

SRE Responsibilités: availability, latency, perf., efficiency, maintenance
change mgmt., monitoring, emer. response
and cap. planning of services.

SRE PRACTITIONER

OB/13:

SLO → user exp (data)

Tech Redn → Lean

TOC → Theory of constraints (not everything - can't be automated)

→ create think-tank time

SLI → parameters to measure

measurement → 360 degree view (monitoring observability)

"Slow is the new fast"

Anti-Patterns → if we do, it will backfire

1. ReBranding Ops:

* SRE → engineering work

require

* build sys. to less human intervention, to fail less often and modify existing sys. to remove emergent failure modes.

* don't bring ops closer to machines → sys. should take care of it.

(Build right alerts) modify existing sys. to remove emergent failure modes.

Automation of alerts (False positives and (a person has spe. disease or condⁿ) when the person actually doesn't have)

False negative → Blind spot / Trap (sys. thinking knowledge)

(a person does not have a spe. disease or condⁿ when the person actually has)

we think everything is OK.
but it is not discovered.

spend more than half of their time

on building better sys. → cultural shift ownership

rather than

conducting or

documenting operational

tasks.

reduce tech debt

focus on user exp.

→ ExLA (experience level agreement)

Recall: 100% if every sig. event results in an alert.

Det'n Time: How long it takes to send notif. in various cond's. ^{classmate}

Reset Time: How long alerts fire after an issue is resolved.

(Turn your SLOs into actionable alerts that are grounded on the customer's path)

2. Users Notice an issue before you do

Precision → quality of being exact
and refers to how close two or more measurements are to each other

Recall → key attributes (more measurements are to each other)

When it happened → Det'n Time → Reset Time

Prec^{ce} is 100% if every alert corresponds to a significant event.

Build solid fail-over, fault tolerance to resolve within SLO.

sentiment analysis
NLP algo.

SRES → direct contribution to getting the cust. what they want.

3. measuring until my Edge

Customer's percept^{ce} → reality

SLO → user exp.

and not what your SLOs read

being close to biz.
and understand their pain

SLOs → keep on changing.

- * Shared inst^{ns}, shared integ^{rs} and shared comm^{ns} → key to success
- * Meeting SLOs - useless if the cust. doesn't have the exp.

4. False positives are worse than no alerts.

↳ monitoring is abt. the steady flow of traffic, not a steady flow of alerts.

5. config. mgmt trap → SRES + more on homogenizing ecosystem

Snowflake - inefficient

design for immutability infra.

change → about

replacing

never updating or patching

- * use config^{ce} mgmt. to consolidate and integrate to immutable infra.

6. DogPile: Mob Incident response

→ procedural framework is a mandate.

Incident Commander: → just calm down

incident command framework
keeps off - disruptive to engg. work

- * minimize damage, make outage as short as possible

7. Point Fixing: → minimize outage with automated alerts and solved paging which and quick workarounds, fast rollback, fail over and fix forward

Reactive Proactive Predictive auto remediation

→ Short term fix
→ permanent fix

class of
elemental file design errors

→ moving to predictive methods

i) understanding ii) goal iii) historically what we spent our time.

④ Auto remediation

closed loop remediations

without
human
intervn

→ goal of SRE

8. Prod Readiness Gatekeeper:

④ Shift left to
build in resilience
by design in the
devp. life cycle

Incident - wait for concerned person to respond

→ increase the length of the time between the creator of a change and its prodn release (gatekeeper)/chokepoint/speed bump.

Systems → 1st line of defence

Secondly should be in-built.

+ SREs provide dev teams with devops frameworks and templated config's to speed up reviews.

9. Human Error → it's a sys. error rather than a human error.

→ prevent the occurrence of the event. → making the sys. more robust.

Root cause is just the place where we decide we know enough to stop analyzing and trying to learn more.

* Q's on case story:

Momo Bank: defense in depth works for reliability
not just for security.

* Chaos Engg: helps in dist. sys. in order to build confidence to the system's capability to withstand turbulence in prod.

