

④ Scalability + cost + risk + complexity → impl<sup>n</sup> of SRE

↓ biz. appl<sup>n</sup>s

## SRE Foundation

07 Aug 2021

⑤ Bring in operations team  
into the ecosystem earlier

DevOps: velocity + stability

**CALMS**

culture  
change

Reactive      Proactive      predictive

↓                  ↓                  ↓

Event  
monitoring  
Incident-  
management

you do

Reliability → clo collaboration without any issues.  
Resilience → how soon can you come back from  
(reactive) an unplanned outage

preventing  
incident  
and  
addressing  
at the  
earliest  
it occurs

Watermelon syndrome: { looks green from outside  
                                  } red from inside

SRE: → permanent fix to user errors

is a **discipline**      ⑥ apply soft. engg → infra + operns problems

→ how can we reduce problems?

→ In companies, L1 → removed / incorporated chatbots.

Goal of SRE: create ultra-scalable and highly distributed  
S/W systems

→ **Prerelease** → better user experience.

SREs → journey / guide Dev. & Platform team

to spread awareness about SREs  
and build the mindset

Goal: <sup>create</sup> ultra-scalable and highly reliable distributed software systems



50%  
(ops)

+ 50%  
(devp.)

+  
issue  
reops,

+  
new.

On-call & manual interventions

features, Scaling or autom'



Monitoring + Alerting + Automation

Part of

SRE

→ Reduce "Totl" → no manual stuff → autom'

- SRE: class that complements DevOps
  - \* SRE and DevOps coexist.
  - \* SRE can work without DevOps also.

part  
of

Platform Engg. Team

- SRE - Engg. discipline - achieve appropriate level of reliability
  - Biz. + performance + oper's insight well → Key to monitoring
  - growing SRE role
- Scalability  
availability  
incident response  
autom'
- Google: 50% → on the aggregate 'ops' count.  
50% → actually closing devp.

\* Auto-healing → red<sup>o</sup> in the volume of incidents

SRE ↔ \* Value stream mapping (VSM) + awareness  
 ↳ collective efforts

\* SRE → experimentation

- Experience Level Agreement (Ex-LA) : new dimension to SLA.

- SRE - specific impl<sup>o</sup>  
 of DevOps with some ext<sup>o</sup>s (Google)



6x Principles:

1) oper<sup>o</sup>s - s/w problem

↳ S/w engg. as a discipline focusses on

designing and building rather than operating and maintaining.

↳ oversee the NFR during s/w devp.

↳ estimations:

40-90% of the TCO are incurred after launch.

"Slowness is ~~not~~ the main load".

2) Service Levels

SLO → service level obj.



SLOs need consequences if they are violated.

↳ availability target for a product or service (more 100%).



↳ - In SRE, services are managed to the SLOs.

Tasks can provide the "wisdom of prod"

3) Tool: "wisdom of prod"

kind of work tool  
to running a prod<sup>n</sup> service  
that tends to  
be manual, repetitive

reduce tool → automate  
whatever can  
be automated.

→ system thinker  
→ level of understanding

4) Autom<sup>n</sup>: → reduce  
manual tasks.

+  
engineering-based approach  
to problems

KISS: Keep it simple, stupid.

5) reduce the cost of failure

(proactive)

MTTD: mean time  
to detect.

detecting  
before  
incident occurs.

look to  
improve MTTR (mean  
time  
to  
repair)

late problem (defect)  
discovery is expensive  
so SRE looks for ways  
to avoid this.

6) shared ownership

SRE's shift left  
and provide  
wisdom of  
prod<sup>n</sup> to devp.  
teams.

(boundaries between  
app.dev and prod<sup>n</sup>  
should be removed.)

SLOs → element of shared ownership

Greg McKeown

Essentialism → Book

Get up early → cognitive bandwidth

no gadget-day

1% improvement per single day

Module-IIService level objective: → binding target  
for a coll'n of SLIs.

- ~~SLOs~~ → tightly related → user exp.
- widely tracked SLO → availability → goal for how well a product or service should operate.
- products | services could have several SLI's.
- SLO is a subject of SLA.
- 100% reliability → expensive, technically complex

\* customer journey exp. (going to Amazon)

\* functional req's (feature teams)

\* NFRs (oper/infra. Team)

\* SLOs: → defined by SRE and product mgmt.

→ internal target

→ improvement over own perf. and enhance the benchmark.

moving target

SLO → become the new SLA. with increased user-experience.

depends on  
SLO's → business

↳ set according to customers.

SLA - 95% (Contractual agreement)

SLO - 96% (Business)

$$\text{Error Budget} = 100 - \text{SLO}$$

or

budget of unreliability (quantitative measure shared between product and SRE teams to balance innovation & stability)

SLO - 95%

$$\text{Error budget} = 100 - 95\% = 5\%$$



Breach SLO

+ increase to error budget

Consequences

calculated every day

Biz. prioritized period

stop deployments

Dev. Team + SRE Team meet daily to outline their improvements

Pins to be exhausted.

>Error budget → balance between reliability & stability

assigned every month or quarter.

difficult to achieve.

Quantitative risk measure to understand sys. behaviour.

most popular

SLOs - Availability (1%)

Latency (2%)

Error rate (1%)

error-budget

does the service deliver in a user-acceptable period of time?

\* Missed SLO: Post revenue, drop in employee productivity, impact on customers, social media backlash  
BIZ. perf.

(Service level indicator)  
SLI → quantitative measure of an attribute of the service  
eg throughput, latency, failures per request

Error-budget:

Not hitting 99.9% of HTTP req. in a month → scalability issue

Bad Good

Error-budget as going over budget

↓ work overtime or respond to out-of-hours issues.

↓

strategically burn the budget to zero every month

good velocity without comparing availability

\* whether it's for feature launches or arch. changes

\*

High-risk deployments or large big-bang changes have more likelihood of issues and more chance of blowout budget being blown.

← Fixed Error Budget:

encourage the Lean preference for small changes to stay within the error budget.

Consequences:  
of missed SLO's

no new deployments → sprint planning may only pull post-mortem action items from the backlog.

\*

VALET dimensions of SLO

(volume / availability / latency / errors / tickets)

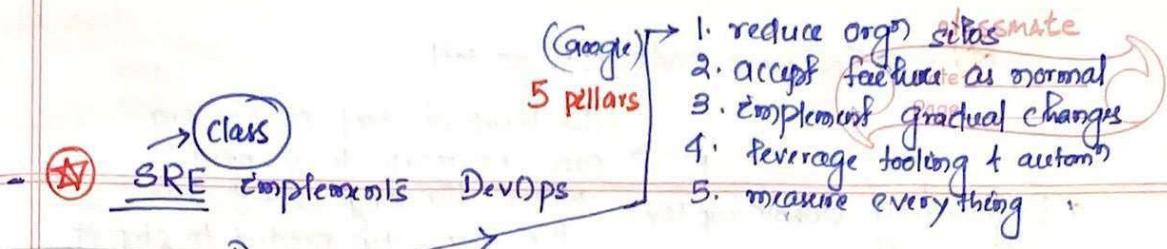
Dev Team + SRE Team

outline their improvements

- SRE is fundamentally more about people and processes than tech.

- purpose of SRE: bridge the gap between platform design, development and operational execution.

- Earlier we catch errors, the cheaper they are to fix.



**DevOps**

Set of principles & culture  
guidelines  
that helps to breakdown  
the silos between  
devp. & operns

**SRE**

a set of practices  
with an emphasis  
of strong engineering  
capabilities that  
implement the  
DevOps practices

**Google:** SRE is a specific empn of DevOps with some extensions.

- measuring reliability:

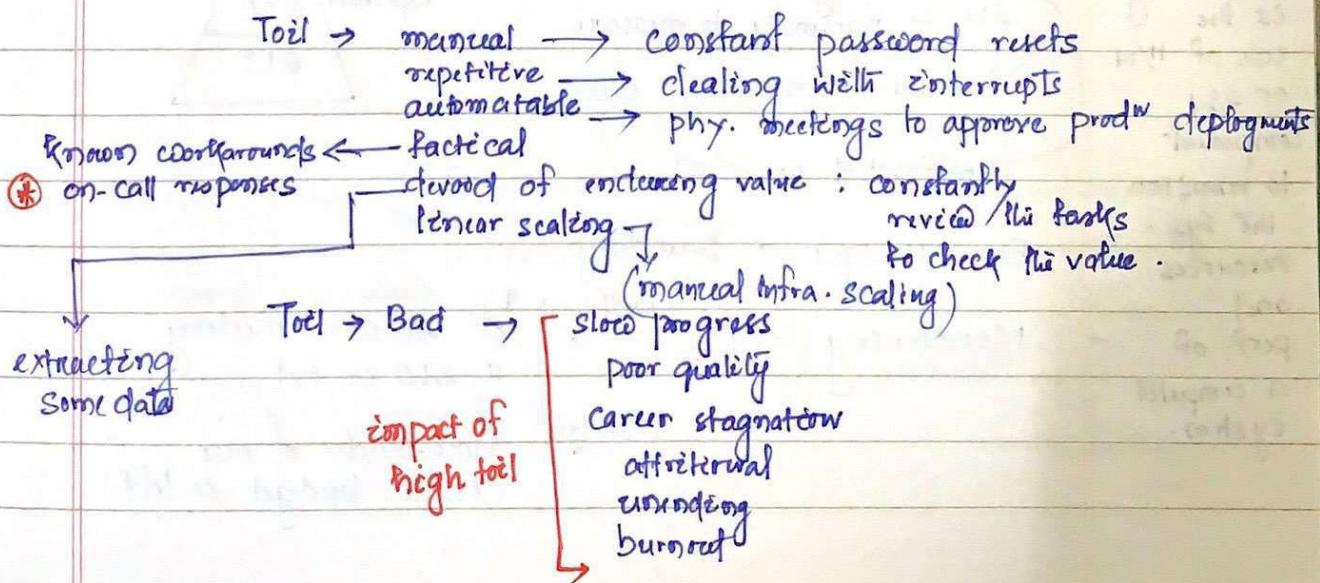
$$\text{availability} = \frac{\text{actual system uptime}}{\text{anticipated sys. up time}}$$

