

# What We Got Wrong

Lessons from the Birth of Microservices at Google

November 6, 2018

# **Part One: The Setting**



Still betting big on the  
Google Search Appliance

“Those Sun boxes are  
so expensive!”

“Those linux boxes are  
so unreliable!”

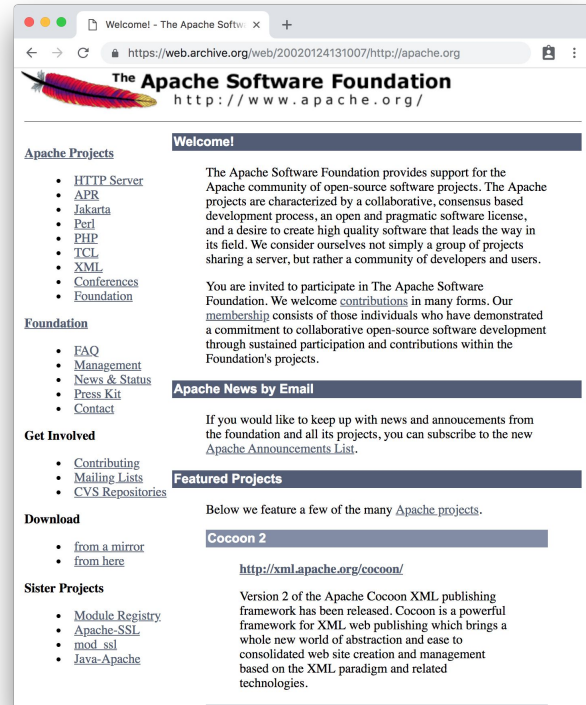
“Let’s see what’s on  
GitHub first...”

*– literally nobody in 2001*

# “GitHub” circa 2001



The screenshot shows the 'Projects' page of the Apache Software Foundation. At the top left is a cartoon illustration of a ram's head. Below it is a 'Table of Contents' with links to 'Introduction', 'Descriptions of Some GNU Projects' (including 'Lists of other GNU Projects' and 'Descriptions of Other Projects'), and 'Introduction'. The 'Introduction' section contains two paragraphs of text and a link to 'Descriptions of GNU Projects'. Below that is a list of software projects including Classpath, Free Film Project, GPKCS-11, GNU Cobol, GNU Enterprise, GNU GLUE, GNU XML, WebDAV and RTSP, GNU Octal, GYVE, Corel Draw, Harmony, and HURD.



The screenshot shows the Apache Software Foundation homepage. The browser address bar indicates the URL: https://web.archive.org/web/20020124131007/http://apache.org. The page features the Apache logo (a colorful feather) and the text 'The Apache Software Foundation http://www.apache.org/'. The main content is organized into several sections: 'Welcome!' with a paragraph about the foundation's support for open-source projects; 'Apache News by Email' with a link to the Apache Announcements List; 'Featured Projects' featuring Cocoon 2 with a link to http://xml.apache.org/cocoon/; 'Sister Projects' listing Module Registry, Apache-SSL, mod\_ssl, and Java-Apache; 'Download' with links for Cocoon 2; 'Get Involved' with links for Contributing, Mailing Lists, and CVS Repositories; 'Foundation' with links for FAQ, Management, News & Status, Press Kit, and Contact; and 'Apache Projects' with a list of projects including HTTP Server, APER, Jakarta, Perl, PHP, TCL, XML, Conferences, and Foundation.

