

We've long relied on our canine friends to provide boozy sustenance in our times of greatest need, ever since 16th-century St. Bernards revived snowbound Alpine travelers with mini-kegs of brandy. Actually, that's [almost entirely made up](#). But hey, this dog can fetch beer from the fridge on command!

Elements Network Sources Timeline Profiles Resources Audits Console

Preserve log Disable cache

Name Path	Method	Status Text	Type	Initiator	Size Content	Time Latency	Timeline	3.3 min	5.0 min
pixel?google_nid=crimtan&google_push=AHNF13KjgCbRr69Yheh_PkPU... cm.doubleclick.net	GET	200 OK	image/png	http://i.ctnsnet.com/int... Redirect	506 B 170 B	115 ms 113 ms			
?h.key=8HB75-GTLRZ-D5A8S-LFUSZ-UUNZF&rt.si=5aefad33-feab-45... 36e4f0de.mpstat.us	GET	204	text/html	On...	233 B	118 ms			
pixel.gif?source=smarttag&fired=user_data_timeout&com... beacon.krxn.net									
pixel.gif?e=11&i=GAWKERV4&k=1%3D1x1%3Ahttp%3A%... v4.moatads.com									
track.gif?&objectType=permalinkTag&objectType=perm... kinja.com/api/analytics/stats									
t?pid=52e9531be79548ced6000008&title=This%20Beer... edge.simplereach.com									
t?pid=52e9531be79548ced6000008&title=This%20Beer... edge.simplereach.com									
t?pid=52e9531be79548ced6000008&title=This%20Beer... edge.simplereach.com									
t?pid=52e9531be79548ced6000008&title=This%20Beer-Fetching%20D... edge.simplereach.com	GET	200 OK	application/j...	reach.js:8 Script	264 B 104 B	125 ms 122 ms			
refresh?callback=jQuery211012063115951605141_1414617120253&... kinja.com/api/analytics/t	GET	200 OK	text/javascript	Main-en-US-4159edb7... Script	799 B 190 B	125 ms 123 ms			
t?pid=52e9531be79548ced6000008&title=This%20Beer-Fetching%20D... edge.simplereach.com	GET	200 OK	application/j...	reach.js:8 Script	263 B 103 B	125 ms 123 ms			
t?pid=52e9531be79548ced6000008&title=This%20Beer-Fetching%20D... edge.simplereach.com	GET	200 OK	application/j...	reach.js:8 Script	264 B 104 B	125 ms 123 ms			
t?pid=52e9531be79548ced6000008&title=This%20Beer-Fetching%20D... edge.simplereach.com	GET	200 OK	application/j...	reach.js:8 Script	264 B 104 B	125 ms 123 ms			
t?pid=52e9531be79548ced6000008&title=This%20Beer-Fetching%20D... edge.simplereach.com	GET	200 OK	application/j...	reach.js:8 Script	264 B 104 B	126 ms 124 ms			
t?pid=52e9531be79548ced6000008&title=This%20Beer-Fetching%20D... edge.simplereach.com	GET	200 OK	application/j...	reach.js:8 Script	264 B 104 B	126 ms 124 ms			

Behavioural Tracking *Architecture for data pipelines*

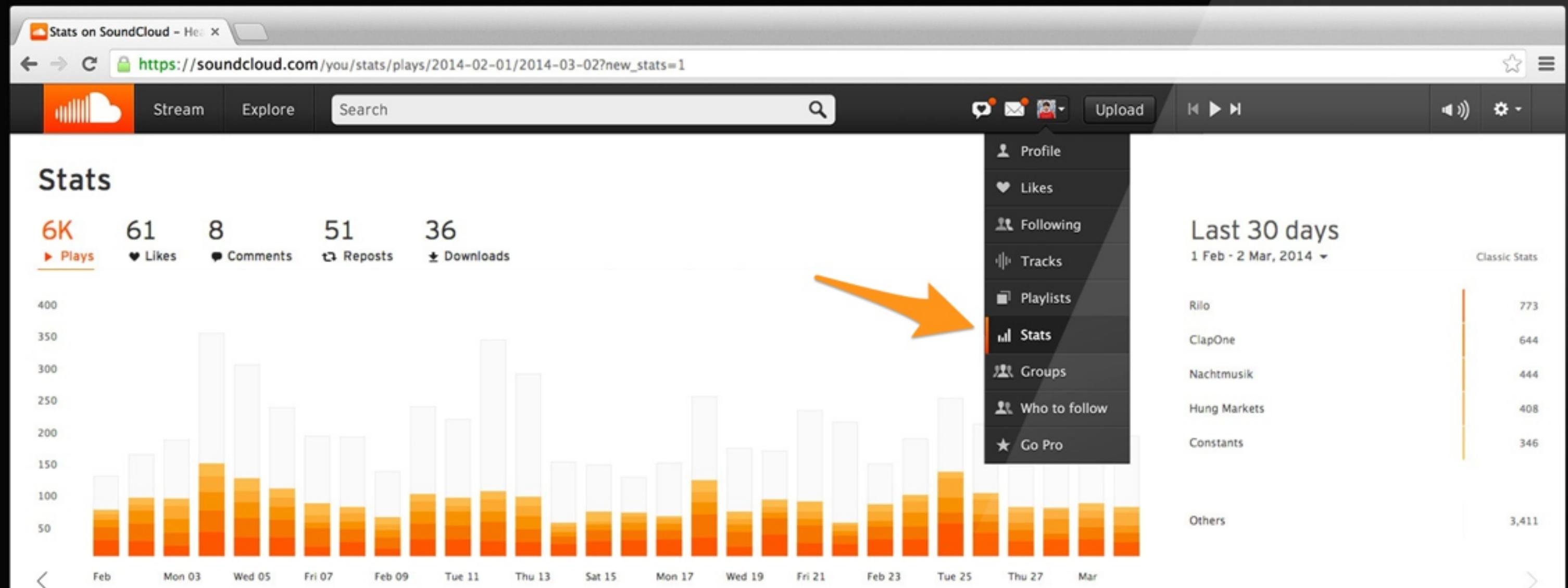
Sean Braithwaite

Outline

- SoundCloud
- Problem definition
- Ethics
- Technical requirements
- Architecture
- Alternative implementation
- Cloud Computing
- Open problems
- Q and A



SOUNDCLOUD



Most played tracks



Rilo

773

Top countries



United States
1,252

Who played the most



detailloop
122

Websites



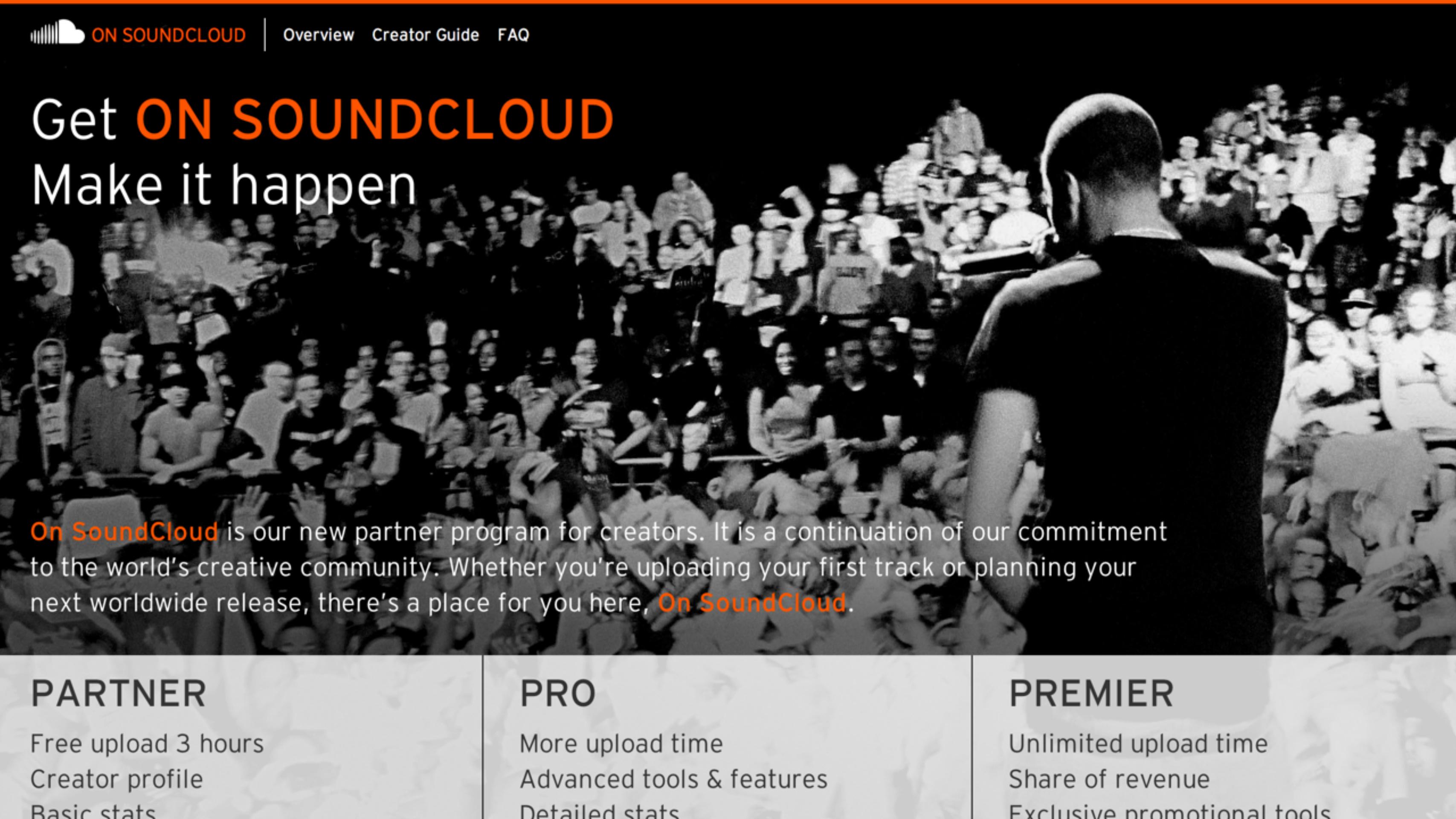
soundcloud.com/seams
1,696

Apps

2	ClapOne	644	2	Germany	1,238	2	gurū	94	2	soundcloud.com/seams/sets/quarters	206
3	Nachtmusik	444	3	United Kingdom	737	3	sanjay.fernandes	71	3	soundcloud.com/stream	134
4	Hung Markets	408	4	France	274	4	Luciano Castello	39	4	soundcloud.com/cecile9/likes	95

Get ON SOUNDCLOUD

Make it happen



On SoundCloud is our new partner program for creators. It is a continuation of our commitment to the world's creative community. Whether you're uploading your first track or planning your next worldwide release, there's a place for you here, On SoundCloud.

