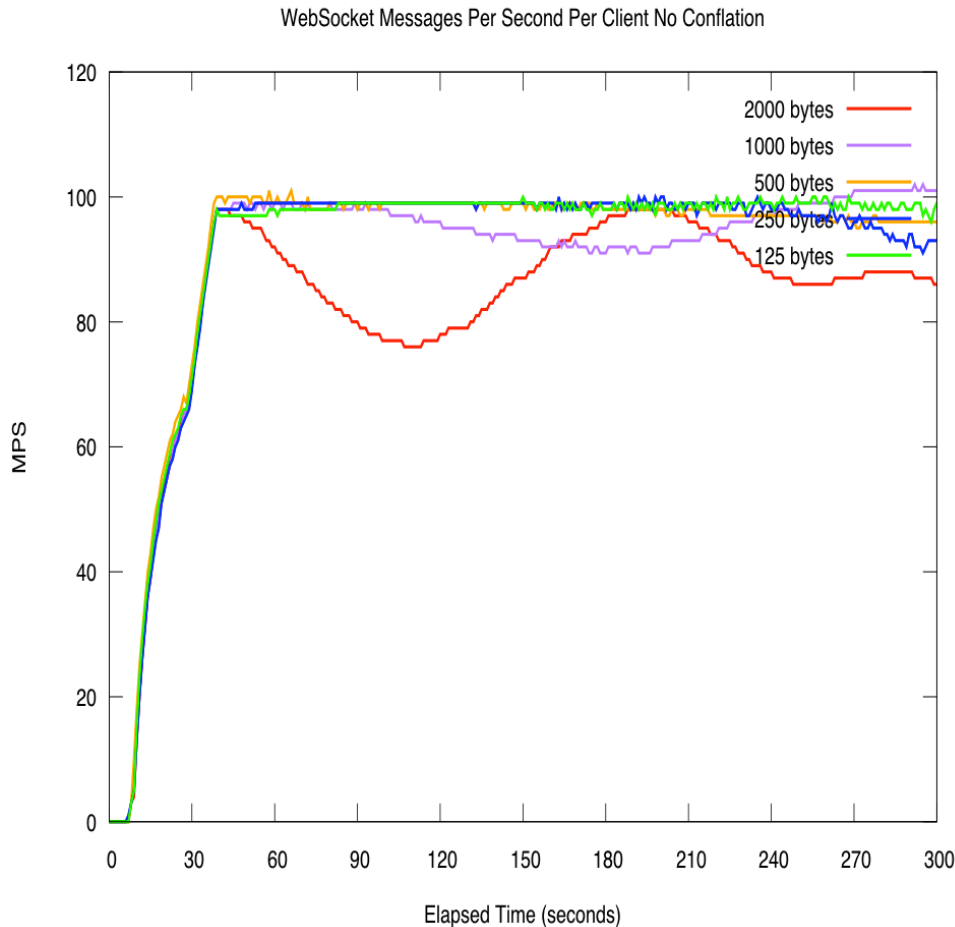WebSocket Bandwidth No Conflation

- Cold start server

- Ramp 750 clients 'simultaneously' at 5 second intervals

- 5 minute benchmark duration

- Clients onboarded linearly until IO (here) or compute saturation occurs.

- What the 'server' sees

Data On Demand

PUSH
TECHNOLOGY

# Throughput. Client view. Non-conflated data.



WebSocket Messages Per Second Per Client No Conflation

- What the 'client' sees

- At and beyond saturation of some resource?

- Things break!

- New connections fail. Good.
- Long established connections ok.
- Recent connections timeout and client connections are dropped. Good.

- Diffusion handles smaller message sizes more gracefully

- Back-pressure 'waveform' can be tuned out. Or, you could use structural conflation!

