# INDEX
## SAN FRANCISCO

Discover. Collaborate. Deploy.

# Fixing Under Fire

Aiton S Goldman

Michael Lipschultz
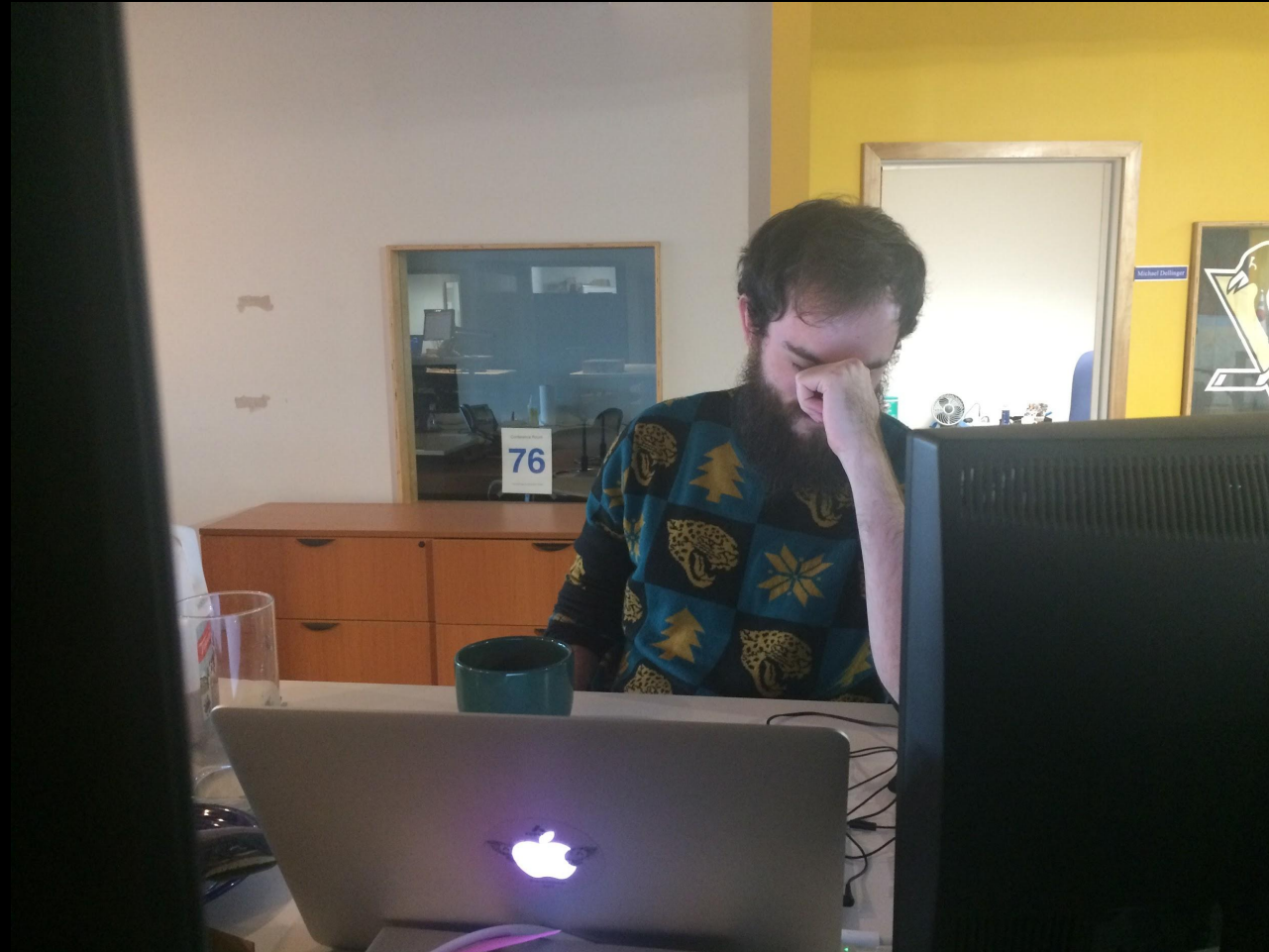
# Fixing Under Fire



we learned a valuable lesson today

theprofoundprogrammer.com

# The first shift….OF TERROR!

# The problem

# Questions we should have asked

# "What is the on call person's duties?"

- The on-call person needs to be allowed time to do more than just fix the current fire