# Engineering constraints

- Must DIY:
  - Wacky scaling requirements
  - Utter lack of alternatives
- Must scale horizontally
- Must build on commodity hardware that fails often



# Google eng cultural hallmarks, early 2000s

- Intellectually rigorous
- Bottoms-up eng decision-making
- Aspirational
- ... and maybe a little overconfident :-/

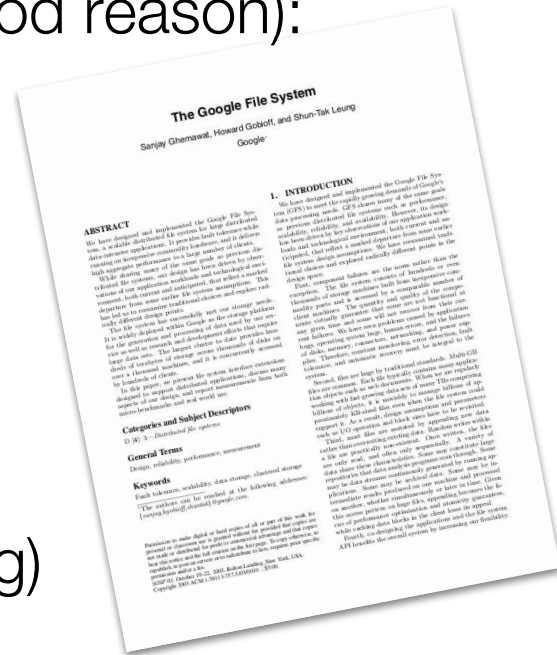
# **Part Two: What Happened**



# Cambrian Explosion of Infra Projects

Eng culture idolized epic infra projects (for good reason):

- GFS
- BigTable
- MapReduce
- Borg
- Mustang (web serving infra)
- SmartASS (ML-based ads ranking+serving)



# Convergent Evolution?

Common characteristics of the most-admired projects:

- Identification and leverage of horizontal scale-points
- Factored-out user-level infra (RPC, discovery, load-balancing, eventually tracing, auth, etc)
- Rolling upgrades and frequent (~weekly) releases

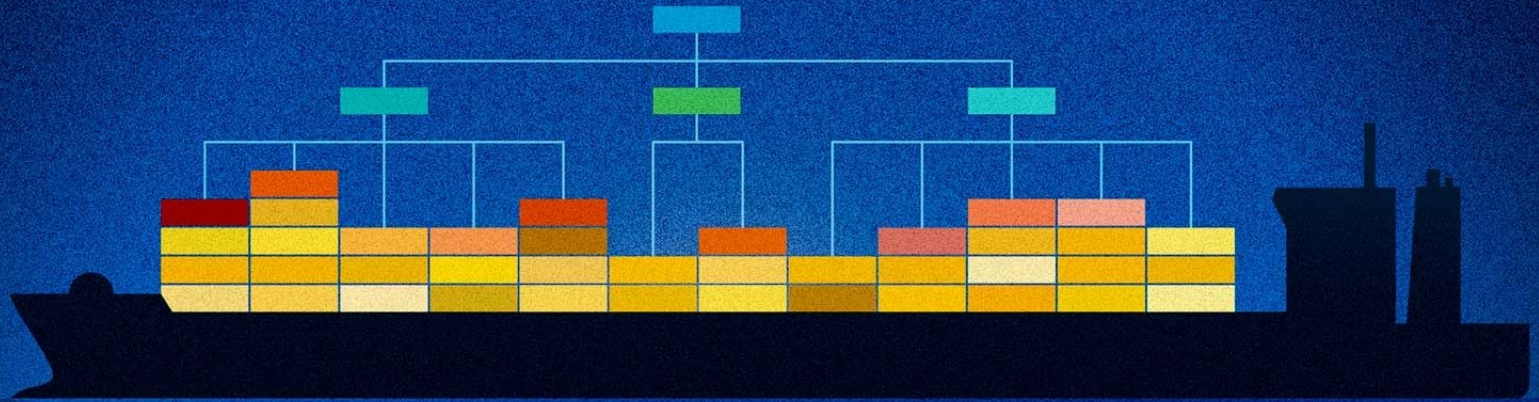
*Sounds kinda familiar...*

# **Part Three: Lessons**

# **Lesson 1**

Know Why

# The only good reason to adopt microservices



You will inevitably ship your org chart



# Accidental Microservices

- “Microservices due to Computer Science,”  
not org charts!
- Ended up with something similar to modern  
microservice architectures ...
- ... but for different reasons (and that  
eventually became a problem)

What's best for Search+Ads is  
best for all!