- Good post mortem
- Blameless
- answers to 5-whys
- details on discovery of the inc.
- timelines (chronology)
- ways to prevent incident
- log of related incidents
- lessons learned & clear act items
- well-timed timelines

Bad postmortem

- finger pointing
- animated language
- missing ownership
- missing context
- key details omitted
- limited audience

Precn / Detn (MTTD)

Incident commander

War-room

Observability
Chaos Engg.

Blameless Post-mortem

Model 2

SLD → Proxy for cust. exp.

Establish SLIs in distributed ecosystems?

↳ varies according to personal customer journey map

Goals

between components in an arch., it can mean producing component and consuming comp.

complex system: understanding - limited

current state

enhancing → enhanced state

Key SLI → measurement of the availability of a sys.

← SLD → goals we set for how much availability we expect out of a sys.

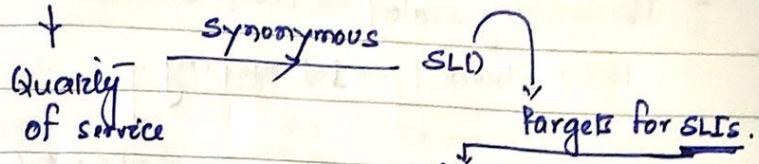
* Setting SLIs and SLOs on system boundaries.

Complex Sys. Traits:

- 1) cascading failures 2) diff. to determine boundaries
- 3) diff. to determine or model behavior
- 4) dynamic network of multiplicity
- 5) emergent phenomena
- 6) relationship - nonlinear
- 7) relationship contains feedback loops.

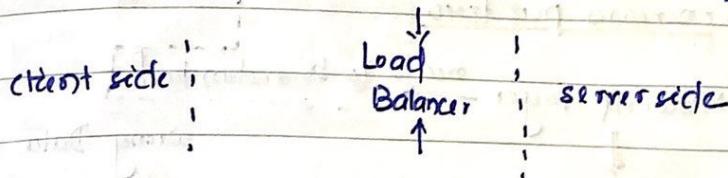
- Relationships: non-linear

- QoS



Kudo's SLO journey:

where the SLO needs to be measured → imp



SLI → identify the options

where SLO carries ~~average~~ weightage
on the overall system

* Elasticsearch HeartBeat → to monitor our internal only APIs and services

System Boundaries: Identify the sys. boundaries within your platform.

→ identify the customer-facing capabilities
that exist at each sys. boundary

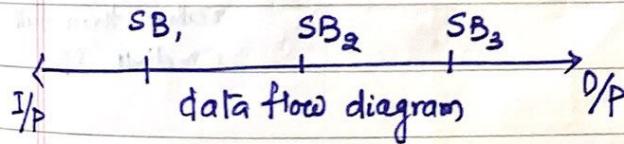
- articulate a defn
of what it means
for each capability to be
available.

Define SLO for
each capability
and start measuring
to get a baseline.

understanding
arch. is a must-
/system

boundaries as a
service owner.

iterate and refine
your system and fine tune
the SLOs over time.



- Horizontally scaled → Load Balancer
- Hard Sharded → 1-1 mapping
- SLO = 0% for a set of users

MTTR → mean time to repair
 ↓
increases per tier

Network Layer → needs to be redundant!

wrong UI → isolated
 wrong service → UI will be impacted but UI caching well reduce impact.
 wrong Data → service + UI fail will go down.

I. System Boundaries (data domain driven app.)

D. cap for sys. boundaries (address cross func issues)

wrong Network → everything goes down

define SLI for each capability in a dis. sys. * MTTR increases with count of impacted Service tiers.

define SLIs targets
 measure baseline → refine & repeat

keep single Service tier')

* Graceful degradation

↓
 significantly decrease the amt. of work or time needed by decreasing the quality of response.

one SLO per capability goal

Large Scale Cultural Shift:

SLO Buy-in

Become a SLO Expert

create standard artifacts for use

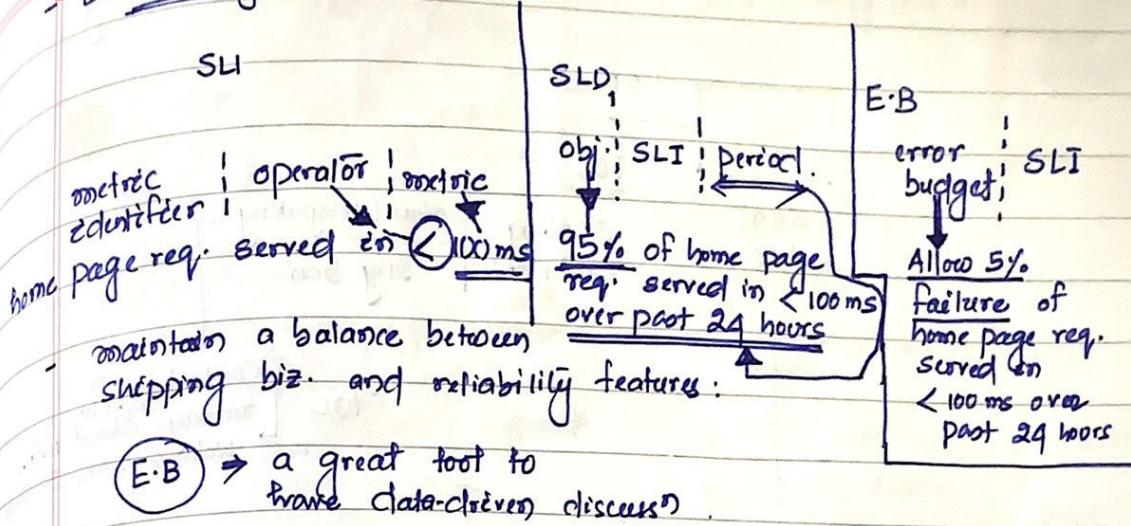
SLO pilot and rollout

communicate (iterate)

e.g.: A search app might only search a subset of data stored in an in-memory cache rather than full on disk DB.

(*) cascading
→ Error budget

Error Budget:



- Burn of E.B. → determine the weakest link in your dist. ecosystem

* use E.B. to improve resilience
* exp. window

E.B. → error-budget → great tool to have data driven discn
Prometheus monitoring → Grafana dashboard

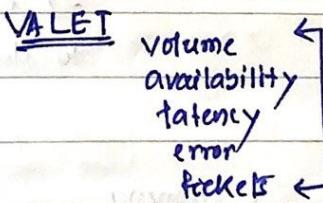
* use Error Budget: turn off heritage components to see impact.

SLO Inflection Point:

* E.B. burndown charts are c/sps for decⁿ making on disaster recovery / datacenter drills.

Home Depot's SLO Journey:

{ good strategy
executive buy-in
easy adoptⁿ patterns
and patience

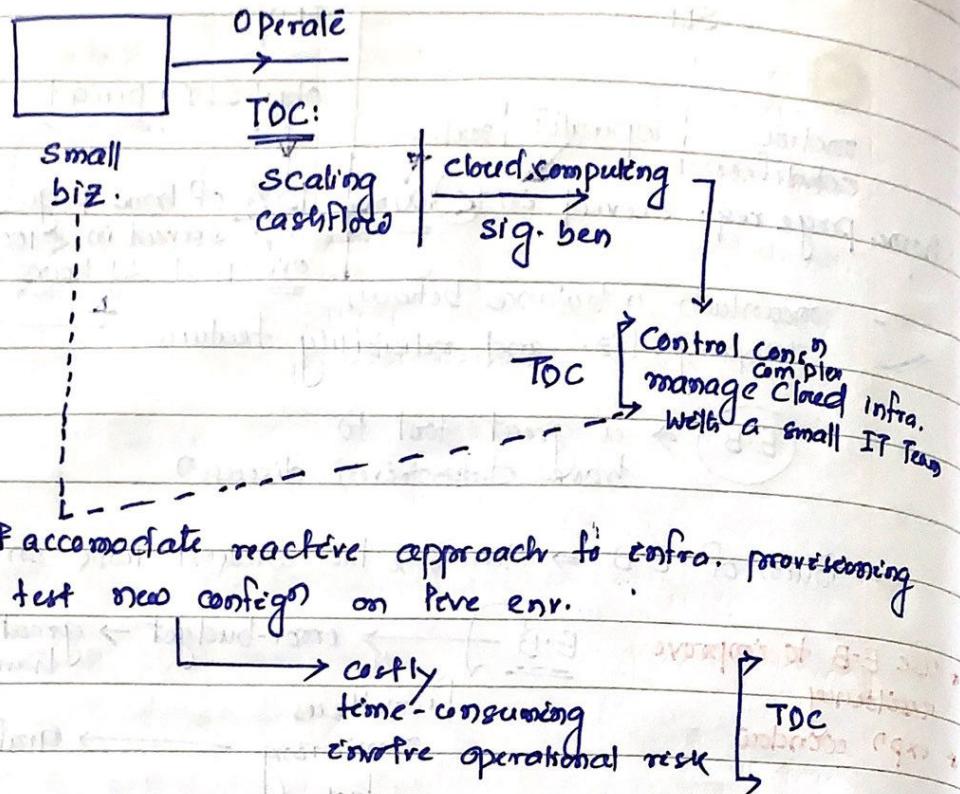


* SLO culture: bring prod. dev. & ops close together

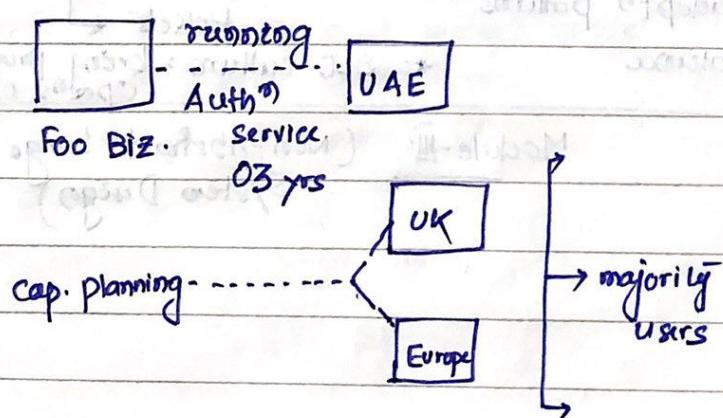
Module-III (Non-Abstract Large System Design)

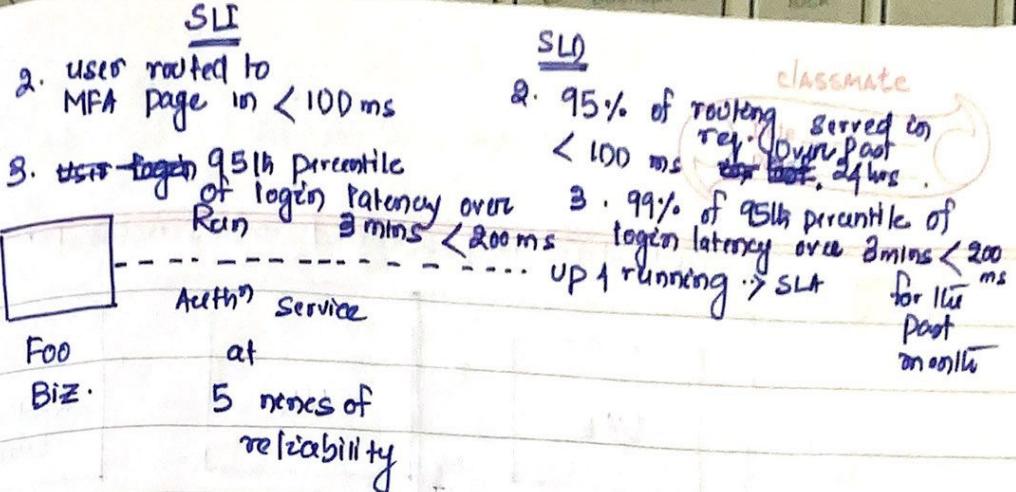
CLASSMATE Date _____ Page _____

Non-Abstract Large Scale Design - Capacity planning



Simulation → eff. and efficient evaluatⁿ and expⁿ of cloud infra. configⁿ and provisioning/de-provisioning of capacity by examining how infra. reacts to demand in the short/med. term.





Current Fact

concurrent auth req/sec - 10K.
cluster - 20 machines in pub. cloud.

