# Large-Scale Systems
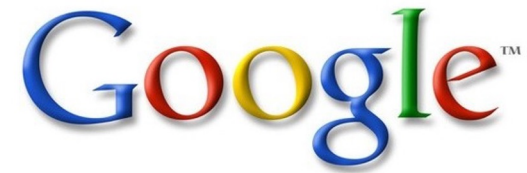## The Unreasonable Effectiveness of Simplicity

Randy Shoup
@randyshoup

# Background

# Goals

- From a systems perspective, what characterizes a scalable, well-engineered system?

- What can (application) systems designers learn from this community?

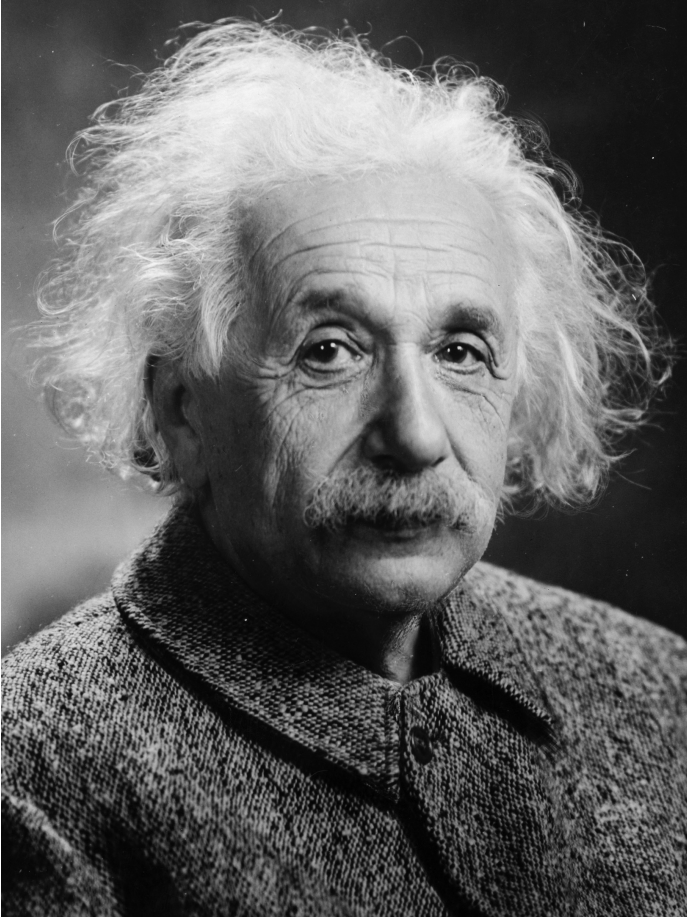- Where are the biggest opportunities to improve (application) systems?

@randyshoup

# Evolving Systems

- eBay
  - 5$^{th}$ generation today
  - Monolithic Perl → Monolithic C++ → Java → microservices

- Twitter
  - 3$^{rd}$ generation today
  - Monolithic Rails → JS / Rails / Scala → microservices

- Amazon
  - Nth generation today
  - Monolithic Perl / C → C++ / Java services → microservices

No one starts with microservices

…

<u>Past a certain scale</u>, everyone ends up with microservices

"Make everything as simple as possible, but not simpler."

# Large-Scale Systems

- Simple Components
- Simple Interactions
- Simple Changes

# Large-Scale Systems

- Simple Components
- Simple Interactions
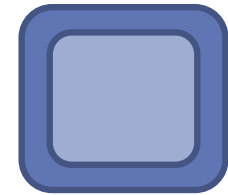- Simple Changes

# Modular Services

- Service boundaries match the problem domain

- Service boundaries encapsulate business logic and data
  - o All interactions through published service interface
  - o Interface hides internal implementation details
  - o No back doors