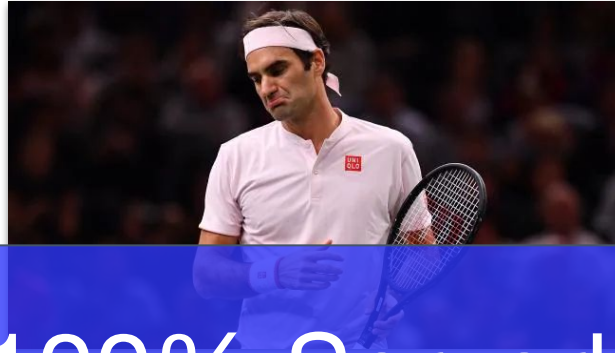
What's best for Search+Ads is  
best for ~~all~~ just the massive,  
planet-scale services

“But I just want to serve 5TB!!”

# Lesson 2

Serverless Still Runs on Servers

# An aside: what do these things have in common?



All 100% Serverless!



```
1 from __future__ import print_function
2
3 import json, urllib2, boto3
4
5
6 def lambda_handler(event, context):
7     response = urllib2.urlopen('https://ip-ranges.amazonaws.com/ip-ranges.json')
8     json_data = json.loads(response.read())
9     new_ip_ranges = [x['ip_prefix'] for x in json_data['prefixes'] if x['service'] == 'cloudfront']
10    print(new_ip_ranges)
11
12    ec2 = boto3.resource('ec2')
13    security_group = ec2.security_group('sg-3xxxxx3x')
14    current_ip_ranges = [x['cidr_ip'] for x in security_group.ip_permissions[0]['ip_ranges']]
15    print(current_ip_ranges)
16
17    params_dict = {
18        u'prefixlistids': [],
19        u'fromport': 8,
20        u'toport': 8080,
21        u'ipranges': [],
22        u'ipprotocol': 'tcp',
23        u'useridgrouppairs': []
24    }
25
26    authorize_dict = params_dict.copy()
```

# About “Serverless” / FaaS

## Numbers every engineer should know

### Latency Comparison Numbers (~2012)

L1 cache reference	0.5	ns			
Branch mispredict	5	ns			
L2 cache reference	7	ns		14x L1 cache	
Mutex lock/unlock	25	ns			
Main memory reference	100	ns		20x L2 cache, 200x L1 cache	
Compress 1K bytes with Zippy	3,000	ns	3	us	
Send 1K bytes over 1 Gbps network	10,000	ns	10	us	
Read 4K randomly from SSD*	150,000	ns	150	us	~1GB/sec SSD
Read 1 MB sequentially from memory	250,000	ns	250	us	
Round trip within same datacenter	500,000	ns	500	us	
Read 1 MB sequentially from SSD*	1,000,000	ns	1,000	us	1 ms ~1GB/sec SSD, 4X memory
Disk seek	10,000,000	ns	10,000	us	10 ms 20x datacenter roundtrip
Read 1 MB sequentially from disk	20,000,000	ns	20,000	us	20 ms 80x memory, 20X SSD
Send packet CA->Netherlands->CA	150,000,000	ns	150,000	us	150 ms

### Notes

-----  
1 ns =  $10^{-9}$  seconds  
1 us =  $10^{-6}$  seconds = 1,000 ns  
1 ms =  $10^{-3}$  seconds = 1,000 us = 1,000,000 ns

### Credit

-----  
By Jeff Dean: <http://research.google.com/people/jeff/>  
Originally by Peter Norvig: <http://norvig.com/21-days.html#answers>