Future Fact

- * diff. SLOs and SLAs per region
- * Conc. auth req/sec - 20K.
- * cluster should grow automatically based on load
- * traffic flow monitoring should be enabled to enforce future predictability.

SLI

Initial
① user.login req served < 100ms

SLD

95% of user.init.log.req served in < 100ms over past 24 hrs.

Sample set
adv. optimizn
conclusn

→ SLOs: Obj | SLI | Period

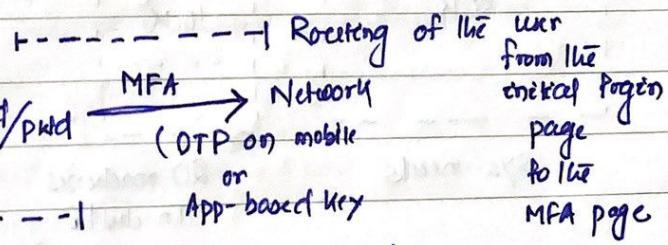
→ reliable network

→ globally dist. user base

→ availability regions - UK

Rest of Europe

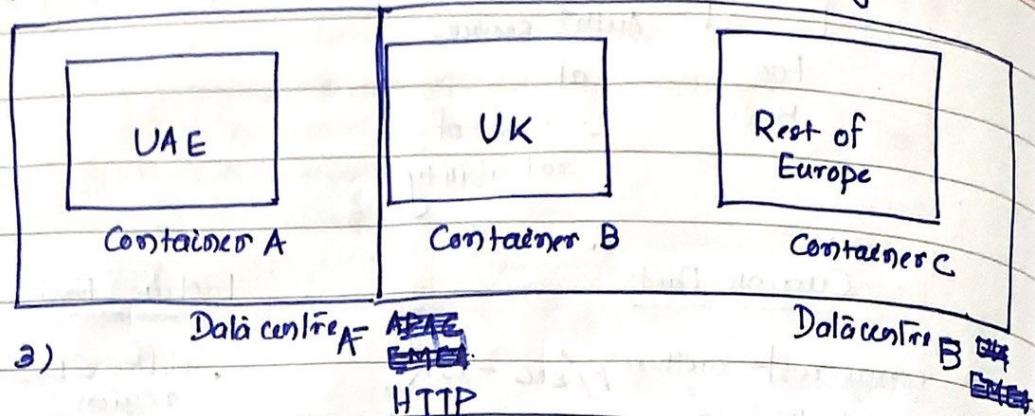
→ H/w: HDD mach.
SSD mach.