PARTNER

Free upload 3 hours
Creator profile
Basic stats

PRO

More upload time
Advanced tools & features
Detailed stats

PREMIER

Unlimited upload time
Share of revenue
Exclusive promotional tools

We've long relied on our canine friends to provide boozy sustenance in our times of greatest need, ever since 16th-century St. Bernards revived snowbound Alpine travelers with mini-kegs of brandy. Actually, that's **almost entirely made up**. But hey, this dog can fetch beer from the fridge on command!

Name Path	Method	Status Text	Type	Initiator	Size Content	Time Latency	Timeline	3.3 min	5.0 min
pixel?google_nid=crimtan&google_push=AHNF13KjgCbRr69Yheh_PkPU... cm.doubleclick.net	GET	200 OK	image/png	http://i.ctnsnet.com/int... Redirect	506 B 170 B	115 ms 113 ms			
?h.key=8HB75-GTLRZ-D5A8S-LFUSZ-UUNZF&rt.si=5aefad33-feab-45... 36e4f0de.mpstat.us	GET	204	application/	Other	233 B	118 ms			
pixel.gif?source=smarttag&fired=user_data_timeout&confid=JMWgL... beacon.krxn.net									
pixel.gif?e=11&i=GAWKERV4&k=1%3D1x1%3Ahttp%3A%2F%2Fgawke... v4.moatads.com									
track.gif?&objectType=permalinkTag&objectType=permalinkTag&obj... kinja.com/api/analytics/stats									
t?pid=52e9531be79548ced6000008&title=This%20Beer-Fetching%2... edge.simplereach.com									
t?pid=52e9531be79548ced6000008&title=This%20Beer-Fetching%2... edge.simplereach.com									
t?pid=52e9531be79548ced6000008&title=This%20Beer-Fetching%2... edge.simplereach.com									
t?pid=52e9531be79548ced6000008&title=This%20Beer-Fetching%20D... edge.simplereach.com	GET	200 OK	application/j...	reach.js:8 Script	264 B 104 B	125 ms 122 ms			
refresh?callback=jQuery211012063115951605141_1414617120253&... kinja.com/api/analytics/t	GET	200 OK	text/javascript	Main-en-US-4159edb7... Script	799 B 190 B	125 ms 123 ms			
t?pid=52e9531be79548ced6000008&title=This%20Beer-Fetching%20D... edge.simplereach.com	GET	200 OK	application/j...	reach.js:8 Script	263 B 103 B	125 ms 123 ms			
t?pid=52e9531be79548ced6000008&title=This%20Beer-Fetching%20D... edge.simplereach.com	GET	200 OK	application/j...	reach.js:8 Script	264 B 104 B	125 ms 123 ms			
t?pid=52e9531be79548ced6000008&title=This%20Beer-Fetching%20D... edge.simplereach.com	GET	200 OK	application/j...	reach.js:8 Script	264 B 104 B	125 ms 123 ms			
t?pid=52e9531be79548ced6000008&title=This%20Beer-Fetching%20D... edge.simplereach.com	GET	200 OK	application/j...	reach.js:8 Script	264 B 104 B	126 ms 124 ms			
t?pid=52e9531be79548ced6000008&title=This%20Beer-Fetching%20D... edge.simplereach.com	GET	200 OK	application/j...	reach.js:8 Script	264 B 104 B	126 ms 124 ms			

Motivation

- Understanding product/content
- Understanding users
- Data is a business

The Ethical Dimension

- It's surveillance
- Subject / observer asymmetry
- Evolving regulations

Technical Requirements

- Heterogeneous production
- Heterogeneous consumption
- Delivery Guarantees
- Insightful

Architecture Overview



Production

Collection

Transmission

Augmentation

Storage

Query

```
where = document.getElementsByTagName('script')[0];
where.parentNode.insertBefore(iframe, where);

try {
    doc = iframe.contentWindow.document;
} catch(e) {
    dom = document.domain;
    iframe.src="javascript:var d=document.open();d.domain='"+dom+"';void(0)";
    doc = iframe.contentWindow.document;
}
doc.open()._l = function() {
    var js = this.createElement("script");
    if(dom) this.domain = dom;
    js.id = "boomr-if-as";
    js.src = '//c.go-mpulse.net/boomerang/' +
    '8HB75-GTLRZ-D5A8S-LFUSZ-UUNZF';
    BOOMR_lstart=new Date().getTime();
    this.body.appendChild(js);
};
doc.write('<body onload="document.close();">');
doc.close();
})();
</script><script>window.kinja.s
type='text/javascript'
try {
    amznads.getAds('3076');
} catch(e) { console.log(e.message);
</script>      <script>
window.Krux || ((Krux = function () {
    Krux.q.push(arguments);
}).q = []);
(function () {
    function retrieve(n) {
        var m, k = 'kx' + n;
        if (window.localStorage) {
            return window.localStorage[k] || "";
        } else if (navigator.cookieEnabled) {
            m = document.cookie.match(k + '=([^\;]*)');
            return (m && unescape(m[1])) || "";
        } else {
            return '';
        }
    }
    Krux.user = retrieve('user');
    Krux.segments = retrieve('segs') && retrieve('segs').split(',') || [];
})();
</script><!-- Included CSS Files --><!--[if IE 9]><link rel="stylesheet" href="http://c.kinja-static.com/assets/style.css" type="text/css" /><![endif]--><!--[if !IE]>--><link rel="stylesheet" href="http://c.kinja-static.com/assets/style.css" type="text/css" />
```

Client side implementation

- Instrument application
- Schedule + Plan
- Handle Failure

Client side implementation

- Instrument application logic/state
 - Schedule + Plan for HTTP Calls
 - Handle Failure

Production

Collection

Transmission

Augmentation

Storage

Query

```
where = document.getElementsByTagName('script')[0];
where.parentNode.insertBefore(iframe, where);

try {
    doc = iframe.contentWindow.document;
} catch(e) {
    dom = document.domain;
    iframe.src="javascript:var d=document.open();d.domain='"+dom+"'";void(0)
    doc = iframe.contentWindow.document;
}
doc.open()._l = function() {
    var js = this.createElement("script");
    if(dom) this.domain = dom;
    js.id = "boomerang";
    js.src = '//c.go-mpulse.net/boomerang/' +
    '8HB75-GTLRZ-D5A8S-LFUSZ-UUNZF';
    BOOMR_lstart=new Date().getTime();
    this.body.appendChild(js);
};
doc.write('<body onload="document.close();">');
doc.close();
})();
</script><script>window.kinja.so
type='text/javascript'>
try {
    amznads.getAds('3076');
} catch(e) { console.log(e.message);
</script>      <script>
window.Krux || ((Krux = function () {
    Krux.q.push(arguments);
}).q = []);
(function () {
    function retrieve(n) {
        var m, k = 'kx' + n;
        if (window.localStorage) {
            return window.localStorage[k] || "";
        } else if (navigator.cookieEnabled) {
            m = document.cookie.match(k + '=([^\;]*)');
            return (m && unescape(m[1])) || "";
        } else {
            return '';
        }
    }
    Krux.user = retrieve('user');
    Krux.segments = retrieve('segs') && retrieve('segs').split(',') || [];
})();
</script><!-- Included CSS Files --><!--[if IE 9]><link rel="stylesheet" href="http
[endif]--><!--[if !IE]>--><link rel="stylesheet" href="http://c.kinja-static.com/assets/sty
```