**08/08:**

**SRE** focusses mostly on  
post-prod responsibilities

**Goal of SRE:** to create ultra-scalable  
and highly reliable  
distributed sys. systems .

**Wisdom of prod** → wisdom gained  
from something  
running in prod





Engineers work 100% on tool.

The level of tool in an org. can increase to a point.

\* Engineering Bottleneck: where the orgn won't have the capacity needed to stop it.

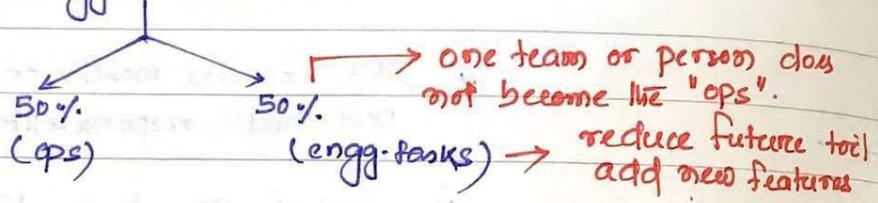
What can be done about Tool?

- ext. automs (scripts, automation tools)  
 - int. automs (automs delivered as part of the service) enhancing the service to not require intervention  
 ↗ need tech.  
 process engineering

- identify tool / effort estimation (operational work)
- Add Tool to the product backlog



- making engg. time available



- Slack: committed to the emp. of reliability over feature velocity

System monitoring is the use of HW or SW components to monitor the sys. resources and perf. of a computer system.

Monitoring & SLI: ways for engg. to communicate quantitative data about systems.

SLO → target (internal) based on user experience  
 SLI → parameter to measure

actual quantitative data

→ estimated expected



Error Budget → boundary

Monitoring SLI's → will tell if we are meeting a SLO or not

also tell how much of our error budget is left.

### SLI Measurement:

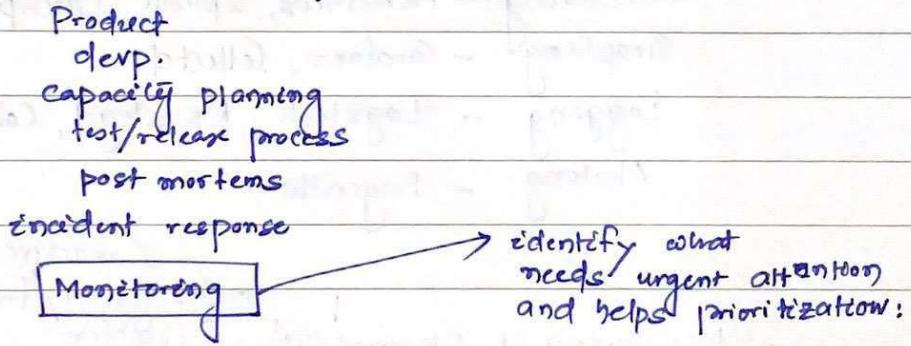
$$SLI = \frac{\text{no. of good events}}{\text{Total no. of events}}$$

may need client-side data coll<sup>10</sup>) along with server-side metrics.  
 (Prometheus)

beyond (consider other sources) → load balancer  
 determined by the app.  
 front-end monitoring/perf infra.

- ★ - SLI → per month
- SLI → tells us whether we are meeting SLO.

### Meeky Dickerson's Hierarchy of service reliability:



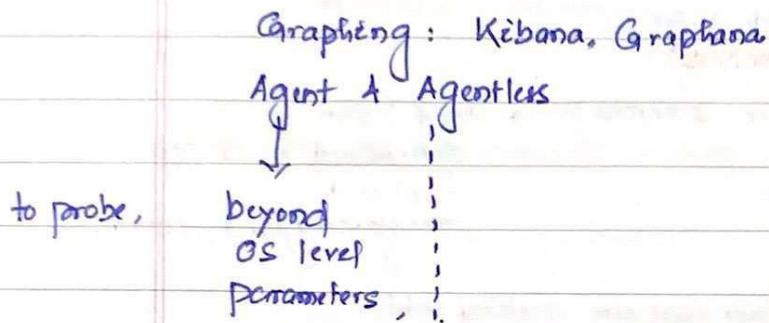
- Telemetry: highly automated process for monitoring.  
 by which measurements are made at remote points and transmitted to receiving equipment

classmate

APM: App<sup>10</sup> Perf. monitoring → strives to detect and diagnose app's perf. problems to maintain an expected level of monitoring and maint. of perf. and availability of S/W apps

Dynatrace      App Dynamics

### - Monitoring Anatomy:

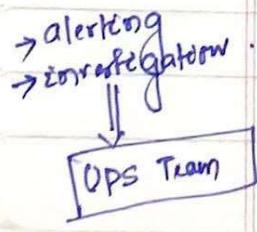
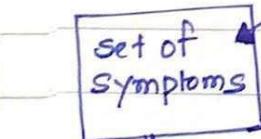


### - SLI Supporting Tools:

- Monitoring - Prometheus, Splunk, Catchpoint, Nagios
- Graphing - Grafana, Collectd
- Logging - Logstash, Rsyslogd, Collectd
- Alerting - Pagerduty

### - Monitoring & Observability:

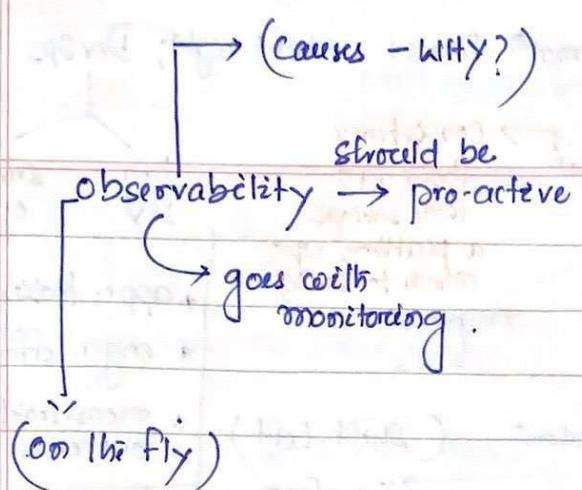
monitoring myth: We can't anticipate the wrong cases with distributed, complex services running at scale



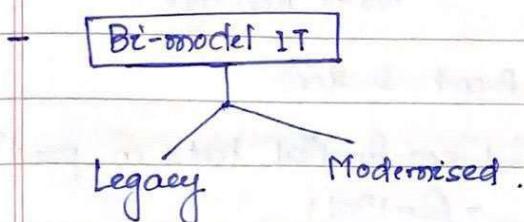
Externalizing all the ops of a service allows us to infer the internal state of that service. (observable)