- Service boundaries encapsulate architectural -ilities
  - o Fault isolation
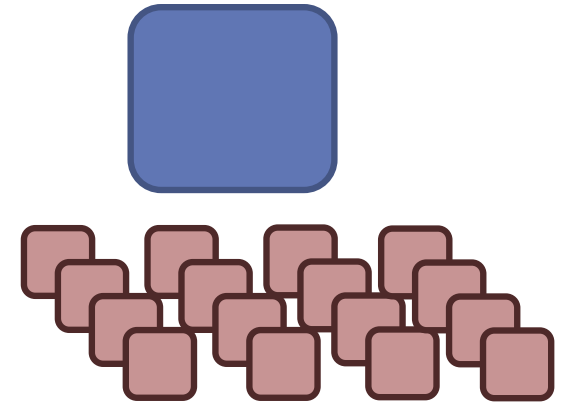  - o Performance optimization
  - o Security boundary

@randyshoup

# Orthogonal Domain Logic

- ## Stateless domain logic
  - o Ideally stateless pure function
  - o Matches domain problem as directly as possible
  - o Deterministic and testable in isolation
  - o Robust to change over time

- ## "Straight-line processing"
  - o Straightforward, synchronous, minimal branching

- ## Separate domain logic from I/O
  - o Hexagonal architecture, Ports and Adapters
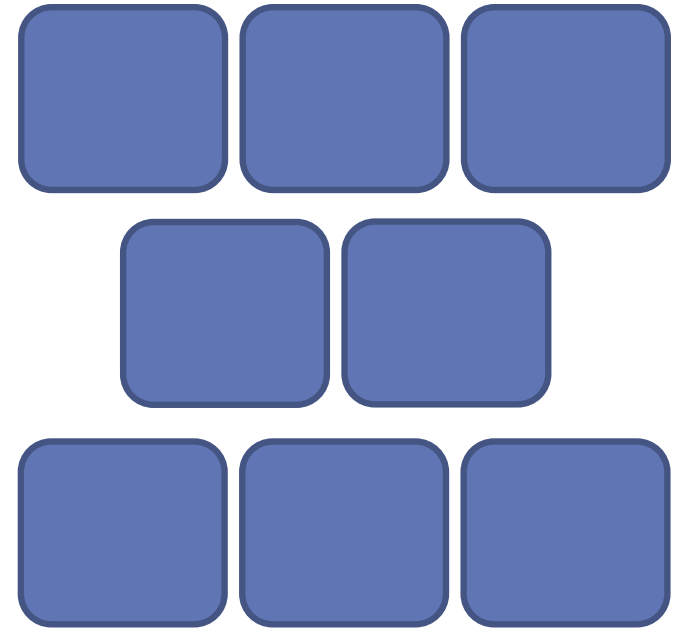  - o Functional core, imperative shell

@randyshoup

# Sharding

- ## Shards partition the service's "data space"
  - o Units for distribution, replication, processing, storage
  - o Hidden as internal implementation detail

- ## Shards encapsulate architectural -ilities
  - o Resource isolation
  - o Fault isolation
  - o Availability
  - o Performance

- ## Shards are autoscaled
  - o Divide or scale out as processing or data needs increase
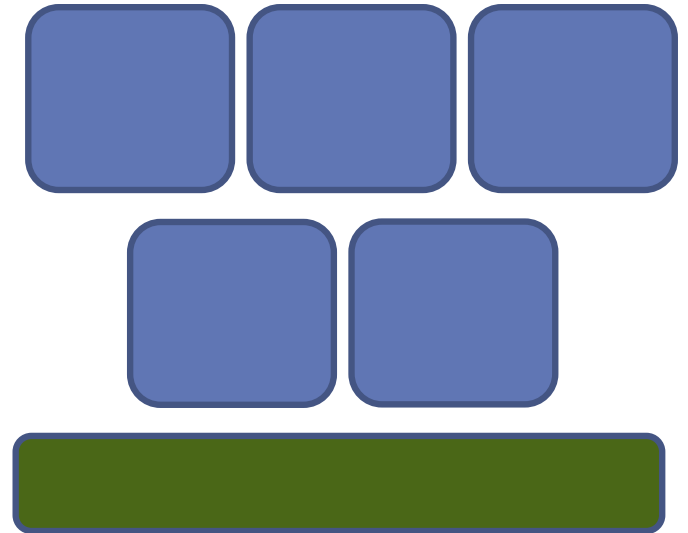  - o E.g., DynamoDB partitions, Aurora segments, Bigtable tablets