HTTP Specifics

- Generate ASYNCHRONOUSLY
- POST is better than GET
- Cookie store usage

HTTP Specifics

- Generate ASYNC HTTP requests
 - POST is better than GET
 - Cookie store unique user ID

```

85
86     where = document.getElementsByTagName('script')[0];
87     where.parentNode.insertBefore(iframe, where);
88
89     try {
90         doc = iframe.contentWindow.document;
91     } catch(e) {
92         dom = document.domain;
93         iframe.src="javascript:var d=document.open();d.domain='"+dom+"'";void(0);";
94         doc = iframe.contentWindow.document;
95     }
96     doc.open()._l = function() {
97         var js = this.createElement("script");
98         if(dom) this.domain = dom;
99         js.id = "boomer-if-as";
100        js.src = '//c.go-mpulse.net/boomerang/' +
101          '8HB75-GTLRZ-D5A8S-LFUSZ-UUNZF';
102        BOOMR_lstart=new Date().getTime();
103        this.body.appendChild(js);
104    };
105    doc.write('<body onload="docu');
106    doc.close();
107  })();
108  </script><script>window.kinja.sco
109 type='text/javascript'>
110   try {
111     amznads.getAds('3076');
112   } catch(e) { console.log(e.message);
113   </script>      <script>
114 window.Krux || ((Krux = function () {
115   Krux.q.push(arguments);
116 }).q = []);
117 (function () {
118   function retrieve(n) {
119     var m, k = 'kx' + n;
120     if (window.localStorage) {
121       return window.localStorage[k] || "";
122     } else if (navigator.cookieEnabled) {
123       m = document.cookie.match(k + '=(.*?);');
124       return (m && unescape(m[1])) || "";
125     } else {
126       return '';
127     }
128   }
129   Krux.user = retrieve('user');
130   Krux.segments = retrieve('segs') && retrieve('segs').split(',') || [];
131 })();
132 </script><!-- Included CSS Files --><!--[if IE 9]><link rel="stylesheet" href="http://c.kinja-static.com/assets/stylesheets/tiger-split-7ca07e9d79a959e3ed448de1c35bd9cc.css"><![endif]--><!--[if !IE]>--><link rel="stylesheet" href="http://c.kinja-static.com/assets/stylesheets/tiger-eldeb8fb2alf708836b8a6aalcbfcf07.css"><!--<![endif]--><meta name="dynamic-

```

Platforms Matter

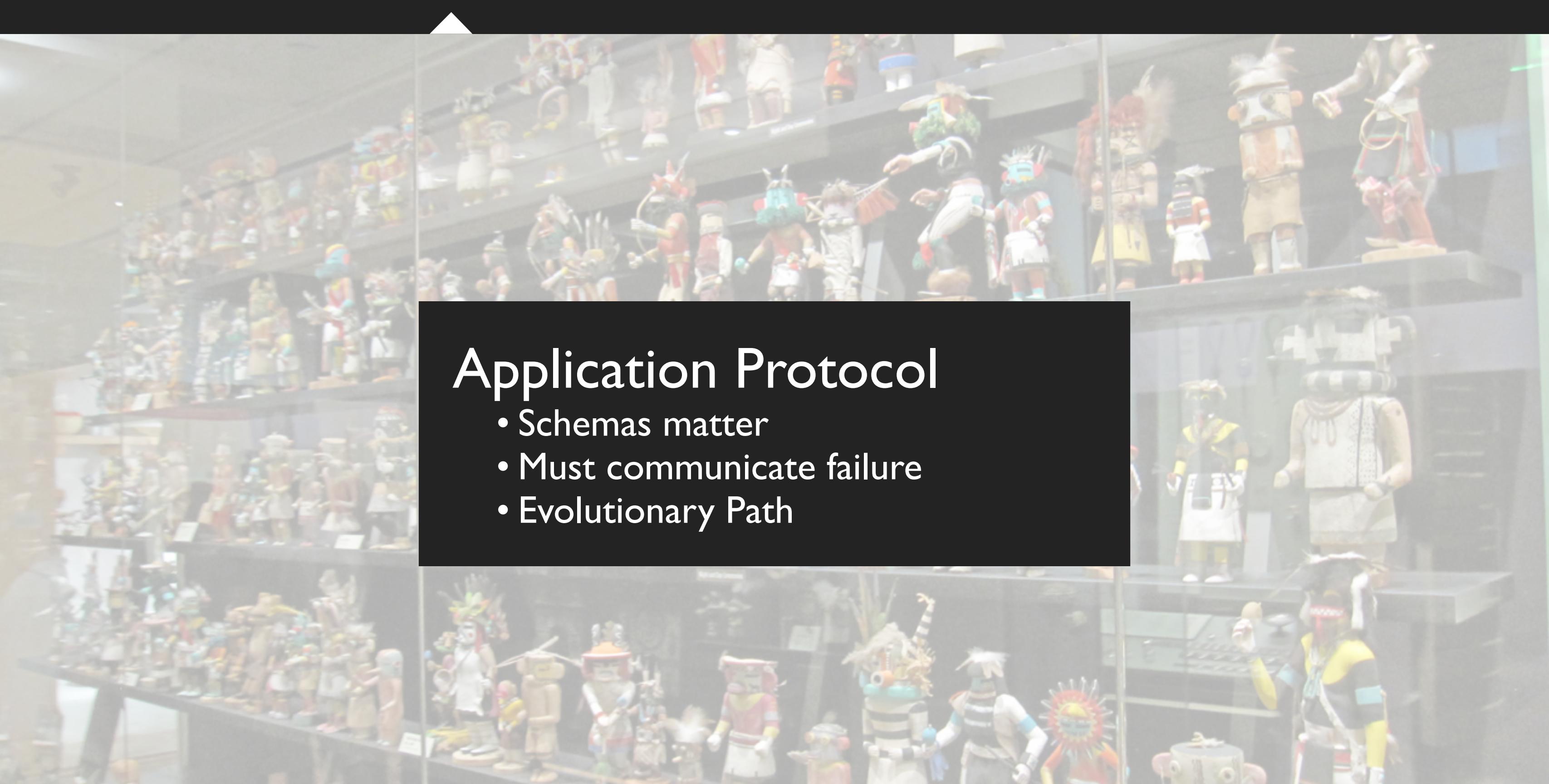
- Battery life on mobile is an issue
- Keep Alive, Scheduling, Batching
- Device ID and local storage

Collection

- Front Door of the pipeline
- Highly Available
- Consistent application protocol