regular feedback or probing is a must.  
observability

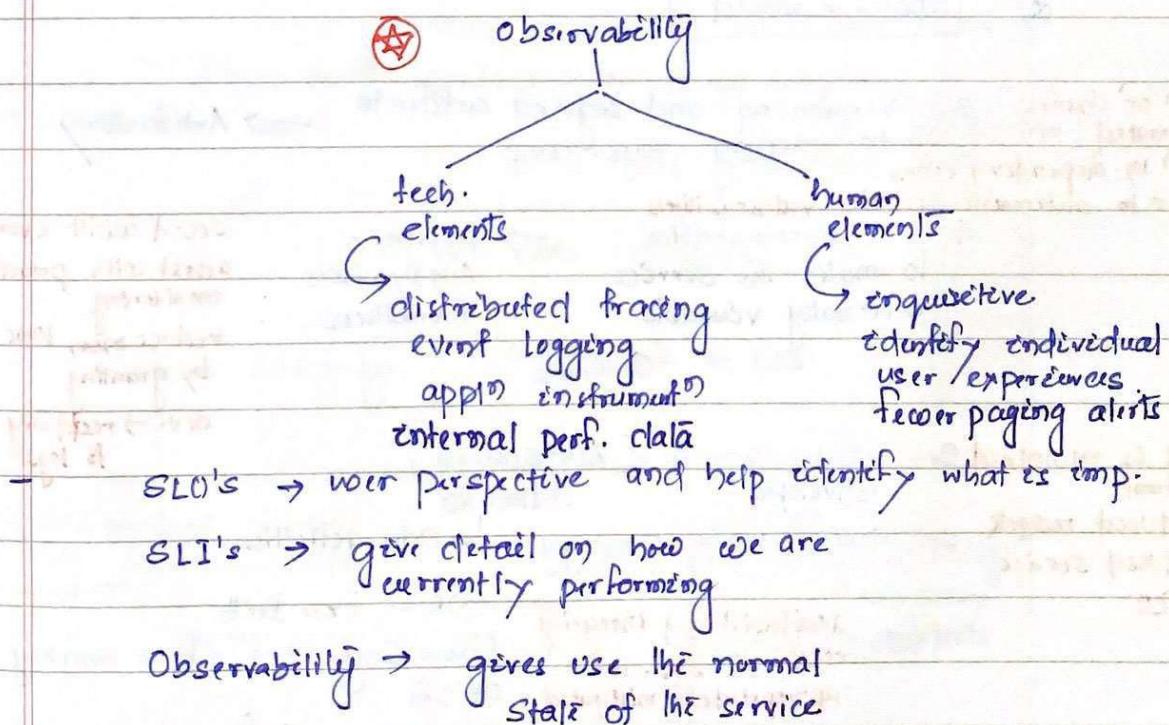
Data collected by monitoring is used in diagnostics/rd.  
 Monitor apps to detect prob. & anomalies.



- Why observability is important?
- 1) rapid rate of service growth
  - 2) dynamic arch.
  - 3) container workload
  - 4) dependences between services
  - 5) customer exp. matters more



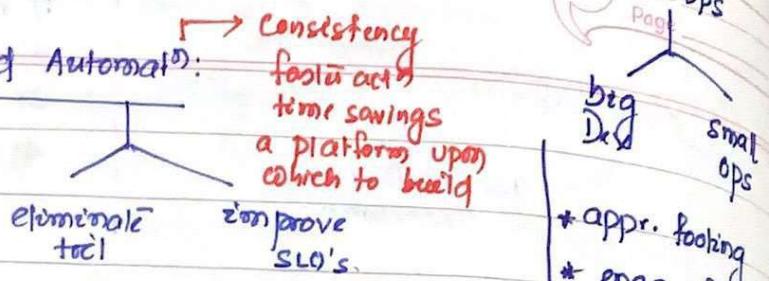
- observability = better alerting
  - ↳ improve signal/noise ratio to focus alerts on key issues



General Automat<sup>ion</sup>) Forces: Shift right, CLASSMATE DevOps Page

### SRE Tools and Automat<sup>ion</sup>:

For SRE, automat<sup>ion</sup> is a force multiplier, not a panacea.



SRE-Led: Service Automat<sup>ion</sup>: (Shift Left)

Env's

1. Infra/Config as code → Terraform, Azure Res. Mgt.  
consistent, repeatable, prod<sup>n</sup> ready  
easy to test and audit changes  
easy to reproduce errors in test env.

Ansible, Puppet, Docker

2. Automated functional & non-functional tests in prod<sup>n</sup>

Systems are well-tested  
but issues can be found in prod<sup>n</sup>.  
extended build pipeline → selenium, cucumber  
functional test

non-functional → JMeter,  
SoapUI,  
Nagios



Spotify model

- prod<sup>n</sup> of change automated  
redn in dependency errors  
easier to determine security vulnerabilities
3. Versioning and signed artifacts to deploy sys. comp. → Artifactory

4. Instrumentation to make the service externally viewable

→ Dynatrace  
AppDynamics  
Prometheus.

sec. & audit events  
assist with proactive monitoring  
reduce mean time to by granting  
dev. → read only access to logs

SLI → OpsGenie

5. future-growt<sup>h</sup> (Auto-scaling)  
Enviro<sup>l</sup> databases

DR - Fire Drills

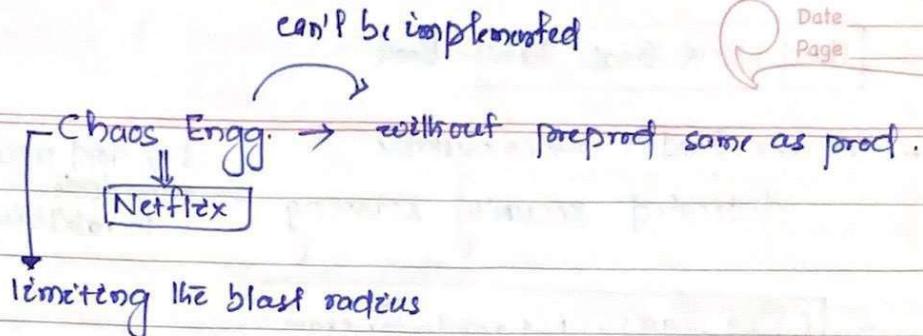
6. anti-fragile strategy: ↴ adm. activities

Chaos Engg. — Chaos monkey  
on call mechanism

availability & integrity  
risks to sys. are appropriately mitigated.  
mitig<sup>n</sup>'s evidenced and tested reducing risk

— PagerDuty,  
Squadcast

toil is minimized  
upfront  
reduced rework  
reduced service TCO



- Hierarchy of Autom<sup>n</sup> Types :

1. None
2. Ext. maintained sys. specific
3. Ext. maintained generic autom<sup>n</sup>

4. Internally maintained system-specific autom<sup>n</sup>      5. Systems that require no intervention

- Secure Autom<sup>n</sup>: (DevSecOps) → must for SRE.
- removes the chance of human error (or well feel sabotage) → build-test-deploy pipeline.

- Pipeline can be validated and checked for compliance
- SRE places extra emphasis on security of prod<sup>n</sup>.

- Secure Build: digitally signing build artifacts avoids possibility of "fake" code.

- Secure Test:

(does not change)  
Test envs are immutable and are created, guaranteeing compliance with code repository.

- Secure Staging: → (immutable)

data security considerations : GDPR, PCI

dedicated to find and uncover vulnerabilities  
Security scanning → or

Proxy security req.

④ dependencies and integr<sup>n</sup> to other services may introduce vulnerabilities.

## Hard Work Beats Talent - Book

- Secure Prod<sup>n</sup>: → immutable  
dedicated security scanning → try and uncover ~~mitigate~~ vulnerabilities.
- (172-195): Not reqd. for exam

### Anti-Fragility and Learning from Failure

Fire-fighting:  
requires  
lot of training  
→ Plan for the unexpected in training

Embracing failure  
will help improve  
MTTD and MTTR metrics.  
↑ optimize monitoring  
mean time to detect (Failure/incidents)

### The Third Way:

\* 1) continual experiment & learning

2) understanding that repetition and practice is the prerequisite to mastery.

introduce faults into the sys. to increase resilience

\* safe experiment and innovation (Hackathons)

\* Pathological  
(Power-oriented)

failure is  
covered up.

Bureaucratic  
(Rule-oriented)

orgn is just  
and sacrificed.

Generative org.  
(Perf. oriented)

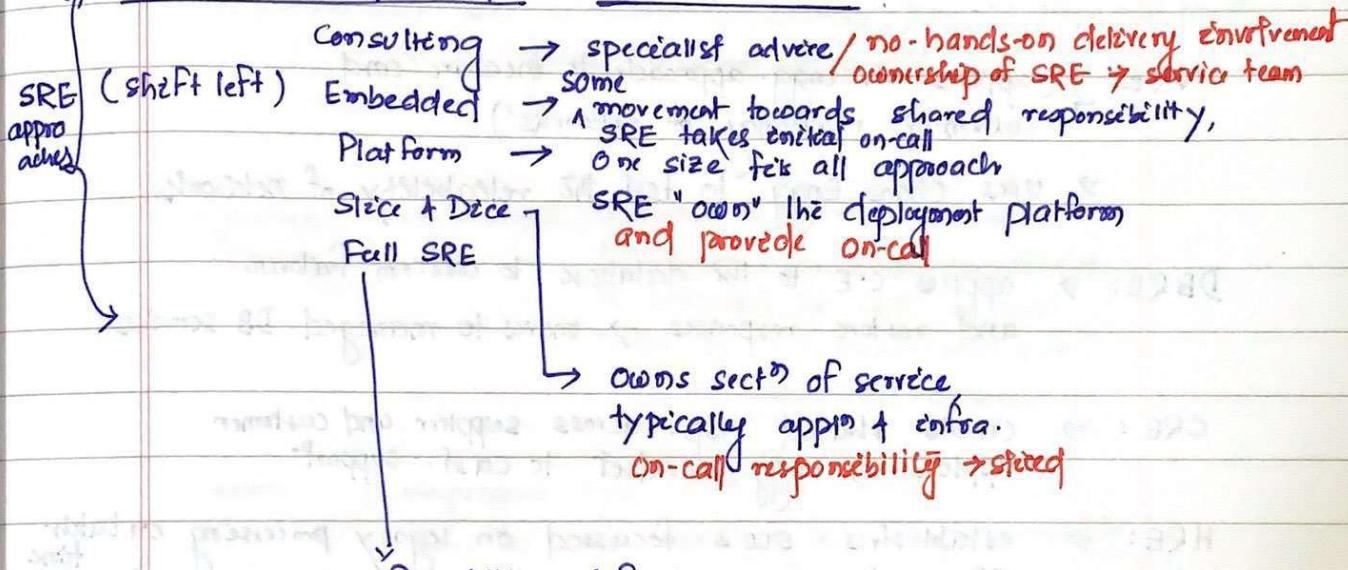
Learning (Failure  
Cause enquiry)

Introduce Fire Drills → annual data center failover test Date \_\_\_\_\_  
Page \_\_\_\_\_  
 ↓  
 loss of facility / techn. / resources / critical 3rd-party vendors  
 or region

e.g. US Army → Chaos Engg. → discipline of experim.  
 ↓  
 minimize the blast radius  
 → outage impact → very minim.  
 e.g. US Army → Chaos Engg. → discipline of experim.  
 ↓  
 minimize the blast radius  
 → outage impact → very minim.