@randyshoup

# Service Layering

- **Common services provide and abstract widely-used capabilities**

- **Service ecosystem**
  - Services call others, which call others, etc.
  - Graph, not a strict layering

- **Services grow and evolve over time**
  - Factor out common libraries and services as needed
  - Teams and services split like "cellular mitosis"

@randyshoup

# Common Platform

- "Paved Road"
  - Shared infrastructure
  - Standard frameworks
  - Developer experience
  - E.g., Netflix, Google

- Separation of Concerns
  - Reduce cognitive load on stream-aligned teams
  - Bound decisions through enabling constraints

Large-scale organizations often invest more than 50% of engineering effort in platform capabilities

# Google Service Layering (2013)

- Cloud Datastore:  NoSQL service
  - Strong transactional consistency
  - SQL-like rich query capabilities
- Megastore:  geo-scale structured database
  - Multi-row transactions
  - Synchronous cross-datacenter replication
- Bigtable:  cluster-level structured storage
  - (row, column, timestamp) -> cell contents
- Colossus:  distributed file system
  - Block distribution and replication
- Borg:  cluster management infrastructure
  - Task scheduling, machine assignment

Cloud Datastore
↓
Megastore
↓
Bigtable
↓
Colossus
↓
Cluster manager

# Large-Scale Systems

- Simple Components

- **Simple Interactions**

- Simple Changes

# Reactive Manifesto

## The Reactive Manifesto

*Published on September 16 2014. (v2.0)*

Organisations working in disparate domains are independently discovering patterns for building software that look the same. These systems are more robust, more resilient, more flexible and better positioned to meet modern demands.

These changes are happening because application requirements have changed dramatically in recent years. Only a few years ago a large application had tens of servers, seconds of response time, hours of offline maintenance and gigabytes of data. Today applications are deployed on everything from mobile devices to cloud-based clusters running thousands of multi-core processors. Users expect millisecond response times and 100% uptime. Data is measured in Petabytes. Today's demands are simply not met by yesterday's software architectures.

We believe that a coherent approach to systems architecture is needed, and we believe that all necessary aspects are already recognised individually: we want systems that are Responsive, Resilient, Elastic and Message Driven. We call these Reactive Systems.

# Event-Driven

- Communicate state changes as stream of events
  - Statement that some interesting thing occurred
  - Ideally represents a semantic domain event

- Decouples domains and teams
  - Abstracted through a well-defined interface
  - Asynchronous from one another

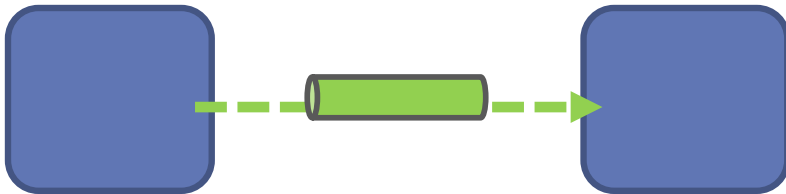- Simplifies component implementation

@randyshoup

# Immutable Log

- Store state as immutable log of events
  - Event Sourcing

- Often matches domain
  - E.g., Stitch Fix package processing / delivery state

- Log encapsulates architectural –ilities
  - Durable
  - Traceable and auditable
  - Replayable
  - Explicit and comprehensible

- Compact snapshots for efficiency

@randyshoup

# Embrace Asynchrony

- Decouples operations in time
  - Decoupled availability
  - Independent scalability
  - Allows more complex processing, more processing in parallel
  - Safer to make independent changes

- Simplifies component implementation

@randyshoup

# Embrace Asynchrony

- Invert from synchronous call graph to async dataflow
  - Exploit asymmetry between writes and reads
  - Can be orders of magnitude less resource intensive

@randyshoup

# Amazon Aurora



- ## Asynchronous redo log writes
  - Sent asynchronously to Aurora storage nodes
  - Acknowledged asynchronously to database instance
  - No distributed consensus round
  - Idempotent, immutable, monotonic

- ## Quorum acknowledgement
  - Log progresses forward once quorum of nodes acknowledges

- ## Reestablish consistency on crash recovery

@randyshoup

Verbitski, et al, 2018, Amazon Aurora: On Avoiding Distributed Consensus, SIGMOD '18.