Application Protocol

- Schemas matter
- Must communicate failure
- Evolutionary Path



HA HTTP services

- DNS offers high level routing
- Load Balancers
- Circuit Breakers

Production

Collection

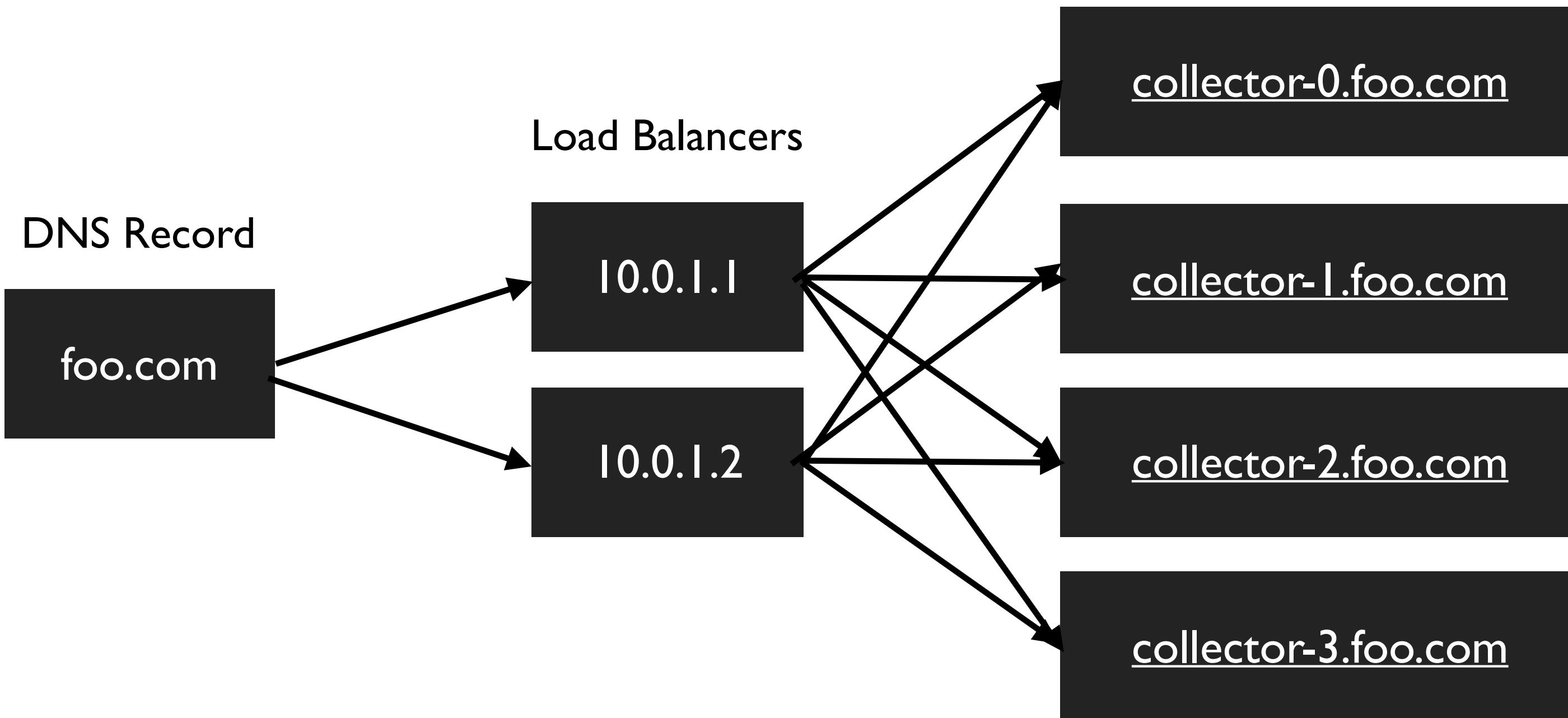
Transmission

Augmentation

Storage

Query

Application Servers



Transmission

- Delivery Guarantees
- Highly Available
- Known Failure Modes

Production

Collection

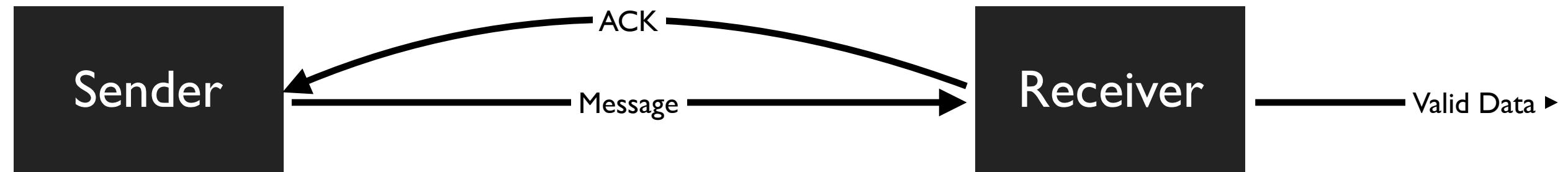
Transmission

Augmentation

Storage

Query

Ideal Case



Production

Collection

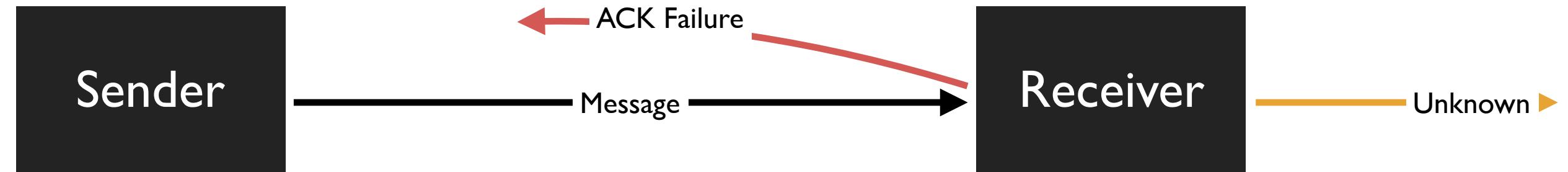
Transmission

Augmentation

Storage

Query

At most once



Production

Collection

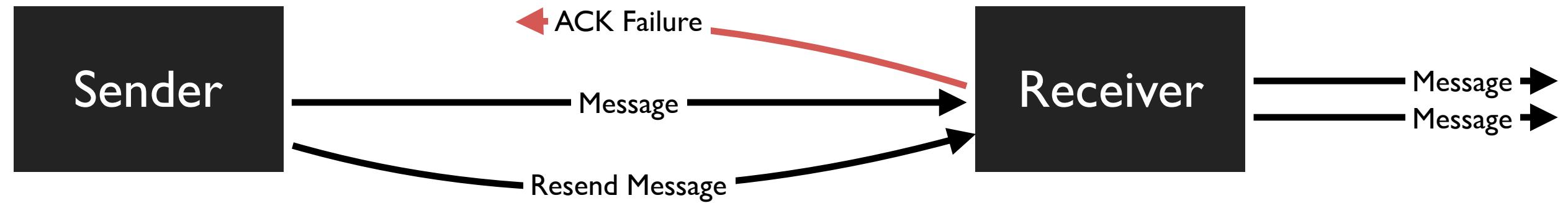
Transmission

Augmentation

Storage

Query