# About “Serverless” / FaaS

## Latency Numbers Every Programmer Should Know

■ 1 ns

■ L1 cache reference: 0.5 ns

■ Branch mispredict: 5 ns

■ L2 cache reference: 7 ns

■ Mutex lock/unlock: 25 ns

■ 100 ns

■ Main memory reference: 100 ns

■ = 1  $\mu$ s

■ Compress 1 KB with Zippy: 3  $\mu$ s

■ = 10  $\mu$ s

■ Send 1 KB over 1 Gbps network: 10  $\mu$ s

■ SSD random read (1GB/s SSD): 150  $\mu$ s

■ Read 1 MB sequentially from memory: 250  $\mu$ s

■ = 1 ms

■ Round trip within same datacenter: 100  $\mu$ s

■ = 1 ms

■ Read 1 MB sequentially from SSD: 1 ms

■ Disk seek: 10 ms

■ Read 1 MB sequentially from disk: 20 ms

■ Packet roundtrip: 100  $\mu$ s

■ = 1 ms

■ = 1 ms

■ = 1 ms

■ = 1 ms

■ = 1 ms

■ = 1 ms

■ = 1 ms

■ = 1 ms

Source: <https://gist.github.com/2841832>



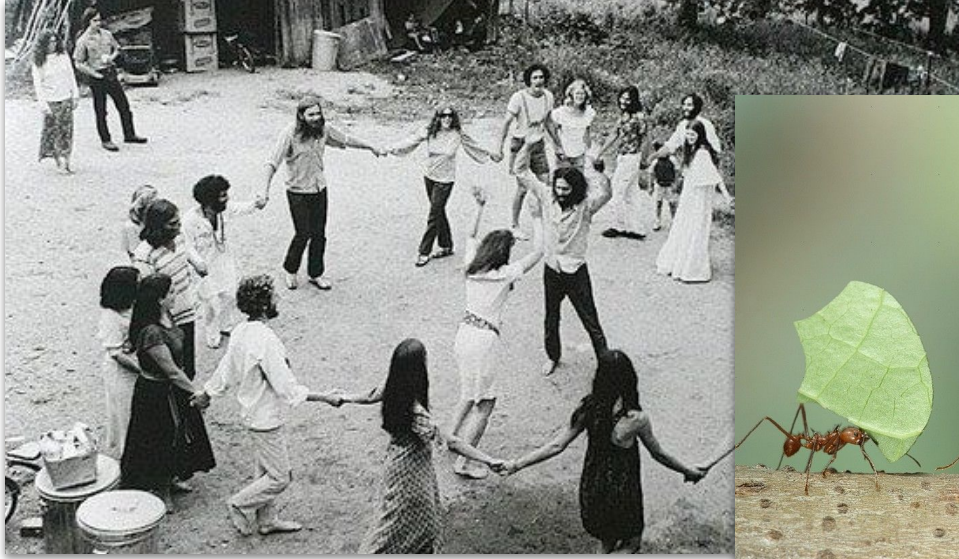
# About “Serverless” / FaaS

- Embarrassingly parallel / stateless situations do exist
- FaaS are great for them
- ... but *caching*