④ minimize impact on customer exp. ⑤ automated self-healing due to  
 fails over-call incidents

Typical SRE Adoption Stages: Org<sup>n</sup> Impact of SRE



Google → recommends → 25% on-call.

Blameless Post-mortem: → P1 or P2 incidents → Self-healing  
Reasons for Blameless PM: monitoring failure, data loss of some kind, etc.  
 SRE & other frameworks      on-call engineer contention

SRE/Agile ← SRE does not stand alone: To do → Backlog  
 def<sup>n</sup> of do → more prod<sup>n</sup> focused.

SRE/DevOps → Pipelines of delivery go further | custom<sup>n</sup>: more widespread & consistent

Agile: Biz. + del.

DevOps: CICD

SRE: Stability + Reliability

ITSM: org<sup>n</sup> learning across value stream

\* SRE can help with ITSM Compliance activities through automation + learning  
+  
Failure → learning opp.      cont. learning

\* ITSM Process models support SRE.  
Trends in SRE:

- 1) Failure as the new normal
- 2) Automat<sup>n</sup> as a service
- 3) Cloud is King
- 4) observe and Learn
- 5) evolution of Network Engineer

NRE: (applies an engg. approach to measure and  
automate reliability of networks)

→ uses Chaos Engg. to test the reliability of networks

DBRE: → applies C-E to the database to confirm failover  
and restore responses → moves to managed DB services

CRE: → creates shared resp. across supplier and customer  
applies an engg. mindset to const. support

HRE: → establishes SLO's focussed on legacy processing by  
R2-platforms onto modern arch.

\* On-call must be undertaken by Ops and engg. teams.  
\* Incident response must be 'sustainable'.

\* Change  
config  
release  
incident-1  
problem

{ ongoing. (ITIL processes) → SREs are involved.

\* Providing 24/7 support req. of SREs  
\* General Autom<sup>n</sup> Focus:  
1) Features are being pushed to those supporting prod<sup>n</sup> in ever increasing numbers.  
2) developers assume that envs are consistent.

3) Testing step introduce false confidence as prod<sup>0</sup> is always unique.

4) monitoring + alerting → focused on things that are known to go wrong.

\* In SRE-Led service autom<sup>0</sup>, overall there is focus on prod. DevOps gains "wisdom of prod". SRE can say "No".

\* Synthetic monitoring: ability to monitor service behaviour by creating scripts to simulate the acts or paths taken by a customer/end-user and the associated outcome.

\* C.E:

- 1) segregate sys. into key comp.
- 2) Test the sys. without key comp. available
- 3) Break the sys. (Non-prod env.s first)
- 4) introduce failure of key components (in-prod)
- 5) introduce DB failure
- 6) introduce total sys. failure in prod

Next Steps:

- 1) Holistic logging
- 2) Identify dependencies
- 3) Improve by error-handling and recovery
- 4) Learn from real failures

Automated Recovery :

1) Create immutable IAC

2) System state + functional code with automated tests

↓  
do not allow SST in prod<sup>0</sup>

self-healing capabilities

Kubernetes / AWS auto-scaling

- detect unresponsive instances of servers or containers
- automatically scale infra/apps up and down
- integrate with monitoring services.

- True North vision: impossible to achieve 100% uptime

- Blameless postmortems: making sure that the same issue does not repeat.

those involved in the failure can give a detailed acc't. of ... without punishment.

- Postmortem Outputs:

details of the incident

follow-up act'

lessons learnt

a timeline of what happened

any supporting info

- key success factors for SRE Adoption:

- i) Exec support and buy-in
- ii) Funding
- iii) Good relationship across spectrum
- iv) Organizing org.

- SRE approach to platform growth: Cont'd + cloud + Net + AS-a-service

- SRE approach to scope growth: Platform SRE  
shift-left  
automating tool

- SRE approach to ticket growth: Tool red<sup>0</sup> + DRY  
(auto-ticket response)  
(clean repeat yourself)  
↓  
so<sup>n</sup> to eliminate repet<sup>n</sup> of prob.

## Chaos Engineering

- 2017: 01 hr of downtime  $\xrightarrow{\text{cost over}} \$100,000$  to biz.
- Rise of microservices & distributed cloud arch. Web has grown increasingly complex.

Companies need a softw.  
and want to avoid next costly outage.

↓  
**Chaos Engg.**

- Chaos Engg: disciplined approach to identify failures before they become outages.
  - \* By proactively testing how a sys. responds under stress, you can identify and fix failures before they end up in the news.
  - \* lets you compare what you think will happen to what actually happens in your systems.
  - \* 'Break things on purpose'  $\xrightarrow{\text{to learn}}$  how to build more resilient sys.

- 2010: Netflix Engg. Tools Team  $\xrightarrow{\text{created}}$  Chaos Monkey

$\Downarrow$   
created to move

on-prem infra.  $\longrightarrow$  cloud infra (AWS)

and the need to be sure that loss of an Amazon instance wouldn't affect the Netflix streaming experience.

2011: Siimonian Army  $\xrightarrow{\text{added}}$  additional failure inj<sup>n</sup> modes  
on top of Chaos Monkey  
 $\Downarrow$   
design a cloud infra. where  
individual components can fail  
without affecting the  
availability of the entire sys.

2012: Netflix shared  
the source code for Chaos Monkey on GitHub.

Failure Inj<sup>n</sup> Testing (FIT)  $\xrightarrow{\text{based on Siimonian Army}}$   
but provided developers more  
granular control over the  
blast radius of their failure inj<sup>n</sup>.

2016: Gremlin: (Founded) by Kolton Andrus

2020: AWS  $\xrightarrow[\text{adds}]{\text{C.E}}$  AWS Well-architected Framework (WAF)

AWS  $\xrightarrow{\text{C.E}}$  Fault Inject<sup>n</sup> Simulator (FIS)

$\Downarrow$   
fully managed service  
for natively running chaos  
exp. on AWS services

2021: Gremlin: publishes 1st ever State of Chaos Engg. report

$\Downarrow$   
growth of C.E  
benefits

frequency of running chaos exp.

C.E runs thoughtful + planned experiments

### Principles:

that teach us how our sys. behave in the face of failure.

These exp.s follow 03 steps.

1. Plan an experiment : Create a hypothesis. What could go wrong?
2. Contain the blast radius : Execute the smallest test that will teach you something.
3. Scale or squash : Find an issue? Job well done. Otherwise increase the blast radius until you are at full scale.

- Companies: National Australia Bank

+

On-prem ... → ... cloud

↓

reduce incident counts

- C.E is a tool to build immunity in our tech. sys. by injecting harm like

latency  
CPU failure or  
network black hole

to find and mitigate potential weaknesses.