Duplicate Case: at least once



Production

Collection

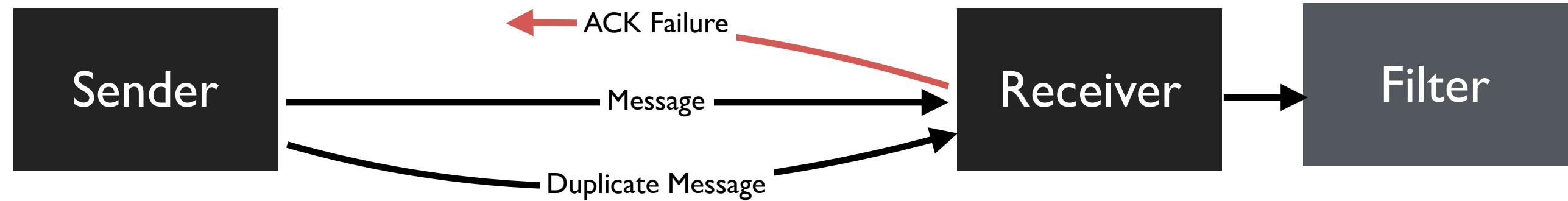
Transmission

Augmentation

Storage

Query

Exactly Once*: De-duplication



RabbitMQ

- Master/Slave
- Complex Topology
- At least once

Flume 0.94

- Master election
- Zookeeper
- At least once

Kafka

- Masterless
- Zookeeper
- At least once

Production

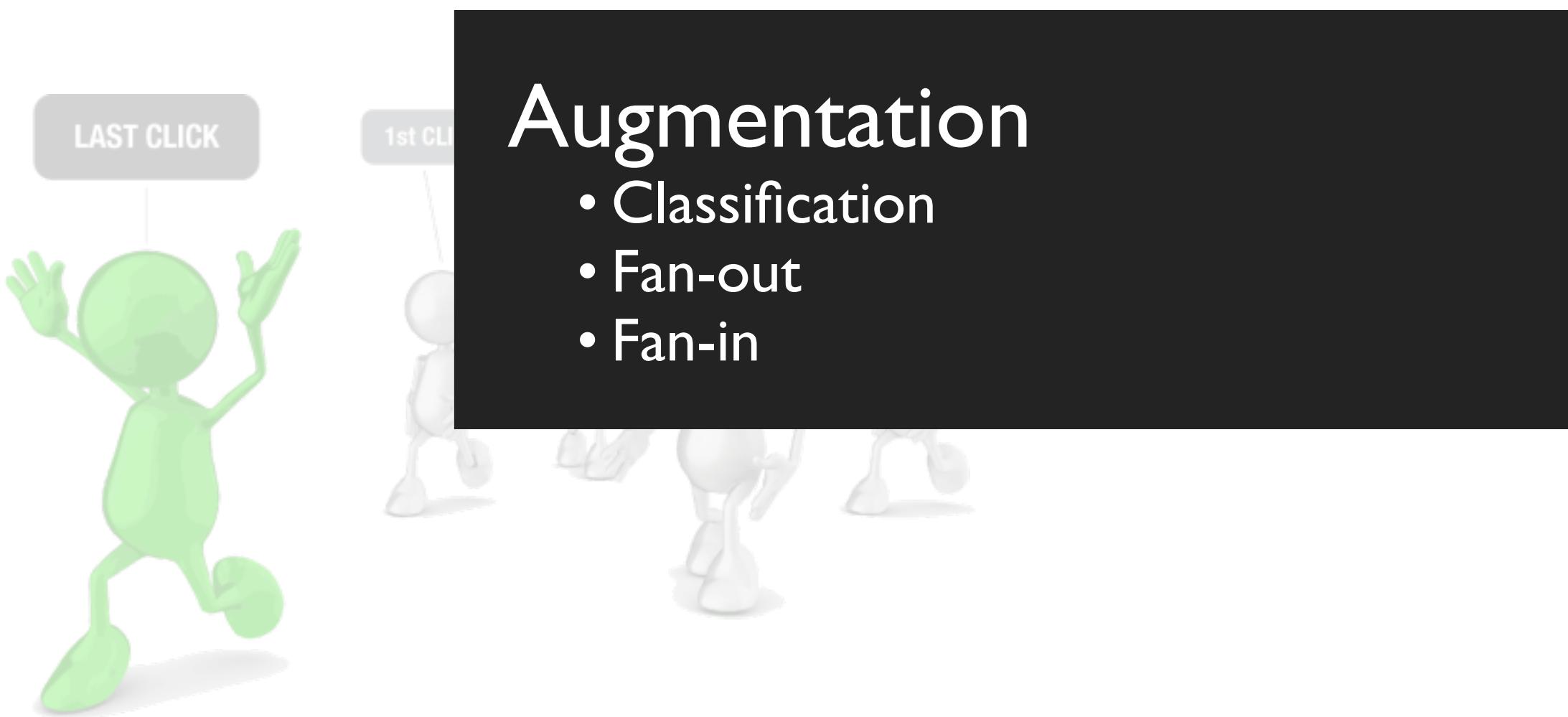
Collection

Transmission

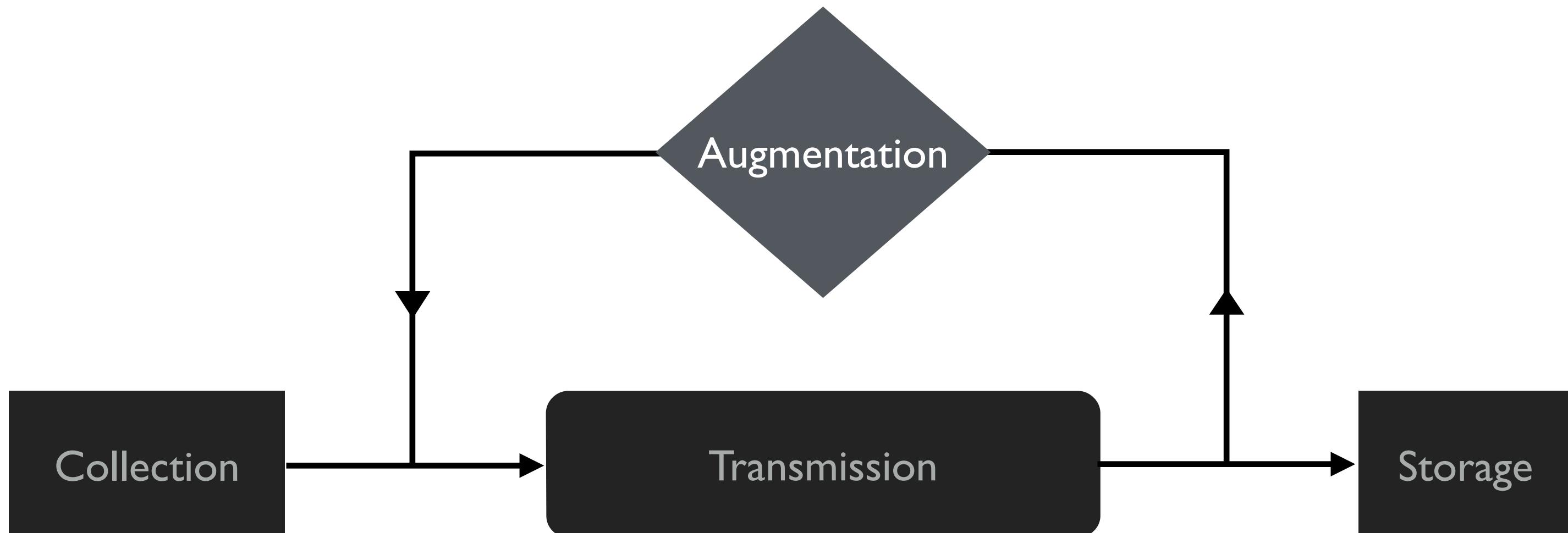
Augmentation

Storage

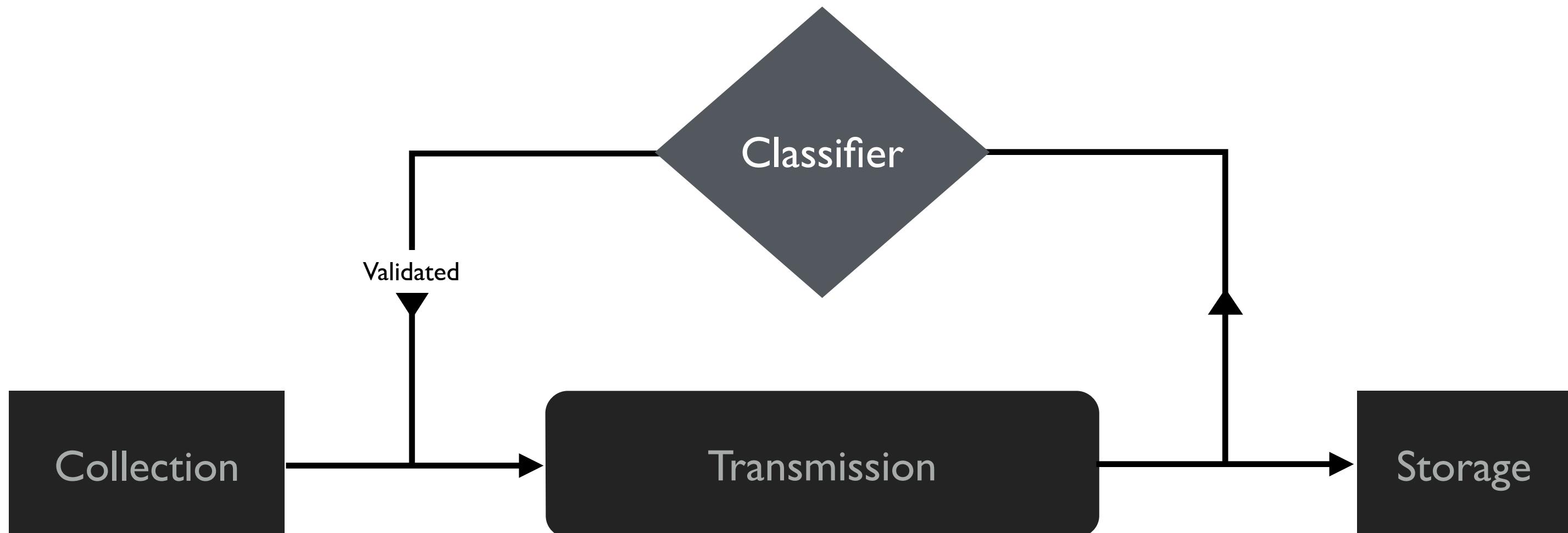
Query



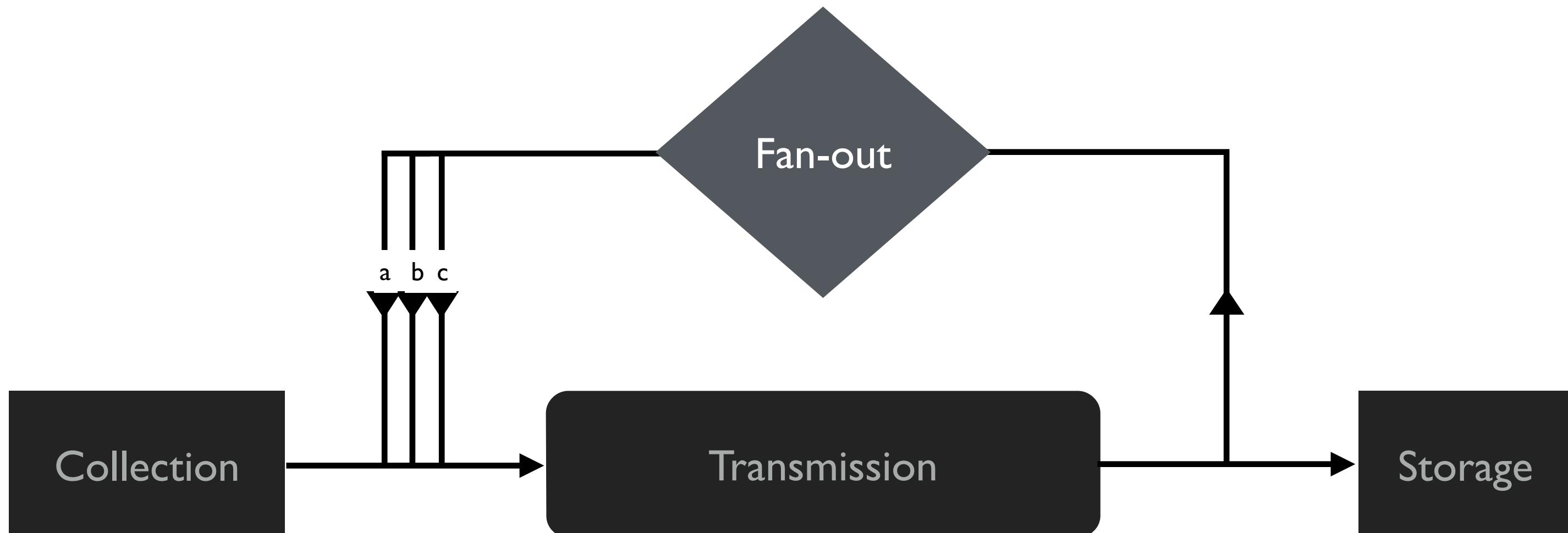
Production Collection Transmission Augmentation Storage Query



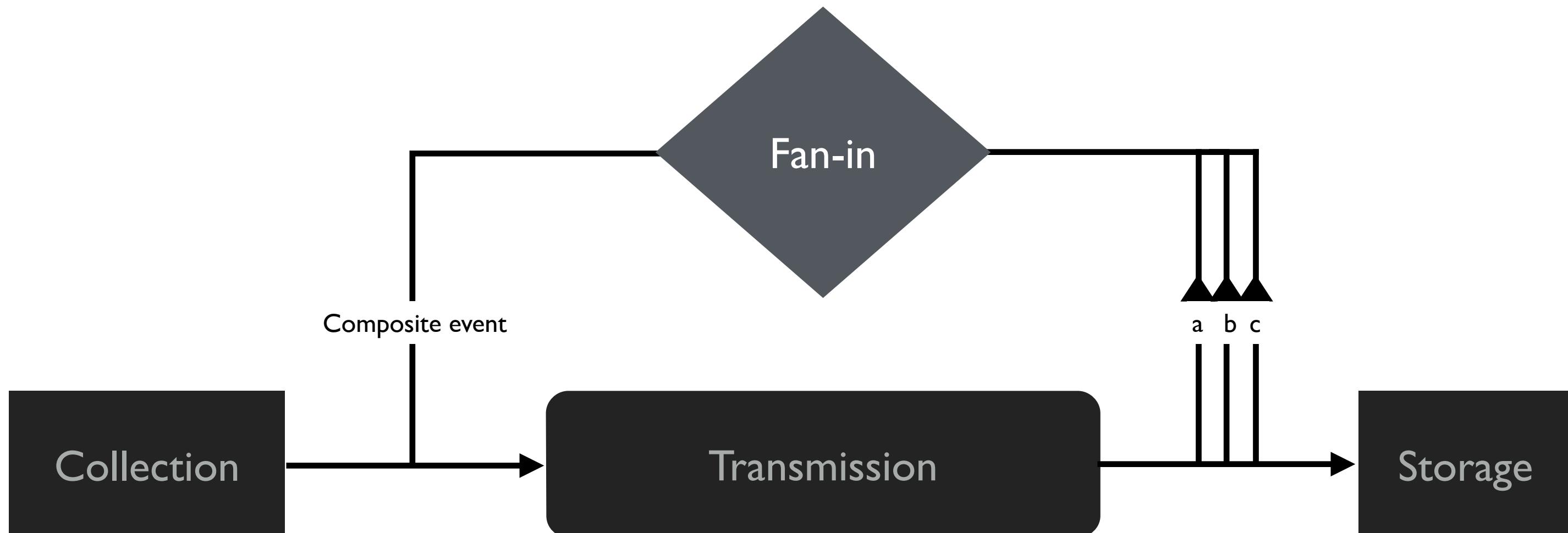
Production Collection Transmission Augmentation Storage Query