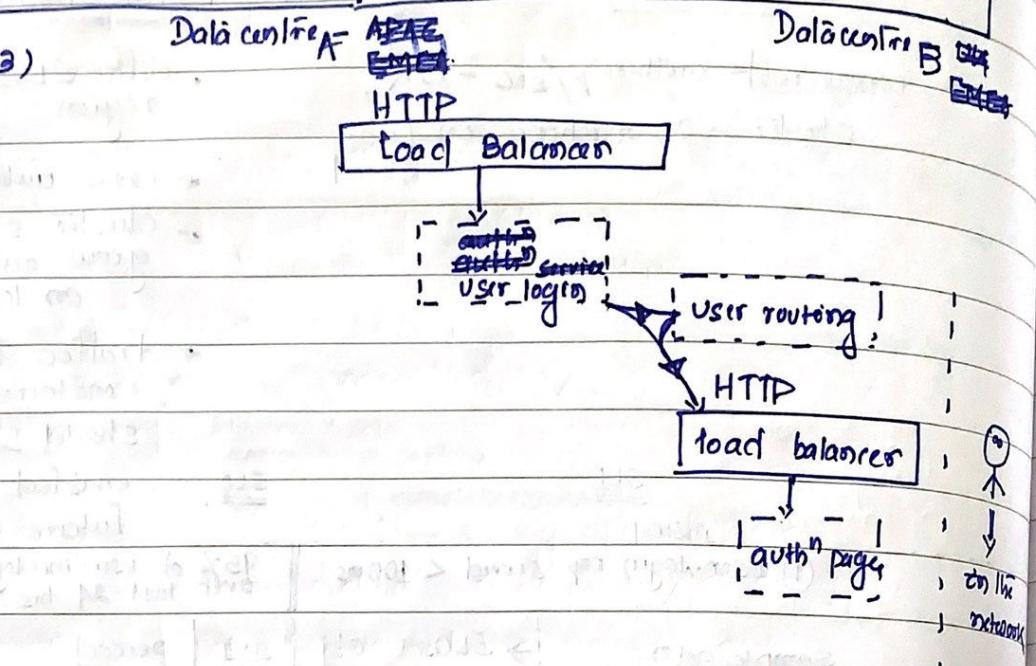
micro-services

Scaling:

2) Replication and sharded microservices per region



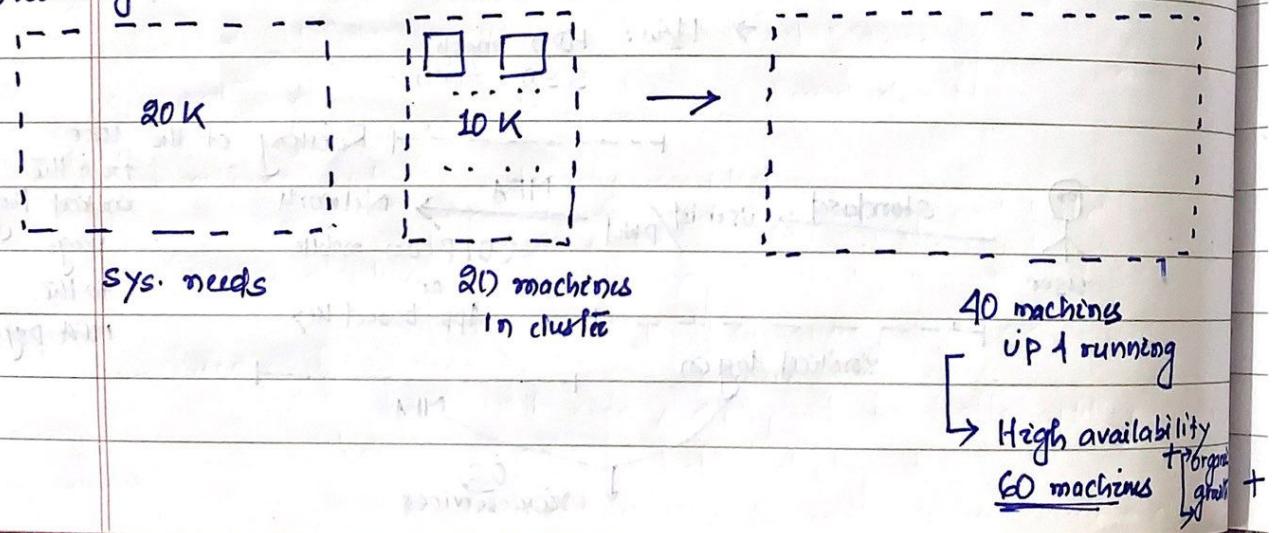
3)



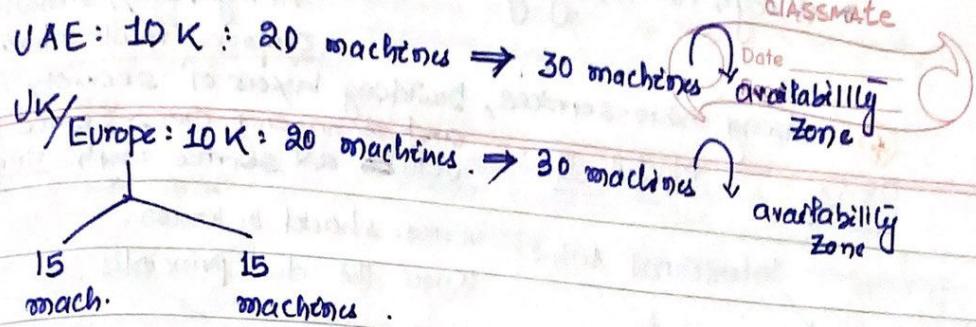
4) How many machines do we need to allocate to each micro-service?

Can we meet our SLO latency req.?