Breaking things on purpose → enable detect unknown issues that could impact our sys. + cust.

is  
di

2021 State of C·E Report: Most common outcomes of C·E

increased availability ( $> 99.9\%$ )

low MTTR

low MTTD

fewer bugs shipped to product  
fewer outages

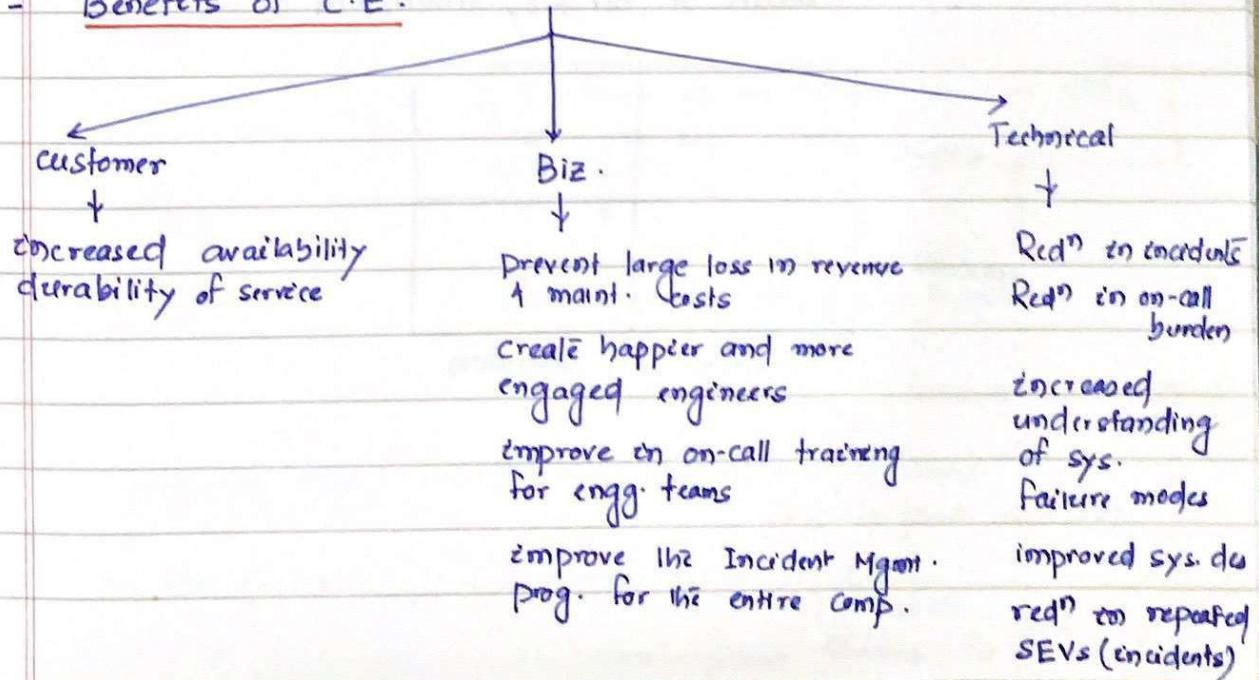
\* C·Exp are more likely to have  $> 99.9\%$  availability.

### - 8 Fallacies of distributed sys.:

- 1) Network - reliable
- 2) Latency - 0
- 3) Bandwidth - infinite
- 4) Network - secure
- 5) Topology - doesn't change
- 6) OI Admin
- 7) Transport cost - zero
- 8) Network - homogenous

These fallacies  
drive the  
design of C·E .

### - Benefits of C·E:

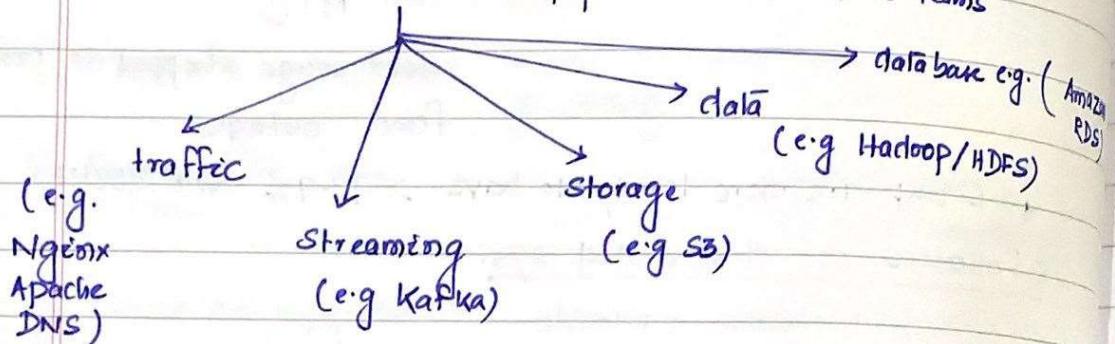


### - C·E For Service Teams:

dedicated Chaos Engg. Teams → small team



Awareness and practice C·E across teams



### - C·E Experiments:

Experiments should be performed in the following order.

- i) Known Knowns → things you are aware of and understand
- ii) Known Unknowns → things you are aware of but don't understand
- iii) Unknown Knowns → things you understand but are not aware of
- iv) Unknown unknowns → things you are neither aware of nor fully understand

↓  
things you are neither aware of nor fully understand

Known	
unknown	
Known	unknown

- Planning for First C-Exp: 'What could go wrong?'

Services  
dependencies  
data stores  
env.s

Create a hypotheses. → chosen the exact failure to inject

+ What happens next?

discuss the scenario and  
hypothesize on the expected  
outcome before running it live.

measure the impact → sys. availability ↑  
durability ↓

KPI → relate customer success  
such as

orders/min

stream starts/sec

If there is an impact to these KPIs,  
halt the exp. immediately.

measure the failure where  
you want to verify your hypothesis.

impact on latency, req/sec or sys. resources.

RollBack Plan: → Backup plan (but sometimes  
backup plan can fail)

Go for it: → Running your 1st C-Exp

(2 outcomes)  
Sys.-resilient Problem (fix it)

## Reducing MTTD for High-severity incidents:

- C-E: disciplined approach to identify failures before they become outages.

Using thoughtful experiments, you 'break things on purpose' to learn how to prevent outages and build more resilient sys.

### C.Experiments

| 4 steps

↓  
define the sys. normal behavior - its steady state  
- based on measurable o/p like overall throughput, error rates, latency and so on.

↓  
hypothesize about the steady state behavior of an experimental group, as compared to a stable control group.

↓  
expose the experimental group to simulated real-world events such as server crashes, malformed responses or traffic spikes .

↓  
test the hypothesis by comparing the steady state of the control group and the experimental group.

- Gremplin's Failure-as-a-service (FAAS) platform was built by Amazon and Netflix Engineers.

- High-severity incident (SEV) used → Amazon, Dropbox, Gremplin

availability drops  
product feature issues  
data loss  
revenue loss  
security risks

Q:  
d:

When  
TTD - Time to def<sup>n</sup> → interval from an incident starts to the time it was assigned to a tech lead on call who is able to start working on resolution or mitigation.

- \* SEV 0s (highest and worst kind of severity)
  - ↓ industry-recommended (able to detect and resolve within 15 mins)
- ⇒ SREs has to improve the MTTD.

- \* High-impact methods to reduce MTTD:

Incident classification and levelling

Tooling

Monitoring

KPIs

Alerting

Observability

C-E

Step 0: Incident classification, including SEV descriptions and levels, SEV timeline and TTD timeline

Step 1: Critical-service monitoring including key dashboards and KPI metric emails

Step 2: Service ownership and metrics, including measuring TTD by service, service triage, service alerting

Step 3: On-call principles, including the pareto principle, rotation structure, alert threshold maintenance and escalation practices

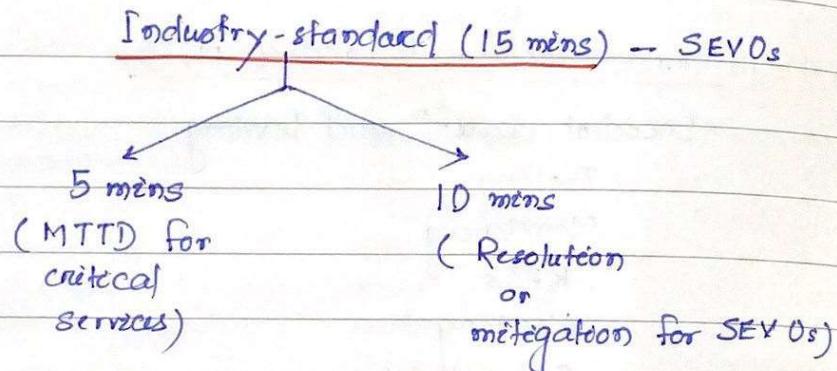
Step 4: Chaos engg., including chaos days and continuous chaos

Step 5: self-healing sys., including when automation incidents occur, monitoring and metrics for self-healing sys. autom.

Step 6: listening to your people and creating a high-reliability culture

#### STEP 0: Incident classification:

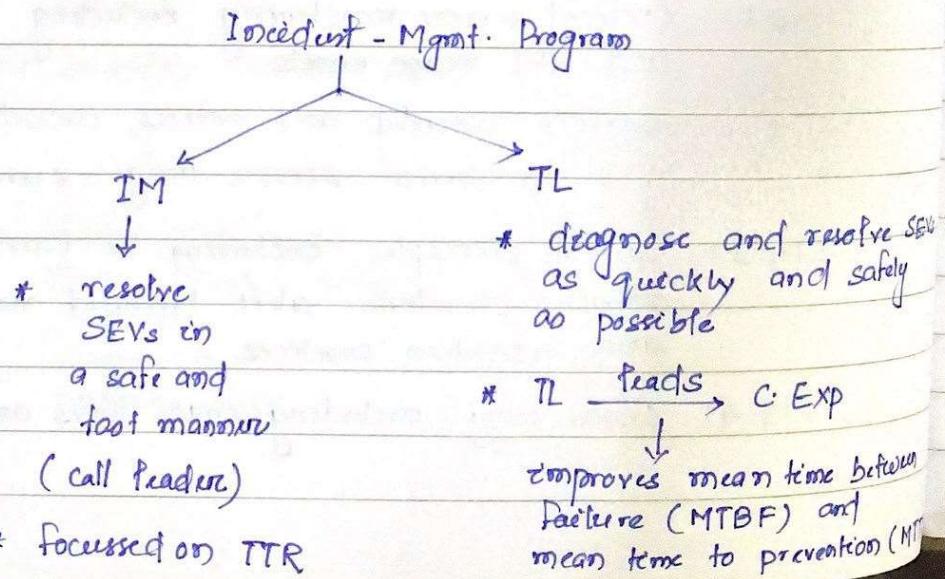
SREs → determine the SEV level of an incident and then triage and route the incident to an appropriate TL.