Production Collection Transmission Augmentation Storage Query



Production Collection Transmission Augmentation Storage Query



Production

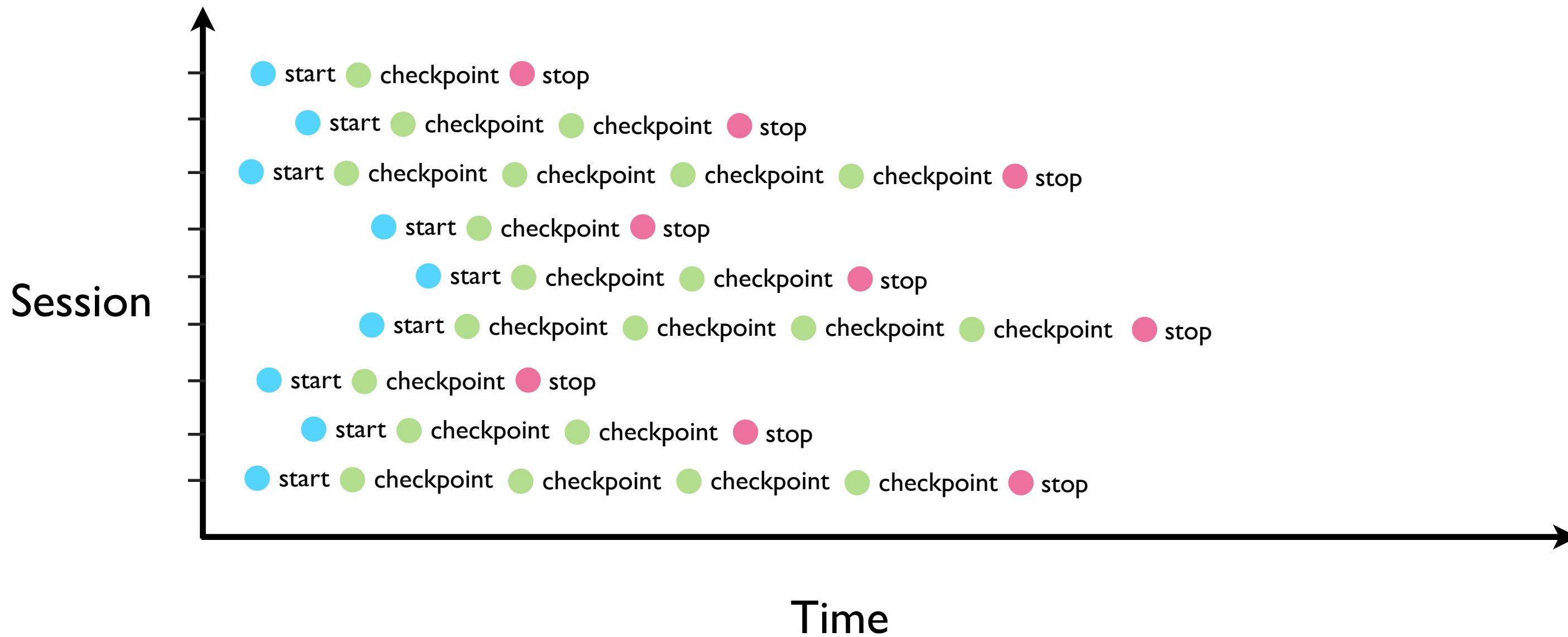
Collection

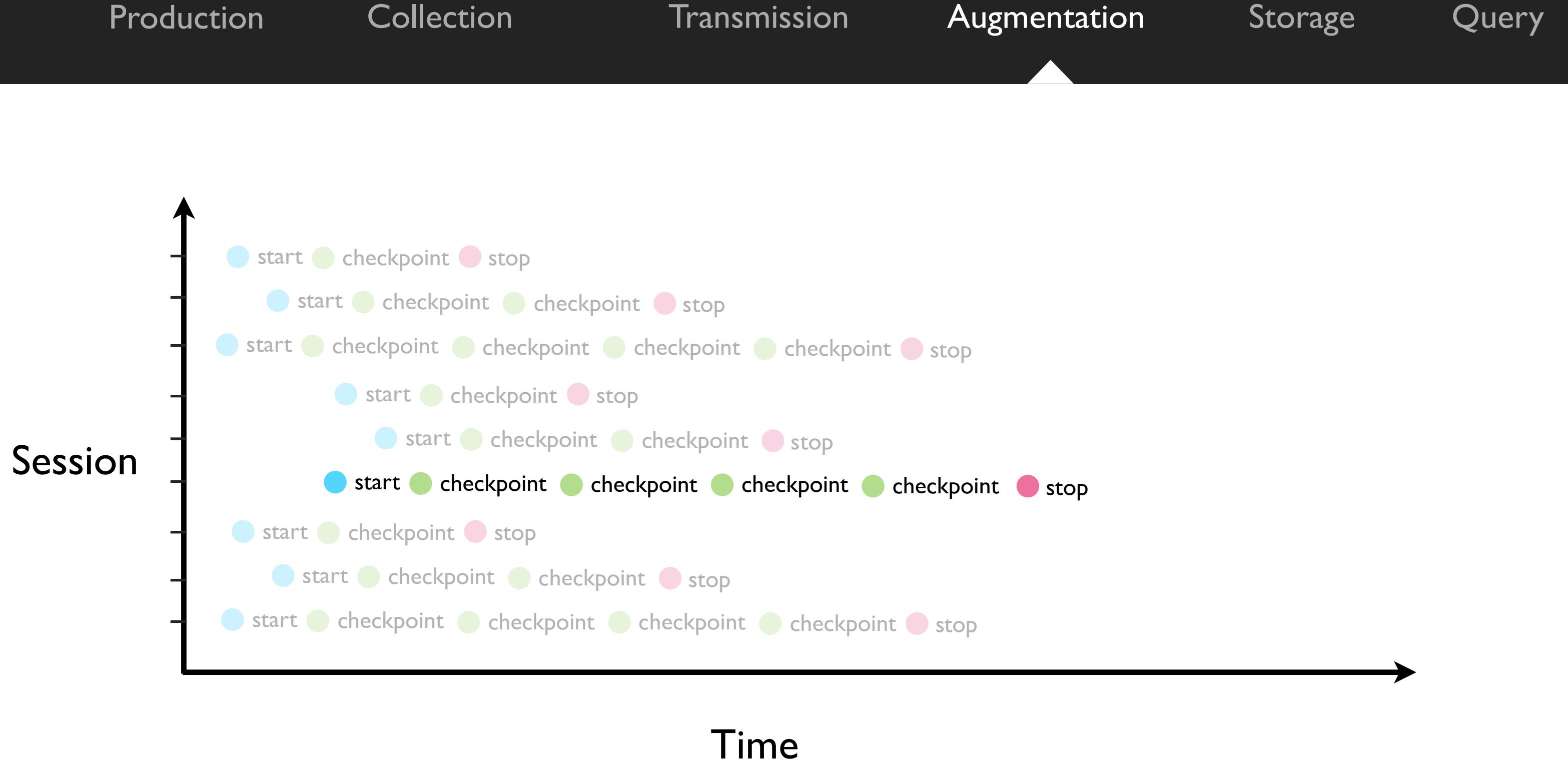
Transmission

Augmentation

Storage

Query





Production

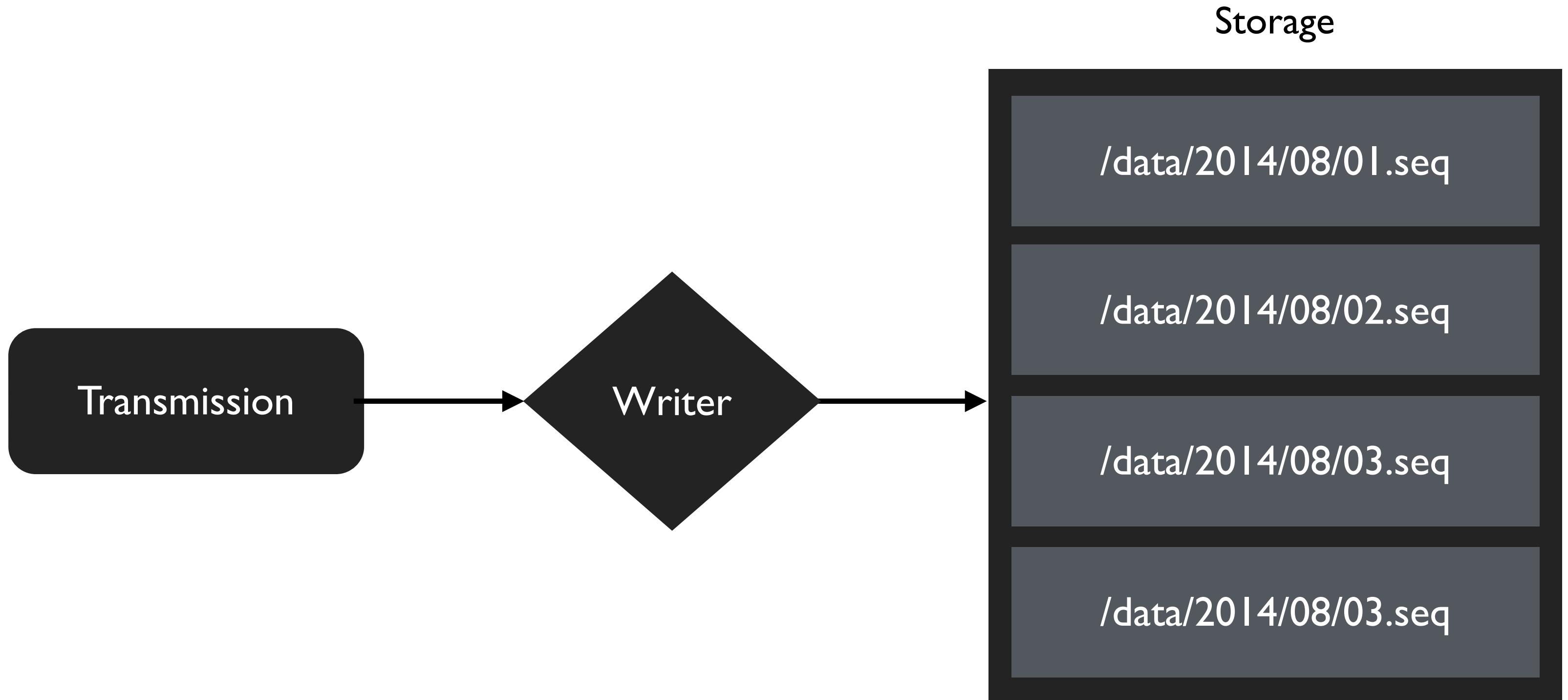
Collection

Transmission

Augmentation

Storage

Query



Production

Collection

Transmission

Augmentation

Storage

Query

Storage

- Source of truth
- Scalable
- Replicated

HDFS

- Self hosted
- Cost effective at scale
- Can run multi tenant
- Non trivial operational cost