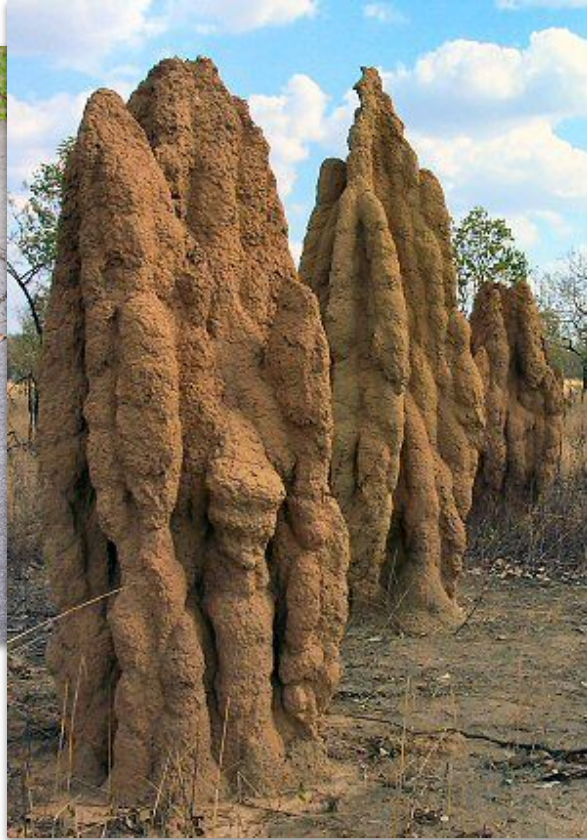
# **Lesson 3**

“Independence” is  
not an Absolute

# Hippies vs Ants



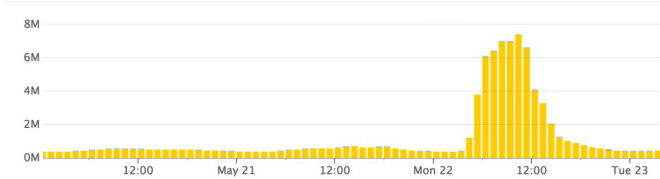
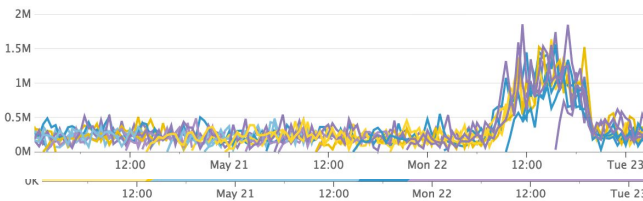
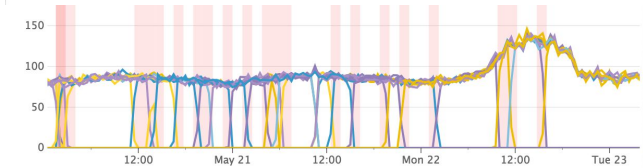
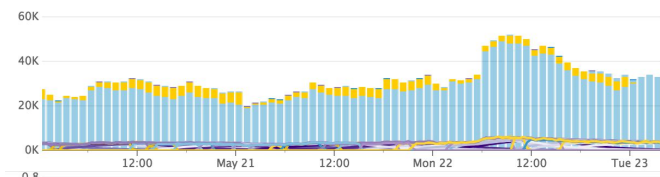
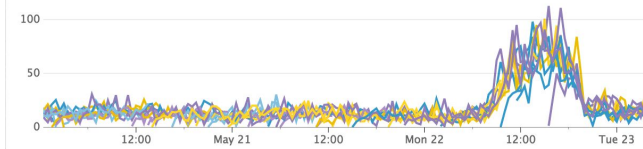
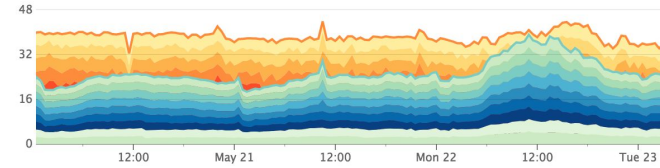
# More Ants!