SEV 0 → catastrophic service impact

SEV 1 → critical device impact

SEV 2 → high service impact

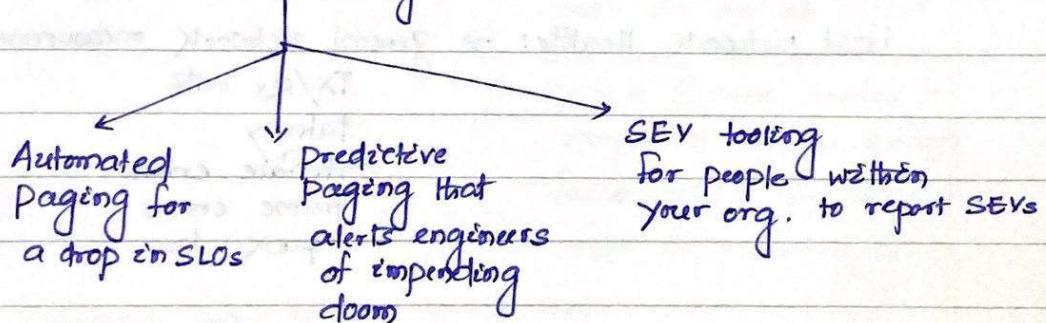


\* One primary IM  
One secondary IM } on-call per week

- TTD: Time to detect  $\Rightarrow$  time interval from when the incident started to the time it was acknowledged by an engg.

TTR: time to resolution (mitigation)  $\Rightarrow$  introduce fix and mitigate impact of fix

- SEV Automation Tooling:



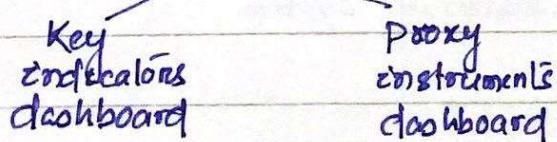
f) Step 1: Org-wide critical service monitoring:

Monitoring  $\rightarrow$  reduce MTID for high-severity incidents

Critical-service dashboard  $\xrightarrow{\text{create}}$  Datadog & Grafana

Google's SRE's 'Four Golden Signals'  
Tom Wilkie's RED Method  $\xrightarrow{\text{good baseline standards for every critical service must have.}}$

High-level service dashboard:  $\xrightarrow{\text{quickly drill down into breadth of the services}}$



System services health: → ancillary services which are reqd. for the infra. to be considered healthy.

↓  
SSH

DNS

Puppet  
LDAP

Host General Health:

↓

cpu time

CPU load

memory usage

I/O latency

Host Network Health: → general network measurements  
TX/RX rates  
latency  
interface errors  
frame errors  
packet loss

- KPI service excel: contain rolled-up daily metrics for each critical service

## ④ Step 2: Service Ownership and Metrics

increasing microservices → managing services has grown increasingly complex.

Service ownership → adv. SRE practice.

RACI (matrix) → Responsible, Accountable, Consulted and Informed model

SOS → Service ownership system

### ③ Step 3 - On-call principles

SREs → new guardians  
+  
review alert architecture once per quarter

### ④ Step 4 - Chaos Engineering

↓  
series of thoughtfully planned experiments  
designed to reveal the weakness in our systems.

Insights from chaos exp. → red<sup>n</sup> in incidents  
red<sup>n</sup> in on-call burden  
increased understanding of  
system failure modes  
improved system design  
faster MTTR for SREs and  
red<sup>n</sup> in repeated SREs.

### ⑤ Effective way to start practicing C.E → Run a Chaos Day

### ⑥ Step 5: Detecting Incidents Caused by Self-Healing Systems:

↓  
e.g.  
Restarting  
a stopped  
service

↓  
dangerous and can cause  
cascading failures and DDoS or  
without adequate testing.

Real-time e.g.: Facebook's F8AR  
LinkedIn's Morse  
Netflix's Winton

Self-healing weaknesses: fast time of detection  
and should not  
be firing regularly

\* Steps: Listening to your people and creating a high-reliability culture

- ▶ Preoccupation with failure
- reluctance to simplify interpretations
- Sensitivity to opera's
- Commitment to resilience
- deference to expertise

### Testing Disaster Recovery with CE:

- DR is the subset of BCP.
  - ↓ focusses on restoring the essential functions of an orgn. as determined by a BIA (Biz. impact analysis)
- DRPs (disaster recovery plans) → playbooks or run books

Black swans: outlier events that are severe in impact.

concerns that trigger problems that you didn't know you had, cause major visible impact and can't be fixed quickly.

by def'n, can't be predicted.

e.g. natural disaster

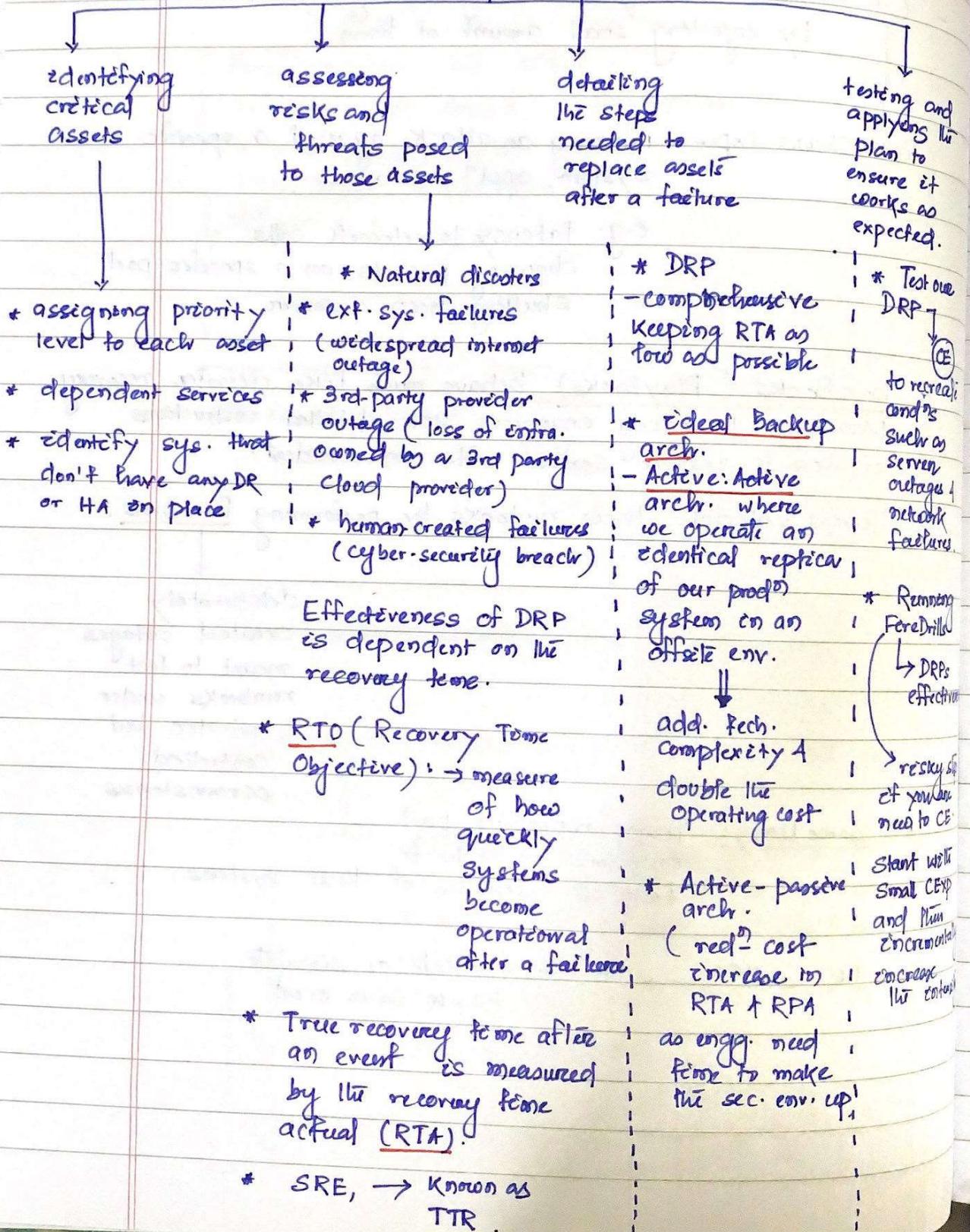
- C.E → practice of systematically testing sys. for failure by injecting small amount of harm.