    - Needs to update playbooks

    - Needs to update logging/metrics

    - Needs to fix root causes

# "How bad are things?"

- The on-call person needs a way to communicate to the rest of the group what is happening

# "How bad are things?"

# THE END

# IT PAGED AT 5 PM!

# The problem

# Questions we should have asked

# "How do we tell upstream services something has gone wrong?"

- Who reported the error?



BAD



GOOD

# "How do I know what to focus on when there is trouble?"

- Is this problem impacting the customer?

- Can I tell if the problem is intermittent?

- Is my service healthy?

# "How do I know what to focus on when there is trouble?"

## Is my service healthy?

# "How do I know what to focus on when there is trouble?"

# "How do I know what to focus on when there is trouble?"

# IT CAME FROM THE USER!

# It was a week like any other ...

… at first

Alert

Alert

Alert

Alert

Alert

Alert

Alert

Alert

Alert

Alert

Alert

Alert

Alert

Alert

Alert

# Questions we should have asked

# "How do we avoid flooding the on call person with information?"

- Avoid double reporting

# "How do we avoid flooding the on call person with information?"

- Avoid double reporting?

HELP!

Something happened..

| Micro Service one | Micro Service two | Micro Service three | Micro Service four |

Metrics Based Alerting

# "How do we avoid flooding the on call person with information?"

- Making things that wake us up configurable

Fire alert when there are 12 failures

Fire alert when successes are less than 5

Metrics Based Alerting

Fire alert when there were failures for at least 3 minutes

Fire alert when the rate of failures is too high

# "Can I tell if it's a specific customer/group of customers causing the problem?"

All logs:

Logs for
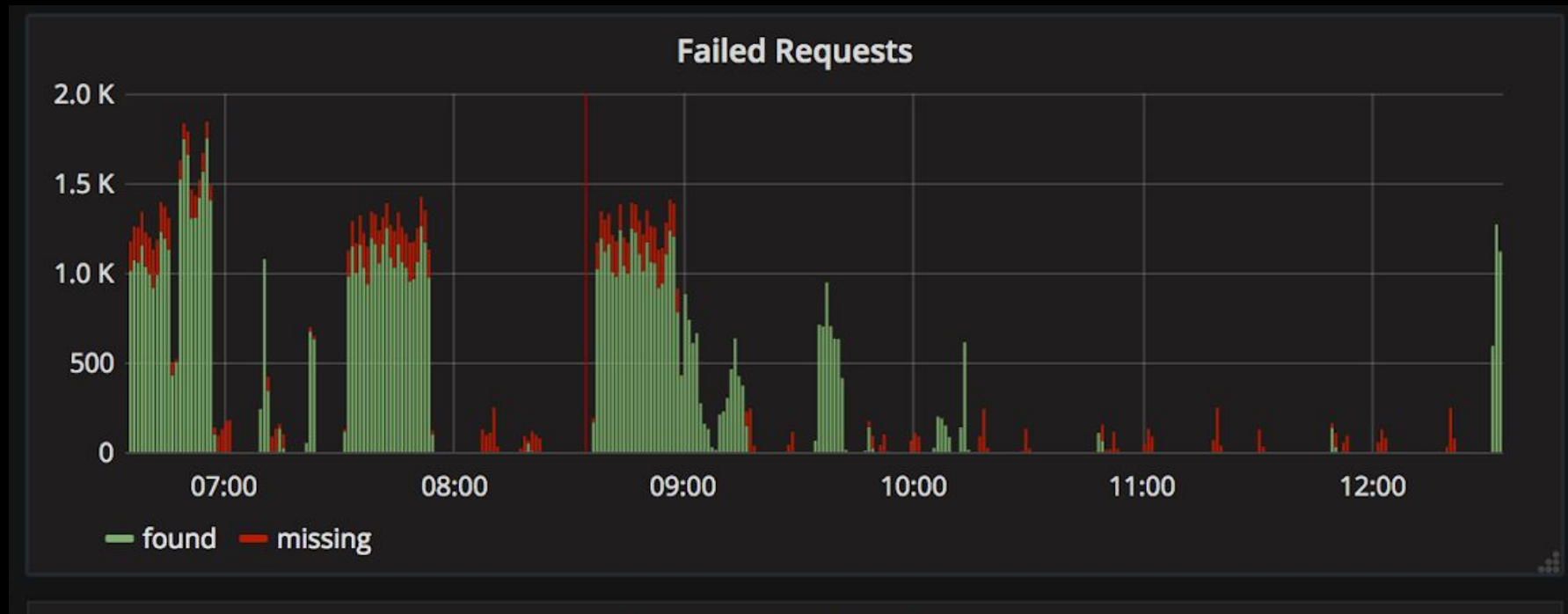specific customer:

# From beyond the proxy!

# The problem

# Questions we should have asked



THE ULTIMATE QUESTION
Which is better
Kirk or Picard

motifake.com

# "Can we trust the external service?"

- What we should have done...

Configurable Retries

Micro Service 2

External Service

# "Can I reproduce a request to the external service?"

- We needed to put as much info (or meta info) as possible into logging with regards to connecting with external services

- We needed to track a request all the way through the stack

# Who you gonna call?

# "Ray, this looks extraordinarily bad"

# Questions we should have asked...

# "Who do we depend on?"

- Watson has over 2000 IBM employees
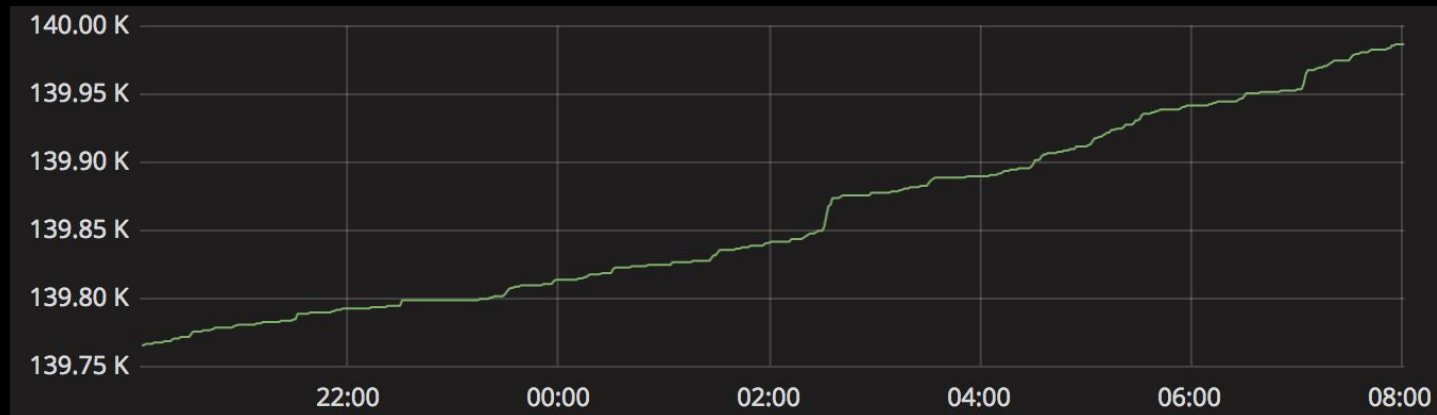
- That doesn't include all the teams outside Watson we depend on

- How do we reach the ones we depend on at 2am?