Provisioning:



Conc. users.



Based on data,

	<u>bandwidth</u>	<u>storage</u>	<u>machines reqd.</u>
users in each zone (UAE)	?	?	?

UK
Europe

* Machines per DC \times No. of DCs \times Infra. Tax + more load spikes.

Advanced optimizations:

* Capacity growth
Privacy req. (GDPR) \rightarrow General Data Protection Regulatⁿ

Peak users

* design for changing landscape: containers, service mesh, API gateway, orchestrator, improves fault tolerance.

* using micro-services, building layers of security context and managing cross-service comms well to a service mesh improves resilience.

08/14: Module-III

Intentional Arch: version should be known.
↓
multi-variance tested + Know the dev. principles.
perf. engineering

* des. for changing landscape:

↓
CI/CD
↓

SonarQube - testing must

more freq. releases; few changes/release

Containers & micro-services

* des. for changing security: IAM

SaaS adoptn

↓
zero trust networking

requires a modern identity and access security paradigm

* multi-grained service arch: microservices → more complex arch.

* Containers: management expose APIs using open standards

↓ Container orchestration: core capability of dev. arch.

Data Plane

Kubernetes: Extending Kub. API 3rd Party Components → Istio (Service mesh) Prometheus (monitoring)

Control Plane

→ Pilot (config)
Mixers (Policy, telemetry)
Citadel (encrypt)
authn)

automated service delivery
KEY app monitoring
distributed tracing
that are tightly coupled.

D3 Load Balancing
Round robin

* Pod: represents a unit of deployment which may consist of either a single container or small no. of containers



Reactive System: → loosely coupled

→ consists of gradually shifting prod. traffic from version A to version B and releasing to a subset of users

← Canary deployment → releasing to small set of users (5%)

→ feature flagging or routing logic can be used.

Service mesh → manages comms between microservices.

→ helps with Canary deployment for traffic splitting.

→ helps → C.D. for traffic deployments

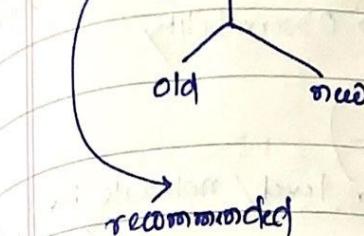
higher level of security for each service inside a cluster

mostly used for error rate and perf. monitoring with Farter rollback

Reactive Systems:

- Manifesto
- Sys. stays responsive under a varying workload.
- Responsible → cornerstone of usability.
- Resilient → resilience achieved by replacⁿ isolation and delegⁿ
- Elastic usage driven

Rolling upgrade:



asynchronous message passing to establish a boundary between components that ensures loose coupling, isolⁿ and locⁿ transparency.

versions of the code at the same time. This requires your app to handle N-1 compatibility.

* max. surge | max. unavailable

Blue-Green:

expensive

④ Router uses select the appⁿ version which will be exposed to the users.

The new version of an appⁿ is deployed in an identical env.

Tech. to understand

- * instant rollout/rollback
- + avoids versioning issue
- + stateful apps result few glitches

④ Threat:

Design for Security:

- Dev & analysis
- Threat intelligence
- adversary emulation & red teaming

MITRE ATT&CK

STRIVE

spoofing
intrusion
tampering

framework of adversary tactics and tech. based on real-world obsns.

Cyber-Security

④ Design for Resiliency:

ability of a sys. to gracefully handle and recover from failures.

④ Circuit-Breaker:

- to protect from cascading failures
- and to block traffic to downstream

④ Supervisor agent - service that continuously monitors the appⁿ/daemons/services and restarts them if they fail

any failure during execⁿ, the ability to gracefull rollback a transⁿ to its original state

← compensating transⁿ: gracefull roll-back

master set: Kubernetes

④ desⁿ for scalability: → horizontal scaling
asynchronous design

ability of a sys.

to handle increase in load without impact on perf.

statelessness → asynchronous design

→ measured using the responsiveness of the system on any given period

* Design for performance: of time, especially during peak loads.)

mechanical

empathy: when you use the tool of an understanding, how it operates but

results in optimum

service level / hypervisor level / load balancer level / network level / zone level / region level

→ web serverless arch. (event driven and cost eff. emplo)

experience more often → observability

go global and be elastic

design for availability: (MTTR)
↓ HA (MTBF)

design for reliability:

→ test recovery procedures

→ automatically recover from failure:

Scale horizontally to increase sys. availability

stop guessing capacity

manage changes using autom.