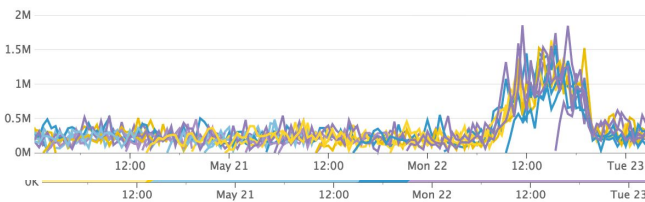
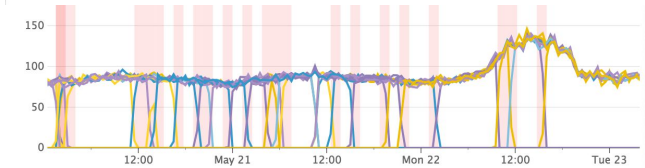
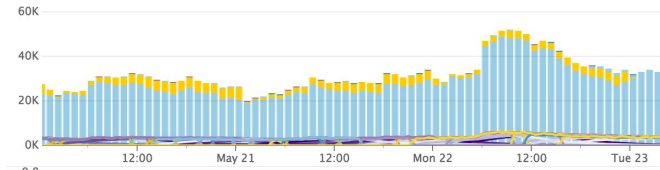
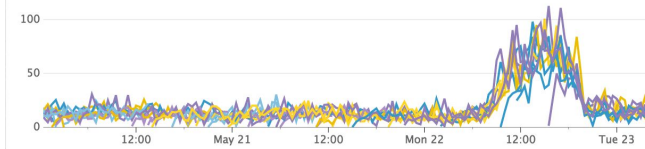
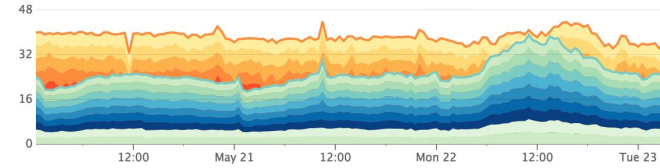
# **Lesson 4**

Beware Giant Dashboards

# We caught the regression!



... but which is the culprit?



# Observability boils down to two activities

1. Detection of critical signals
2. Refining the search space

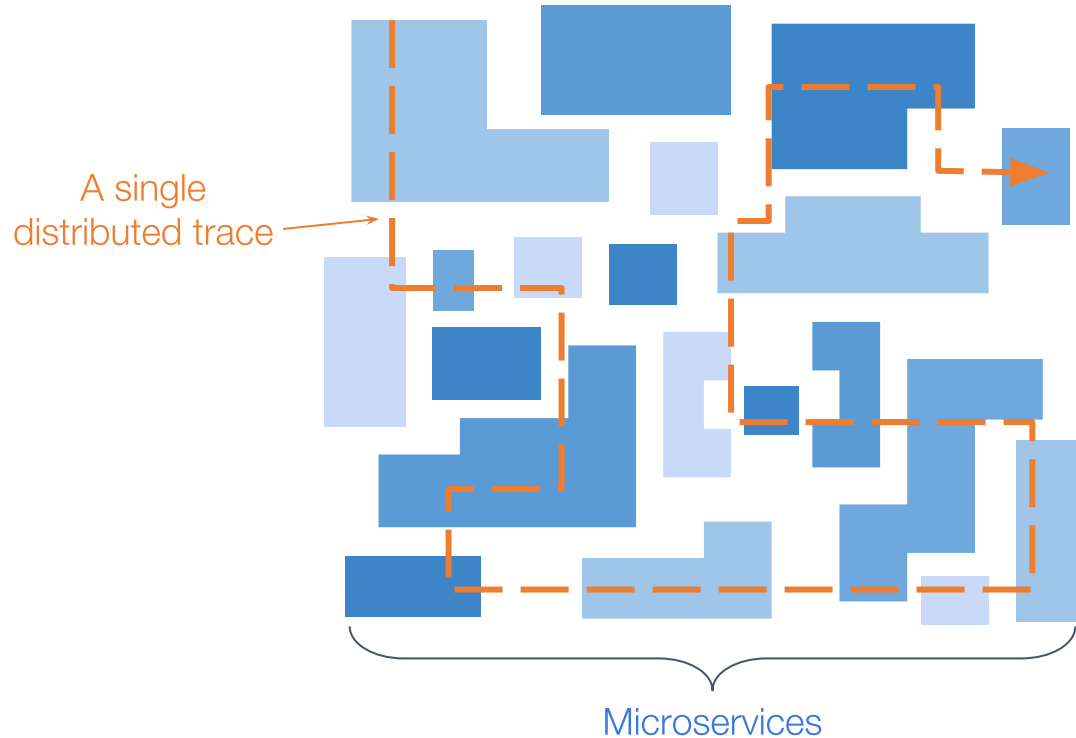
Don't confuse “visualizing the *entire* search space”  
with “*refining* the search space”