# "Who do we depend on?"
## Revisiting "It Paged at 5PM!"

- What information would they need to help?

# "Who depends on us?"

- How do groups that depend on us communicate with us?

- How do we communicate with them?

# A quick fix...OF DOOM!

# Reboot it

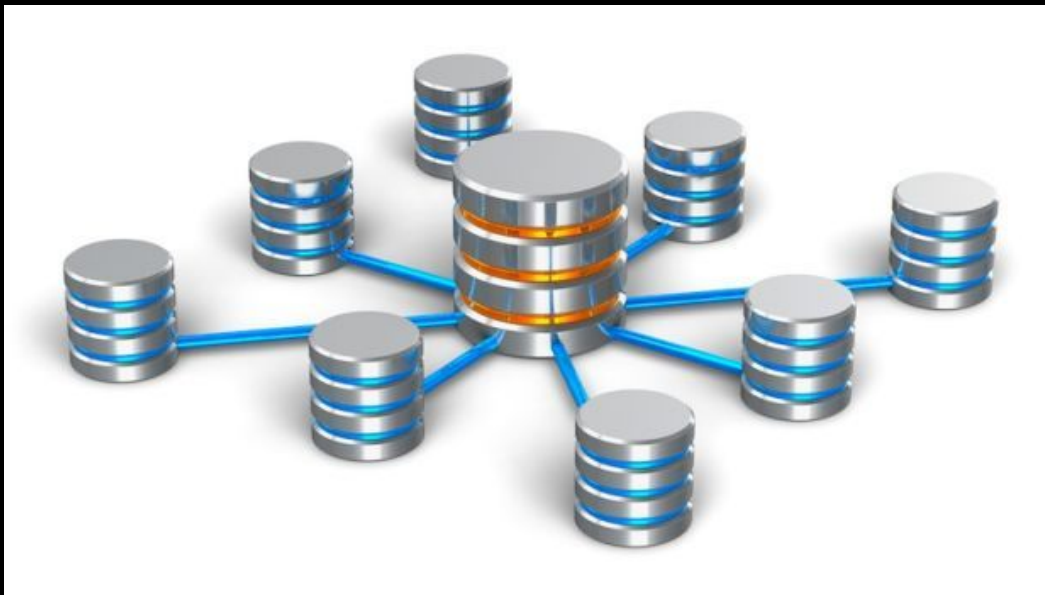

Even if your service is stateless..

● Connection limitations (i.e. to databases)
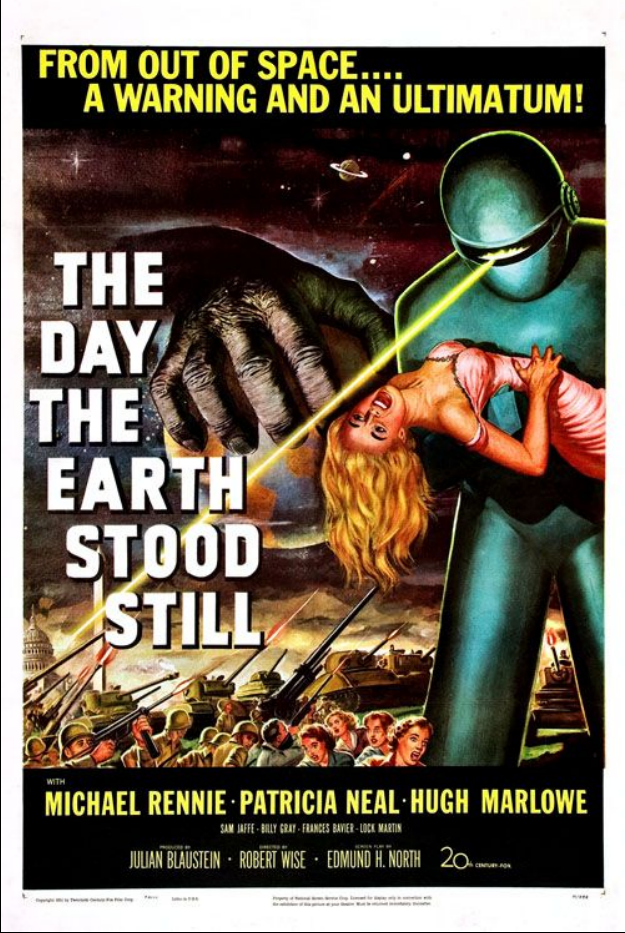
● Quotas

# Rollback



- Do I know what versions to roll back to?