S3

- Managed
- Cost prohibitive at scale
- Network to Map Reduce
- EMR cost

Production

Collection

Transmission

Augmentation

Storage

Query

Query

- Able to handle Billion of Records
- Support for statistics
- Bulk Loadable



Different workloads

- Low vs. high latency
- Common vs custom operations
- Large working sets



Columnar Store

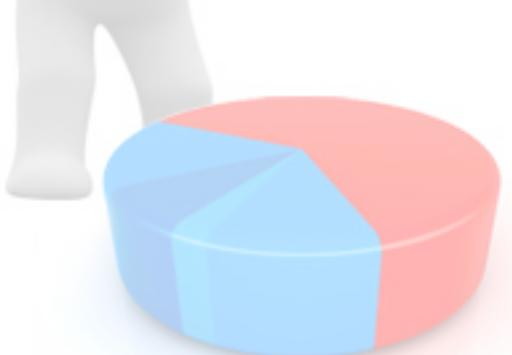
- Redshift/Vertica
- SQL
- Expensive

Hadoop Based

- Pig/Hive/Spark
- Shared tenancy
- Highly Scalable

Elastic Search

- Distributed Lucene
- Extendable
- Can't Bulk Load



Production

Collection

Transmission

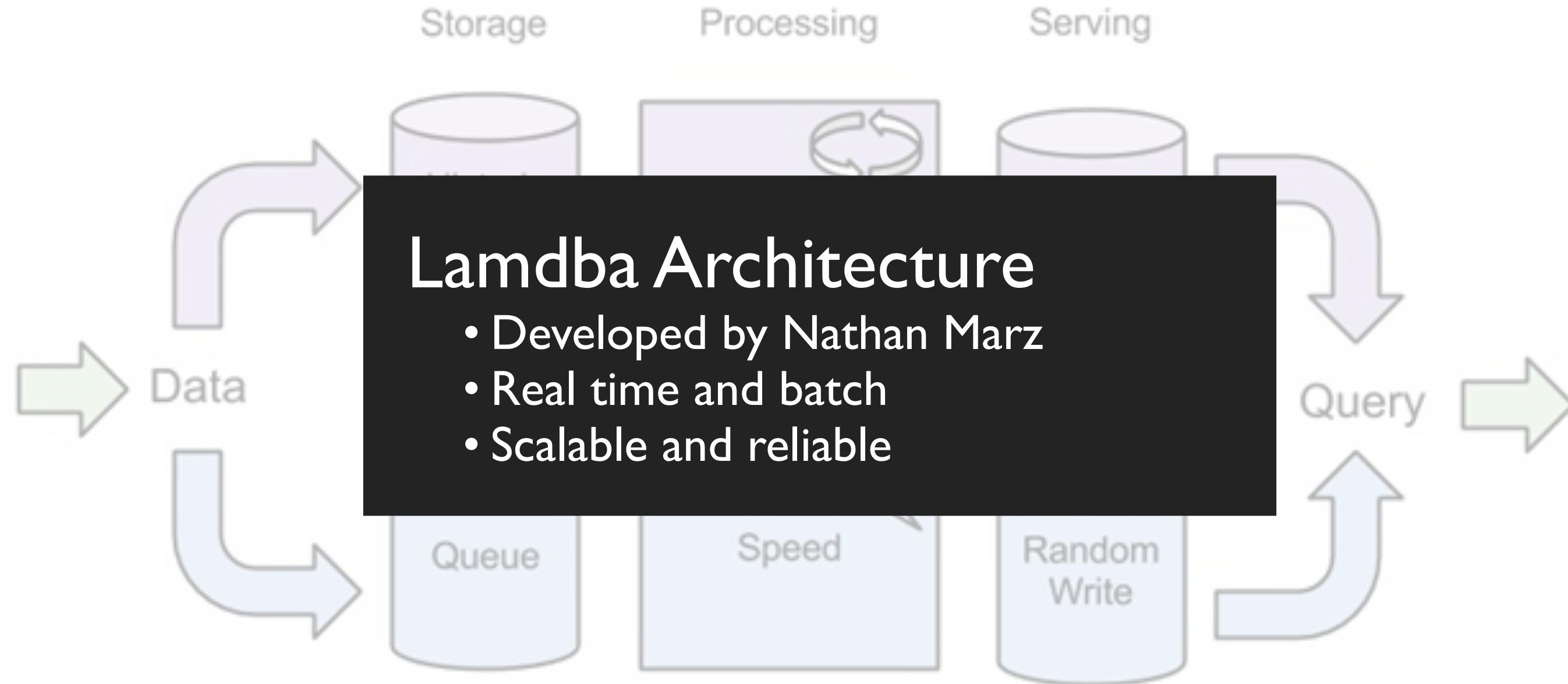
Augmentation

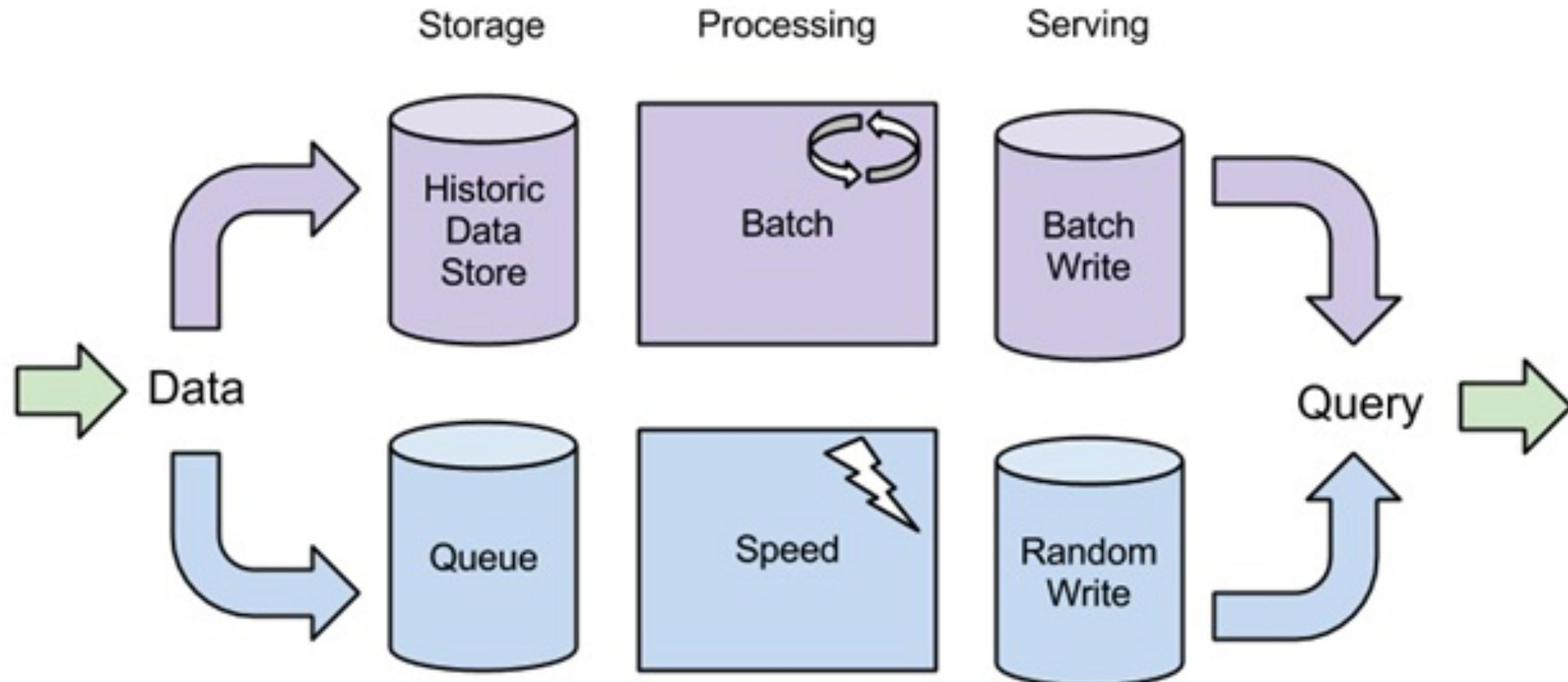
Storage

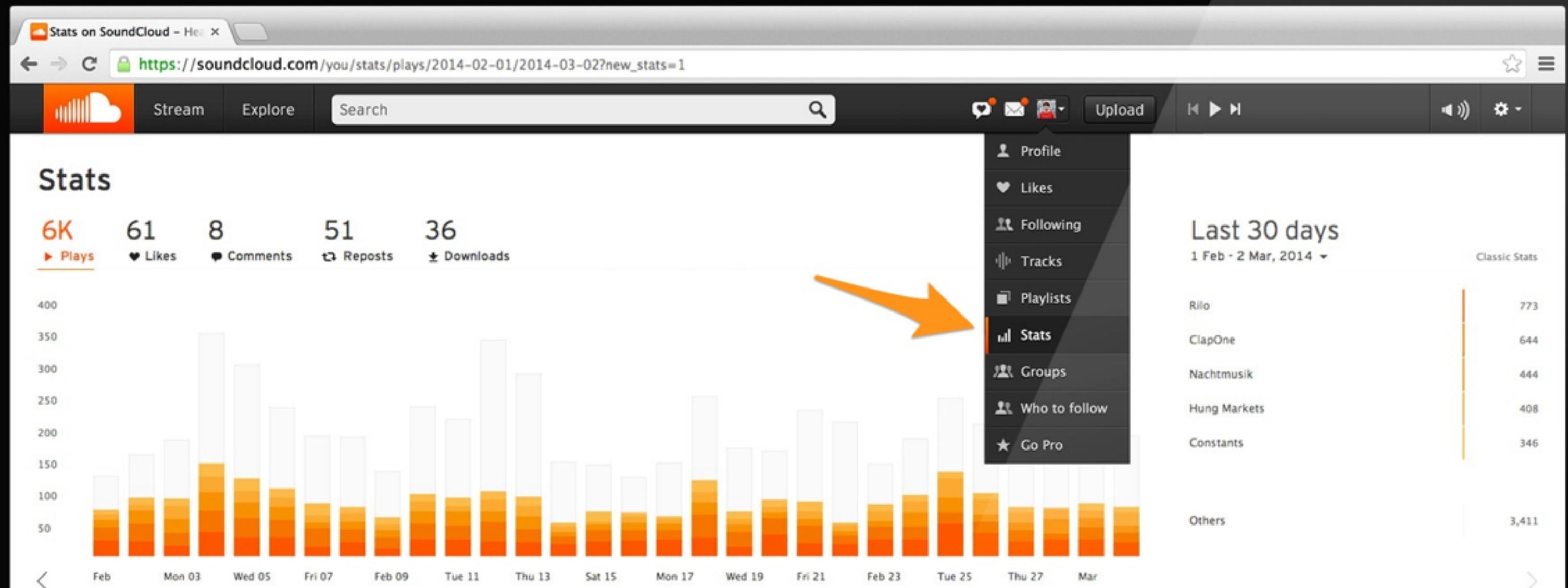
Query

Architecture Summary

- Connected and exchangeable
- Failure modes are a priority
- Source of truth
- Accessible, Insightful







Most played tracks



Rilo

773

Top countries



United States
1,252

Who played the most



detailloop
122

Websites



soundcloud.com/seams
1,696

Apps

2 ClapOne

3 Nachtmusik

4 Hung Markets

644

444

408

2 Germany

3 United Kingdom

4 France

1,238

737

274

2 gurū

3 sanjay.fernandes

4 Luciano Castello

94 soundcloud.com/seams/sets/quarters
71 soundcloud.com/stream
39 soundcloud.com/cecile9/likes

206

134

95



Cloud Computing

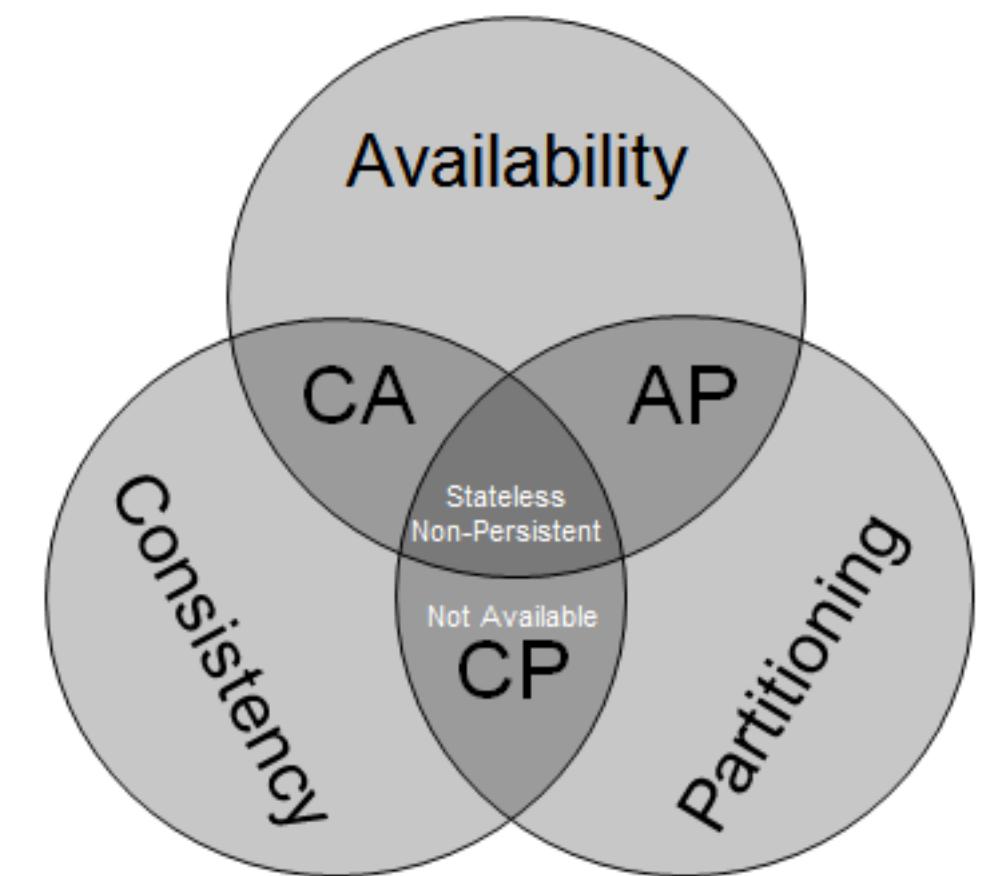
- Commodity Computation
- Outsourced infrastructure
- Has two architectural properties

Immutable Infrastructure

- Atomic software
- Compile to AMI/Docker/X
- Reproducible

Coordination

- Complex distributed system
- Different for every component
- Can't always buy



Conclusion

- Connected set of components