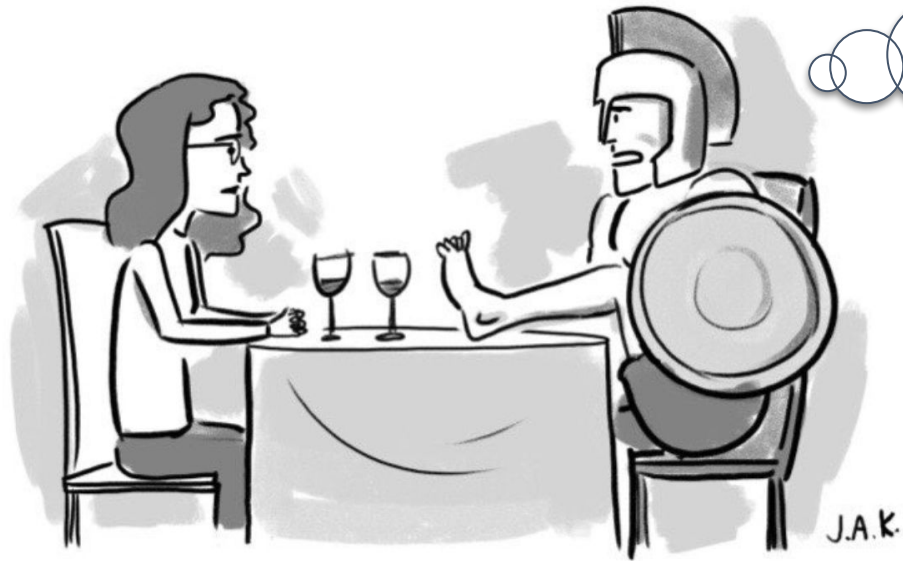
# Lesson 5

You Can't Trace Everything  
(or can you?)