- Do I need to consider DB schema?

# Cleanup the database



- Vacuum the database

- Free up connections

# What did we learn today?

# Lessons Learned

"Communicating internally and externally"



Sharing pager duty reports with the whole team

"How bad are we doing?"
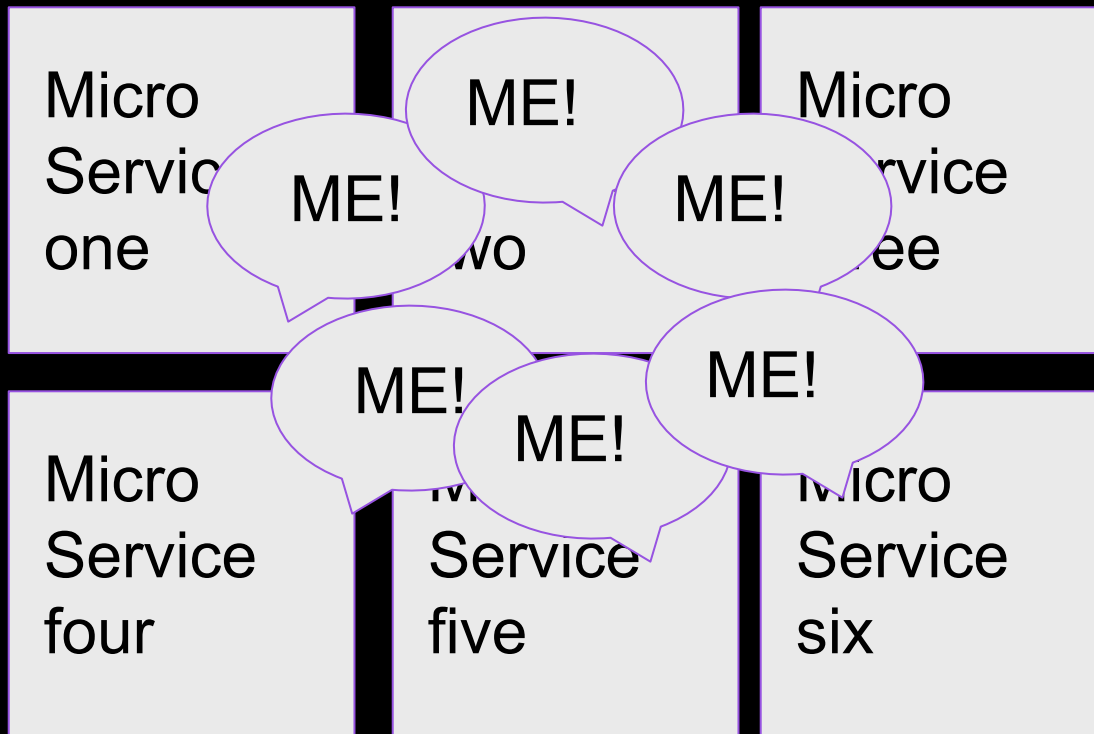
# Lessons Learned

"Communicating internally and externally"