# Netflix Viewing History



Playback API → Durable queues → Request Processor

- Store and process member's playback data
  - 1M requests per second
  - Used for viewing history, personalization, recommendations, analytics, etc.

- Original synchronous architecture
  - Synchronously write to persistent storage and lookup cache
  - Availability and data loss from backpressure at high load

- Asynchronous rearchitecture
  - Write to durable queue
  - Async pipeline to enrich, process, store, serve
  - Materialize views to serve reads

# Walmart Item Availability

- Is this item available to ship to this customer?
  - Customer SLO 99.98% uptime in 300ms

- Complex logic involving many teams and domains
  - Inventory, reservations, backorders, eligibility, sales caps, etc.

- Original synchronous architecture
  - Graph of 23 nested synchronous service calls in hot path
  - Any component failure invalidates results
  - Service SLOs 99.999% uptime with 50ms marginal latency
  - Extremely expensive to build and operate

Scott Havens, 2019, Fabulous Fortunes, Fewer Failures, and Faster Fixes from Functional Fundamentals, DOES 2019.

# Walmart Item Availability

# Walmart Item Availability

- Invert each service to use async events
  - Event-driven "dataflow"
  - Idempotent processing
  - Event-sourced immutable log
  - Materialized view of data from upstream dependencies

- Asynchronous rearchitecture
  - 2 services in synchronous hot path
  - Async service SLOs 99.9% uptime with latency in seconds or minutes
  - More resilient to delays and outages
  - Orders of magnitude simpler to build and operate

# Walmart Item Availability



Scott Havens, 2019, Fabulous Fortunes, Fewer Failures, and Faster Fixes from Functional Fundamentals, DOES 2019.

# Large-Scale Systems

- Simple Components
- Simple Interactions
- Simple Changes

# Incremental Change

- Decompose every large change into small incremental steps

- Each step maintains backward / forward compatibility of data and interfaces

- Multiple service versions commonly coexist
  - Every change is a rolling upgrade
  - *Transitional states are normal, not exceptional*

# Continuous Testing

- ## Tests help us go faster
  - Tests are "solid ground"
  - Tests are the safety net

- ## Tests make better code
  - Confidence to break things
  - Courage to refactor mercilessly

- ## Tests make better systems
  - Catch bugs earlier, fail faster

# Continuous Testing

- Tests make better designs
  - Modularity
  - Separation of Concerns
  - Encapsulation



@randyshoup

"There's a deep synergy between testability and good design. All of the pain that we feel when writing unit tests points at underlying design problems."

-- Michael Feathers

# Continuous Delivery

- Deploy services multiple times per day
    - Robust build, test, deploy pipeline
    - Canary deployment
    - Feature flags
    - Dark launches
    - SLO monitoring
    - Synthetic monitoring

- More solid systems
    - Release smaller, simpler units of work
    - Smaller changes to roll back or roll forward
    - Faster to repair, easier to understand, simpler to diagnose
    - Increase rate of change and reduce risk of change