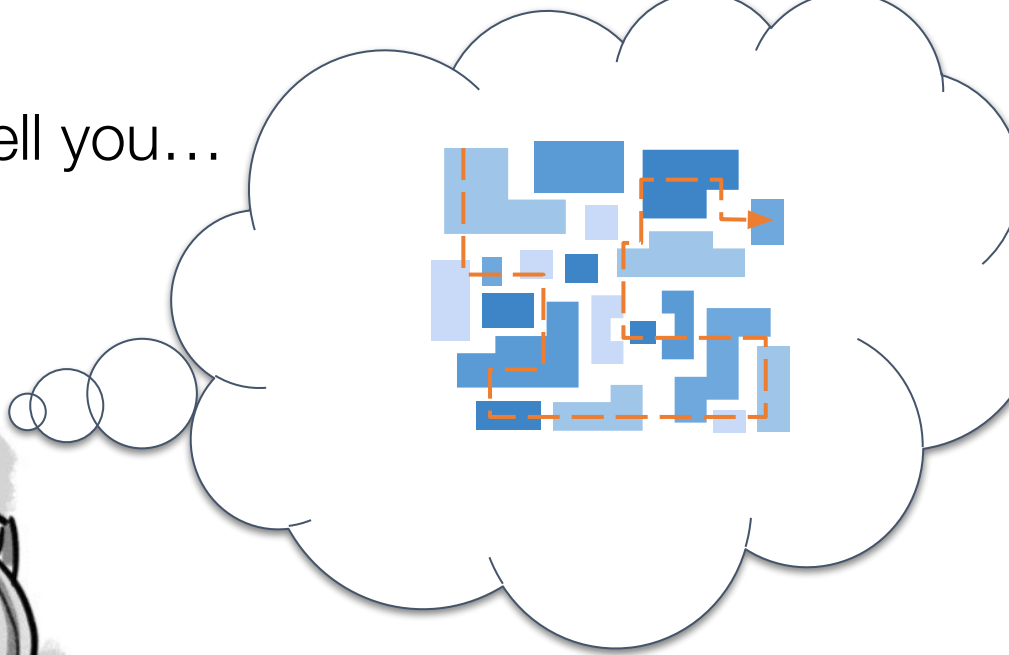
# Distributed Tracing 101



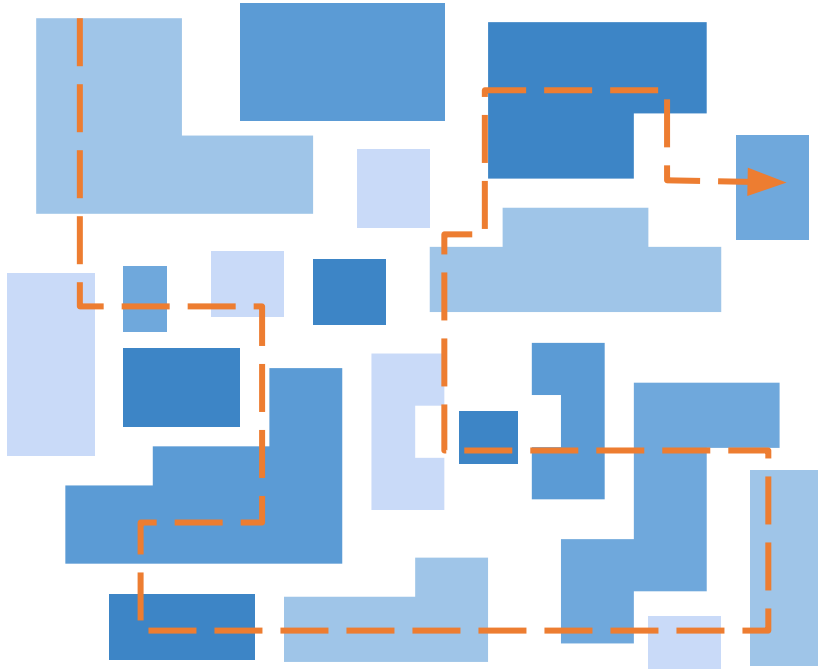
There's something I need to tell you...



"I'm ready to be vulnerable."



# Trace Data Volume: a reality check



app transaction rate  
x # of microservices  
x cost of net+storage  
x weeks of retention

-----  
*way* too much \$\$\$\$

# The Life of Trace Data: Dapper

Stage	Overhead affects...	Retained
Instrumentation Executed	App	100.00%
Buffered within app process	App	000.10%
Flushed out of process	App	000.10%
Centralized regionally	Regional network + storage	000.10%
Centralized globally	WAN + storage	000.01%

# The Life of Trace Data: ~~Dapper~~ Other Approaches

Stage	Overhead affects...	Retained
Instrumentation Executed	App	100.00%
Buffered within app process	App	100.00%
Flushed out of process	App	100.00%
Centralized regionally	Regional network + storage	100.00%
Centralized globally	WAN + storage	on-demand

**Almost Done...**

# Let's review...

- Two drivers for microservices: what are you solving for?
  - Team independence
  - “Computer Science”
- Understand the appropriate scale for any solution
- Services can be too small (i.e., “the network isn't free”)
- Hippies vs Ants
- Observability is about *Detection* and *Refinement*
- We *can* trace everything



# Thank you

Ben Sigelman, Co-founder and CEO

twitter: @el\_bhs

email: [bhs@lightstep.com](mailto:bhs@lightstep.com)



PS: Working at LightStep is fun!  
[lightstep.com/culture](https://lightstep.com/culture)