- Failure modes are a priority
- Build vs Adopt
- Components are volatile
- Architecture is “stable”

Coordinates: (x,y,z)

Velocity Components: (u,v,w)

Time : t
Pressure: p
Density: ρ
Stress: τ
Total Energy: Et

Heat Flux: q
Reynolds Number: Re
Prandtl Number: Pr

continuity:

$$\frac{\partial \rho}{\partial t} + \frac{\partial(\rho u)}{\partial x} + \frac{\partial(\rho v)}{\partial y} + \frac{\partial(\rho w)}{\partial z} = 0$$

- Momentum:

$$\frac{\partial(\rho u)}{\partial t}$$

Open Problems

- Anomaly detection
- Data completeness
- Anonymous authentication

- Momentum:

$$\frac{\partial(\rho v)}{\partial t}$$

dx

dy

dz

dy

$$\frac{1}{Re_r} \left[\frac{\partial \tau_{xx}}{\partial x} + \frac{\partial \tau_{xy}}{\partial y} + \frac{\partial \tau_{xz}}{\partial z} \right]$$

$$\frac{1}{Re_r} \left[\frac{\partial \tau_{xy}}{\partial x} + \frac{\partial \tau_{yy}}{\partial y} + \frac{\partial \tau_{yz}}{\partial z} \right]$$

- Momentum

$$\frac{\partial(\rho w)}{\partial t} + \frac{\partial(\rho uw)}{\partial x} + \frac{\partial(\rho vw)}{\partial y} + \frac{\partial(\rho w^2)}{\partial z} = - \frac{\partial p}{\partial z} + \frac{1}{Re_r} \left[\frac{\partial \tau_{xz}}{\partial x} + \frac{\partial \tau_{yz}}{\partial y} + \frac{\partial \tau_{zz}}{\partial z} \right]$$

Energy:

$$\frac{\partial(E_T)}{\partial t} + \frac{\partial(uE_T)}{\partial x} + \frac{\partial(vE_T)}{\partial y} + \frac{\partial(wE_T)}{\partial z} = - \frac{\partial(up)}{\partial x} - \frac{\partial(vp)}{\partial y} - \frac{\partial(wp)}{\partial z} - \frac{1}{Re_r Pr} \left[\frac{\partial q_x}{\partial x} + \frac{\partial q_y}{\partial y} + \frac{\partial q_z}{\partial z} \right]$$



Questions?