Do we know how to contact other groups?
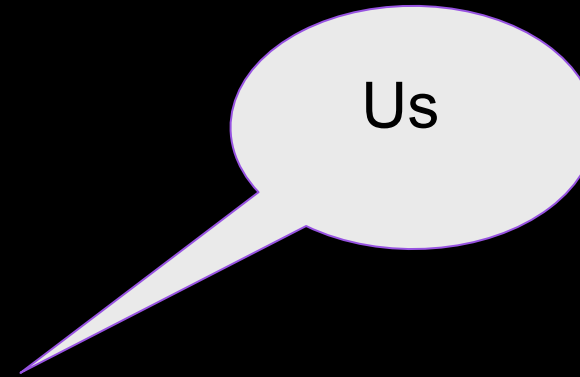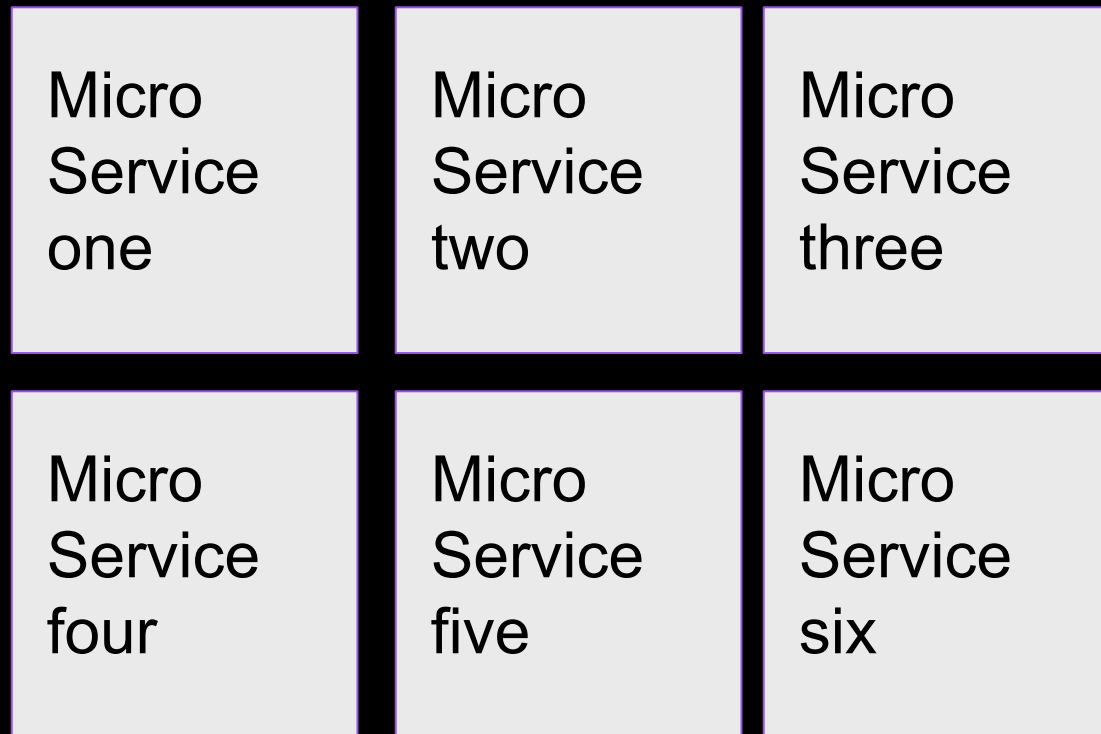Do they know how to contact us?

# Lessons Learned

"Are we treating our microservices as parts of a whole?"

# Lessons Learned

"Are we treating our microservices as parts of a whole?"

| | | |
|---|---|---|
| Micro Service one | Micro Service two | Micro Service three |
| Micro Service four | Micro Service five | Micro Service six |

Us

- Alerting
- Passing errors upstream
- Tracking requests through entire stack

# Lessons Learned

"What is the information I need first when there is a problem?"

# Lessons Learned

"I found the problem, now how do I fix the cause?"

# Final Thoughts

- Your first 4-6 months will suuuuuccckkk

- Use our questions to point yourself in the right direction

# Thanks



Michael Keeling

Anastas Stoyanovsky

Wentao Jiang

Chuck Gala

Joe Runde

# Extra Special Thanks to Eric Kaun

"Before pager duty, these pants were white!"