@randyshoup

In the limit, production monitoring and software testing become the same thing

@randyshoup

# Software Delivery

- State of DevOps Surveys
  - 8 yearly surveys from 2014-2021
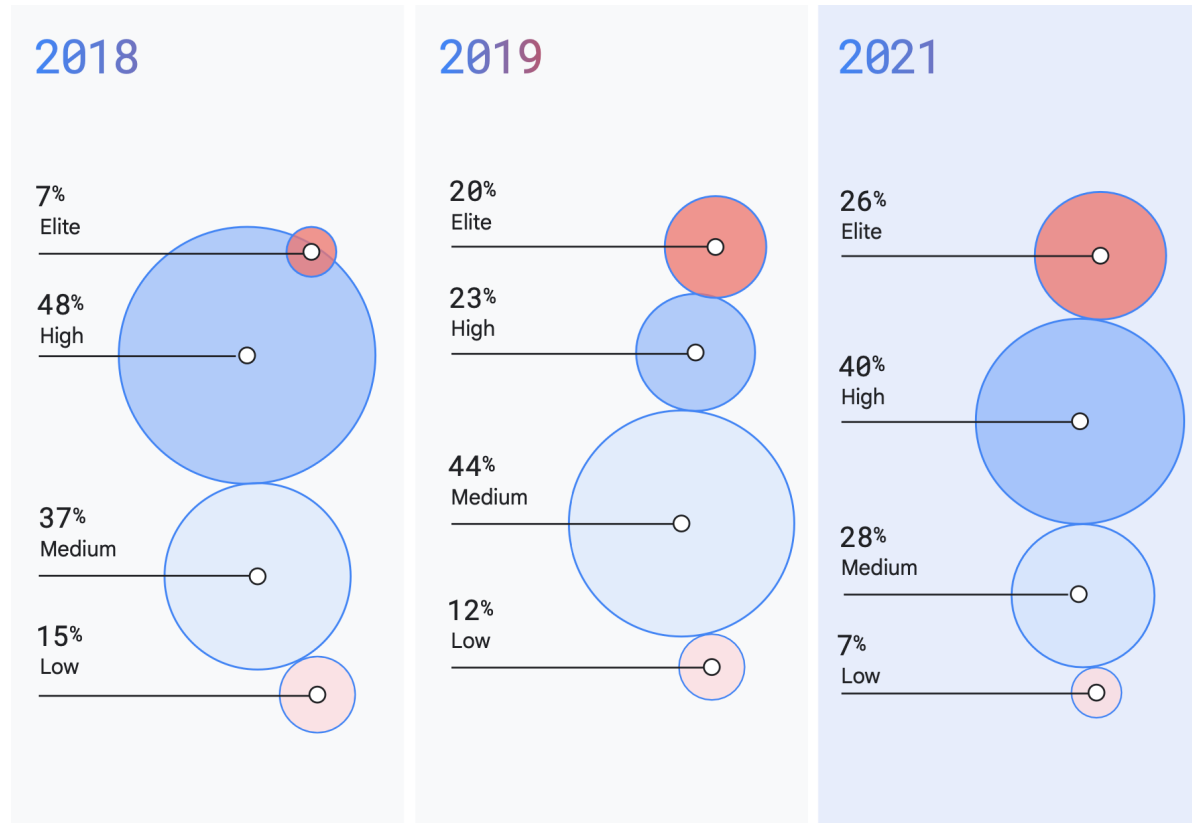  - 31,000 survey responses
  - Rigorous scientific methodology

- Summarized in *Accelerate*

THE SCIENCE OF DEVOPS

## ACCELERATE

Building and Scaling High Performing Technology Organizations

Nicole Forsgren, PhD
Jez Humble *and* Gene Kim

@randyshoup

# Software Delivery



**2018**

7% Elite

48% High

37% Medium

15% Low

**2019**

20% Elite

23% High

44% Medium

12% Low

**2021**

26% Elite

40% High

28% Medium

7% Low

State of DevOps Report, 2021

# Continuous Delivery

- Cross-company *Velocity Initiative* to improve software delivery
  - Think Big, Start Small, Learn Fast
  - Iteratively identify and remove bottlenecks for teams
  - "What would it take to deploy your application every day?"

- Doubled engineering productivity
  - 5x faster deployment frequency
  - 5x faster lead time
  - 3x lower change failure rate
  - 3x lower mean-time-to-restore

- Prerequisite for large-scale architecture changes

# Large-Scale Systems

- Simple Components

- Simple Interactions

- Simple Changes

# Thank you!

@randyshoup

linkedin.com/in/randyshoup

medium.com/@randyshoup