↓  
Run Chaos exp → running an attack against a specific system, appln or service.

e.g. latency to network calls  
dropping packets on a specific port  
shutting down a server

- \* Run Books (Playbooks) behave much like disaster recovery plans by providing engineers with detailed instructions on how to restore service after an incident.
- \* Teams validate their runbooks by performing Fire Drills.  
↓  
deliberately created outages meant to test runbooks under realistic but controlled circumstances.
- \* Game Days: Teams reproduce past incidents in order to test the resilience of their systems.
- \* DR: use → C.E to recreate or simulate Black swan event

## \* Create and validate Disaster Recovery Plan



\* A real incident is a true test and the best way to understand if something works. However, a controlled testing strategy is much more comfortable and provides an opportunity to identify gaps and improve.

\* In an ideal scenario, RTO & RTA will be zero.

\* RPO (Recovery Point Objective) → amount of data loss or data corruption that we are willing to tolerate due to an event.

$RPO = \text{Last Backup Time}$

- Time of failure (at which we assume that any data stored on the failed system is lost)

\* Actual amount of data loss is measured by the recovery point objective (RPA).

\* GetLab:

↓  
accidentally deleted data

from their prod DB.

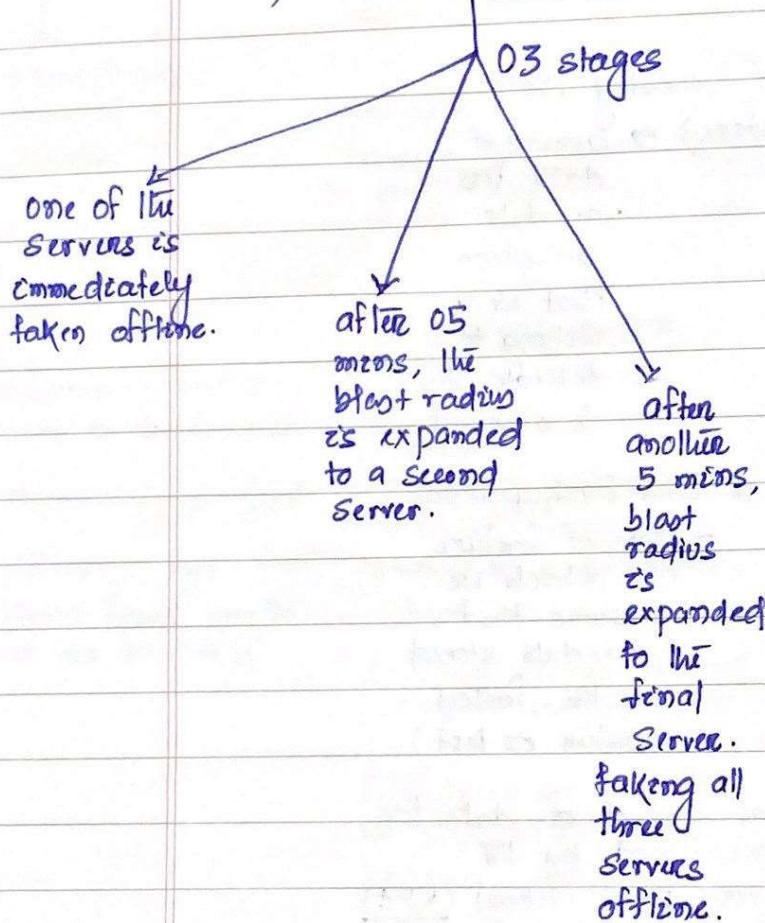
6-hour outage  
(RTA)

5000 proj.  
5000 comments  
700 user acts.  
(RPA)

\* Red<sup>o</sup> in RPO → more freq. backups or replicating continuously  
increasing operating cost and complexity ← to an offsite loc<sup>o</sup>

- Validating the plan by running a Fire Drill:

i) create a scenario in Grenition



(Scenarios let us accurately simulate real-world disaster conditions by incrementally increasing both the blast radius and the magnitude of an attack.)

e.g. Simulate losing a DC or region consisting of 03 servers

#### Black Hole Attack:

Block all incoming and outgoing network traffic on each server for the duration of 100 seconds.

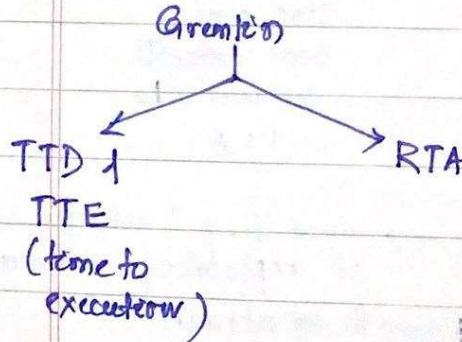
#### Shutdown Attack:

(more realistic simulation)

The Shutdown Grenition issues a sys-call to shut down or reboot the OS on which the target is running on.

\* Once 30 mins has passed, the scenario ends and all 03 servers resume normal service.

\* During the scenario, track scenario start time and end time.



- Post-mortem: analyses of key metrics (RTA + RPA)
- \* DRP - living document that requires frequent updating and testing

### How to Install and Use Gremlio locally with Docker Desktop

- 1) Install Docker Desktop
- 2) Verify Kubernetes cluster is up and running  
    ↓  
    Kubectl get nodes
- 3) Add the Gremlio Helm repository.
- 4) Create a Gremlio namespace  
    ↓  
    Kubectl create namespace gremlio
- 5) Find your Gremlio Team ID and Team Secret
- 6) Install the Gremlio Helm chart
- 7) Verify that your clusters are now up and running  
    ↓  
    Kubectl get pods -n gremlio
- 8) Run your first Gremlio scenario with Docker Desktop.
- 9) Run a Gremlio scenario to validate reliability

## Gremtian Scenarios Walk Through:

- \* Gremtian performs CE → through SaaS based platform.
- \* Scenarios → Gremtian feature that allows you to map real world outage scenarios to CE exp.
- \* Pre-requisites:
  - i) Gremtian account
  - ii) AWS Elastic Kubernetes Service (EKS) cluster
- \* 1) Create and run a custom scenario  
 2) Run a recommended scenario

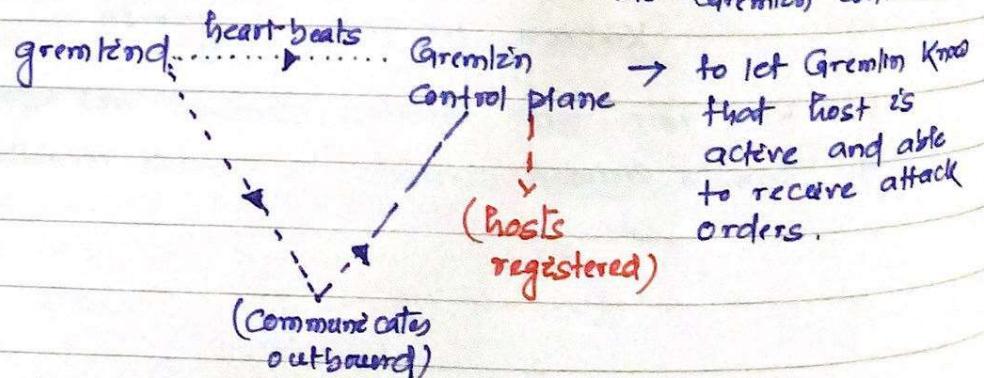


allows users to define their own scenarios.

map real-world scenarios to CE attacks.

## Daemon & Client Overview:

- \* Architecture (Client <--> Daemon → Control Plane)
- \* Gremtian Daemon (gremlond) → binary installed on the OS or available inside the Gremtian container.



- \* Gremtian client (gremten) → Gremtian CLI that is responsible for creating the local impact with the host.

- \* gremlind + gremlin  
(daemon)      (client) → this unit is referred to as a targetable client to the platform.

\* Client Lifecycle:

- Gremlin clients (infrastructure and app<sup>n</sup>) that have been authenticated to the Gremlin Control Plane appear in the infra. clients and app<sup>n</sup> clients lists.
- You can run attacks on active clients.
- Idle client: if there is no activity for the past 5 mins.
- If Gremlin doesn't hear from these idle clients for a period of 24 hours, clients are removed from the list.
- If a client starts communicating with Gremlin again while still within the 24 hour idle window, the client is reactivated and returned to the active state.

- \* Logs → /var/log/gremlin  
daemon log entries → daemon.log file
- ↓
- log entries in this file may indicate events where the daemon is not able to communicate with the control plane and eventually trigger the Dead Man Switch.

Each attack on the host logged /var/log/gremlin/executions

↓  
unique  
attack  
execution ID

\* `du -sh /var/log/gremlin` → log size

### Bandwidth usage

Idle state



daemon uses very  
little bandwidth.