PUSH
TECHNOLOGY

# Throughput. Server view. Replace-conflated
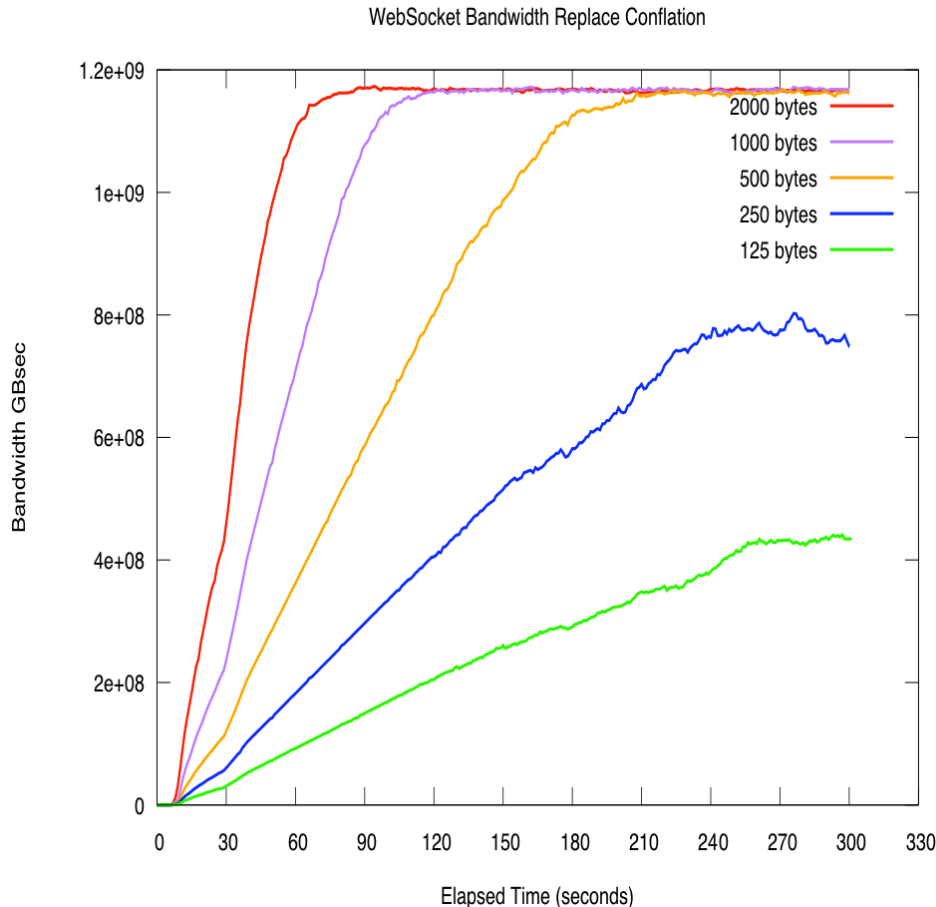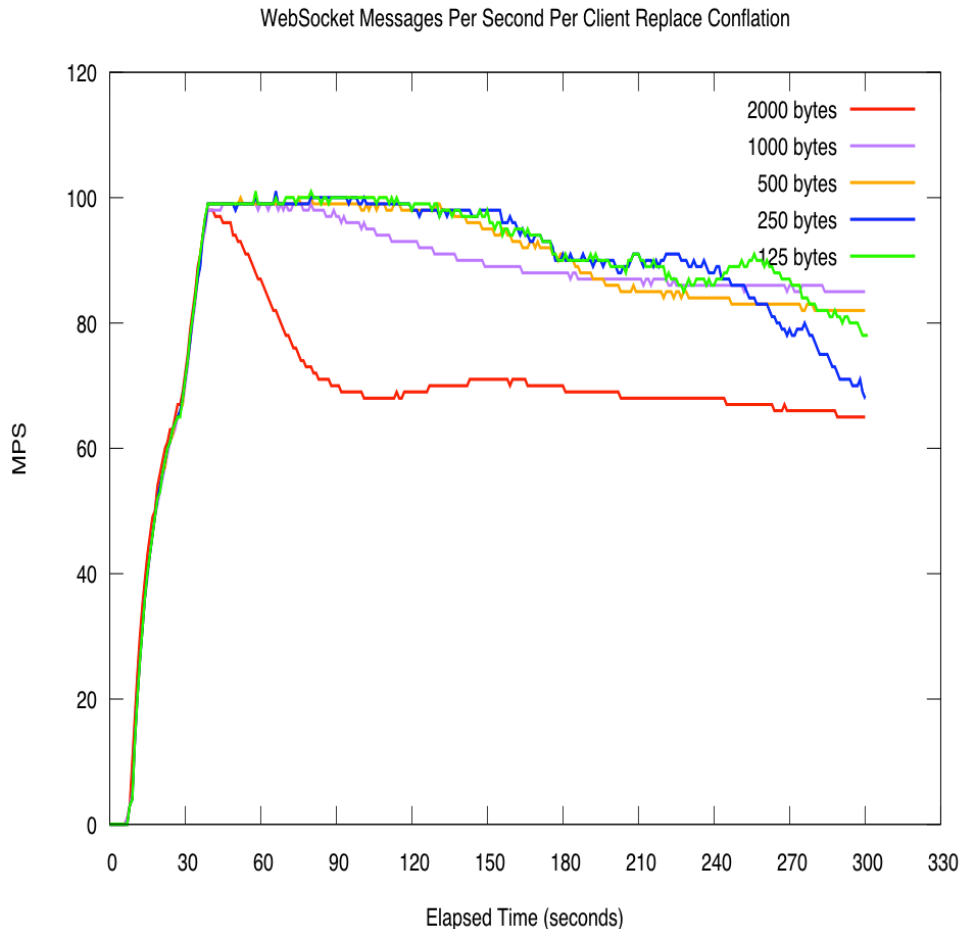


WebSocket Bandwidth Replace Conflation

- Cold start server

- Ramp 750 clients 'simultaneously' at 5 second intervals

- 5 minute benchmark duration

- Clients onboarded linearly until IO (here) or compute saturation occurs.

- Takes longer to saturate than non-conflated case

- Handles more concurrent connections

- Again, 'server' view

Data On Demand

PUSH
TECHNOLOGY

# Throughput. Client view. Replace-Conflated



WebSocket Messages Per Second Per Client Replace Conflation

- What the 'client' sees

- Once saturation occurs Diffusion adapts actively by degrading messages per second per client

- This is good. Soak up the peaks through fairer distribution of data.

- Handles spikes in number of connections, volume of data or saturation of other resources using a single load adaptive mechanism

Data On Demand

PUSH
TECHNOLOGY

# Can failure be fun?

- Where do **HA algorithms fail** humans?

- Can *consistent experiences* be delivered elapsing short, medium and longer term failures?

- Does an **M2H eye for an M2M** guy even work?

Data On Demand

PUSH ℗ ™
TECHNOLOGY

# Who knows. Let's figure it out! ☺

# My 'number'

# #6752

Data On Demand

**PUSH** TECHNOLOGY

- Thank you for listening to, having me

- Le twitter: @darachennis

- Work. http://www.pushtechnology.com/

- Fun. http://github.com/darach/

- Grab me any time on Weds/Thurs.

- Like Erlang? Big Data?
London User Group meetup.
Thurs March 7th, QE2 18:30 – 20:00.
(Google for registration details).

# Questions?

Darach@PushTechnology.com