Attack state



Slight increase  
in overall bandwidth  
consumptn during  
attacks.

The bandwidth used  
is not affected  
by the type of  
attack being run.

- \* When Services discovery is enabled, the Gremlin daemon will send additional data and the bandwidth consumed will depend on how many services are discovered.

Measurement of actual bandwidth consumed by Gremlin for particular install

use a tool → iptraf or netlog

### Command Line Interface:

- \* `attack` → starts a localized gremlin attack

e.g. gremlin attack cpu -f 60 -c 1 → attack on  
cpu-1 completed

- \* `attack-container` → Starts an attack against a locally hosted Docker Container

→ In order to run attack-container, the container must reside on the host and

gremlein attack-container a1a9ee7eb256 -1 60 -c1

- \* check → inspects the state of your Gremlein install

e.g. gremlein check os → o/p is not sent to Gremlein. Helpful to support team if we need to triage an issue with install.

- \* **enct** → sets up the client config.  
A prompt will appear to enter the Team ID  
Team Secret  
→ registers client to control plane.

e.g. gremlein enct

- \* logout → removes the Gremlein Client from the Gremlein Control Plane.

- \* rollback → interrupts an active attack. It can also revert the last impact.

e.g. gremlein rollback

- \* rollback-container → interrupts an active attack against a local Docker container.

- \* status → displays the status of any attacks.

e.g. gremlein status

gremlein status -v <GUID> → display the status of a specific attack

- \* version → state the installed version of Gremlein

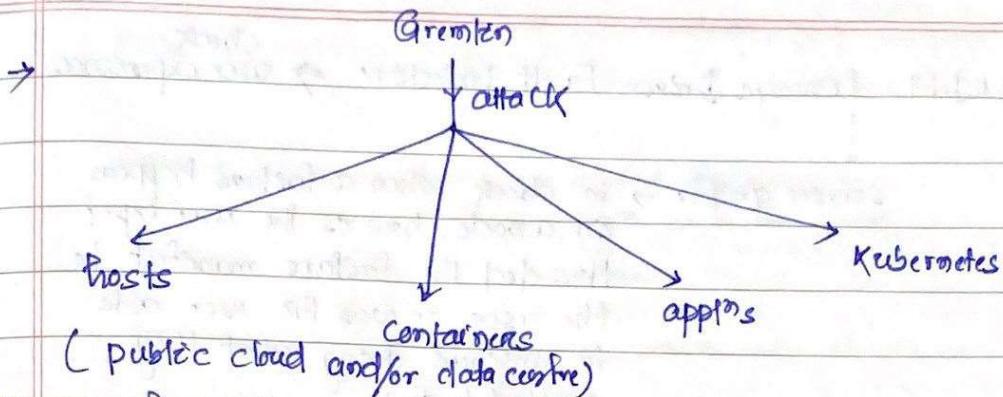
must be specified.

e.g. gremlein version

- C.E: Build  
Fail  
Learn  
Repeat → (BFLR)
- C.E → to better understand our systems.  
 ↓  
 ↗ Shovel  
 Prepare engineers for prodn.
- critical role in C.E
- Tech.  
People ↑  
Process
- C.E → uncover new info  
 ↓  
 Buy-in from Leadership
- Reliability → Technical Debt



- Reliability - journey



Attacks at infrastructure category : { Resource - CPU, disk, I/O, memory  
state - process killer, ShutDown,  
network time travel

Payer

Black Hole  
DNS  
Latency  
Packet Loss

Attacks at appn Payer  $\Rightarrow$  using Gremion's alpha libraries

attacks can be configured webapp

API  
CLI

- Call Leader: Person responsible for managing an incident.

$\hookrightarrow$  (judgement call)

- C.E  $\rightarrow$  cozy to do runtime validation.

- FIT (Failure Injection Testing)  $\xrightarrow{\text{root}}$  developed by Netflix in 2014

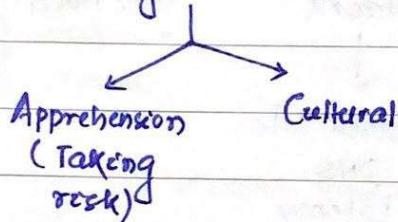
$\downarrow$   
allows an engg.  
to add a failure scenario  
to the request header of  
a class of reqs at the  
edge of our service.

- LDFI: Lineage Driven Fault Injection  $\Rightarrow$  check user experience

↓

Service graph  $\rightarrow$  to check when a failure happens in a node, how is the user exp.? How did the failure manifest to the user or was the user able to continue doing what they wanted to do?

- Challenges in CE:



- C.E : History, Principles and Practice :

Every brief outage can impact a company's bottom line, so the cost of downtime is becoming a KPI for many engg. fun.

CEO of British Airways  $\rightarrow$  tens of thousands of stranded passengers in May 2017

$\downarrow$   
cost the company 80 million pounds

C.E  $\rightarrow$  preventive medicine

$\downarrow$   
disciplined approach to identify failures before they become outages.

C.E  $\xrightarrow{\text{literally}}$  breaking things on purpose

## CLASSMATE

migration from phy. infra to cloud  
Page 10  
and the loss of an  
Amazon instance shouldn't  
affect  
the  
Netflix  
Streaming  
experience.

2010 → Netflix Engg. Tools  $\xrightarrow{\text{created}}$  Chaos Monkey.

2011 → Simeon Army → added failure  
engg modes on top of  
Chaos Monkey

2012 → Netflix shared the source code for Chaos  
Monkey on GitHub.

2014 → Netflix created new role called Chaos Engineer.  
FIT (new tool based on the concepts of Simeon Army)

2016 : Gremlins

2020: AWS adds CE to the reliability pillars of the  
AWS (Reliable-Architected Framework) - WAF.

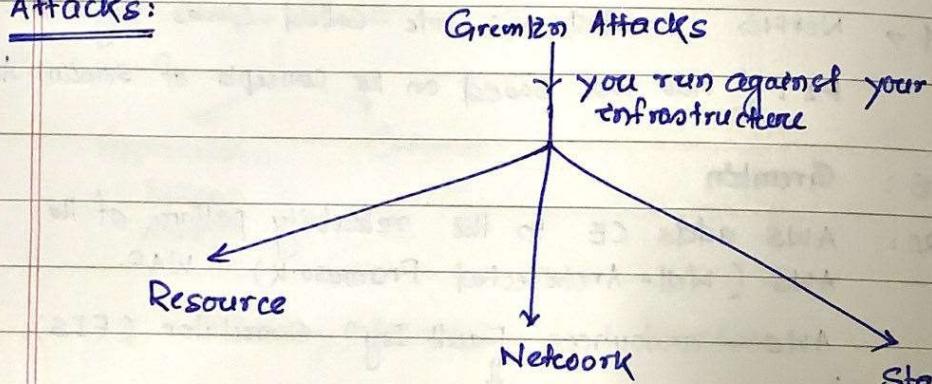
AWS - announces Fault Inj<sup>n</sup> Simulator (FIS)  
 $\Downarrow$

running C. Exp on AWS services

2021: State of CE Report

- \* Recommended config for an critical (Gremlin) Latency attack (100 ms 60 secs)
- \* Gremlin disk attack
- \* Game Day participants and roles
- \* What envs do you conduct Chaos Engg?
- \* CPU attack

### Attacks:



generates high load for one or more CPU cores.

\* CPU

\* memory

\* IO

\* Disk

\* Black Hole

\* Latency

\* Packet Loss

\* DNS

\* Shut Down

\* Time travel

\* Process Killor

Memory - allocates a specific amount of RAM. (malloc)

\* Black Hole  
→ drops all matching network traffic.

Shut down - performs a shut-down on the host OS to test how your sys. behaves when losing one or more cluster machines

IO - generates large amounts of IO req.s to file sys. that are mounted on block devices.

\* Latency  
→ injects latency onto all matching egress network traffic

Time travel - changes the host's sys. time which can be used to stimulate a DST

Disk - fills up a target block device

by writing to

files at the dest<sup>(n)</sup> (dest<sup>(n)</sup> → --dir  
no. of files written → --workers)

### \* Packet Loss

→ includes packet loss  
into all matching  
egress network traffic

### \* DNS: Blocks access to DNS servers

↳ Blockhole attack  
against port 53

### Process Killer:

Kills the specified  
process which can  
be used to simulate  
app or dependency  
crashes.

↳ doesn't work for  
PID ↴  
shutdown  
attack

### Game Days:

- \* are like firedrills - an opportunity to practice a potentially dangerous scenario in a safer environment.
- ↓  
capstone to measure reliability
- \* Running a GameDay tests our company from engagement to incident resolution across team boundaries and job titles.
- \* verifies the sys. at scale, ultimately in Prod<sup>10</sup>.
- \* Planning Out the GameDay is a great opportunity to collaborate with other parts of the company, share context and learn about the sys. as a whole.
- \* Communicate → Command Center
- \* End Goal → to run in Prod<sup>10</sup>
- \* coined by Jesse Robbins (Amazon)