→ Zero Trust.

deployment

of service

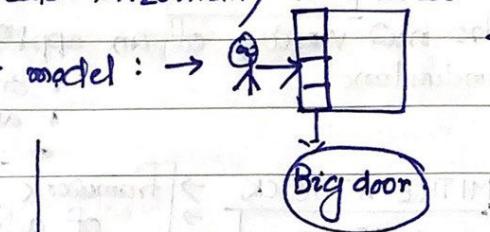
such

as

Kubernetes

pods

Castle and Moat model: →



TD based model

Full Stack Observability

high-level
overviews
of
sys.
health

Observability: → granular insight into the complete failure modes of the system.

→ uncover deeper, systemic issues.

* An observable sys. furnishes ample context about its inner workings, unlocking the ability to uncover deeper, systemic issues.

practitioner's view on observability:

That well

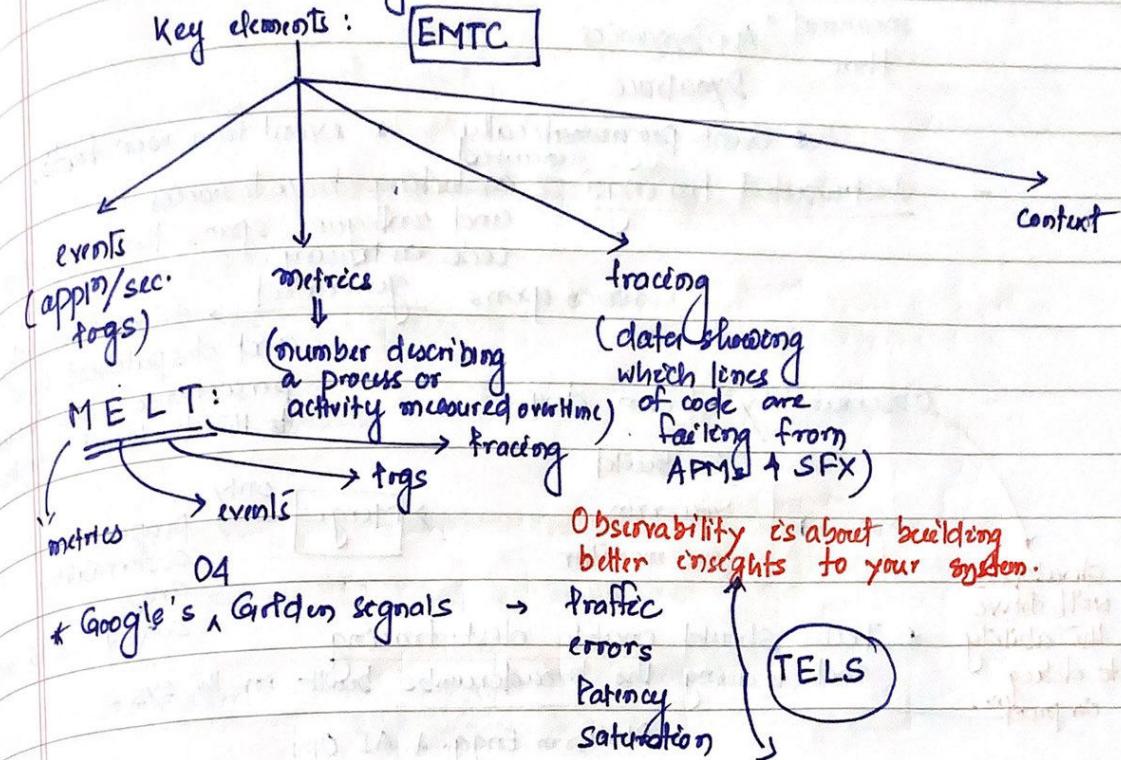
agg'd is the biggest enemy. A kill variety - making the info totally useless!

↓

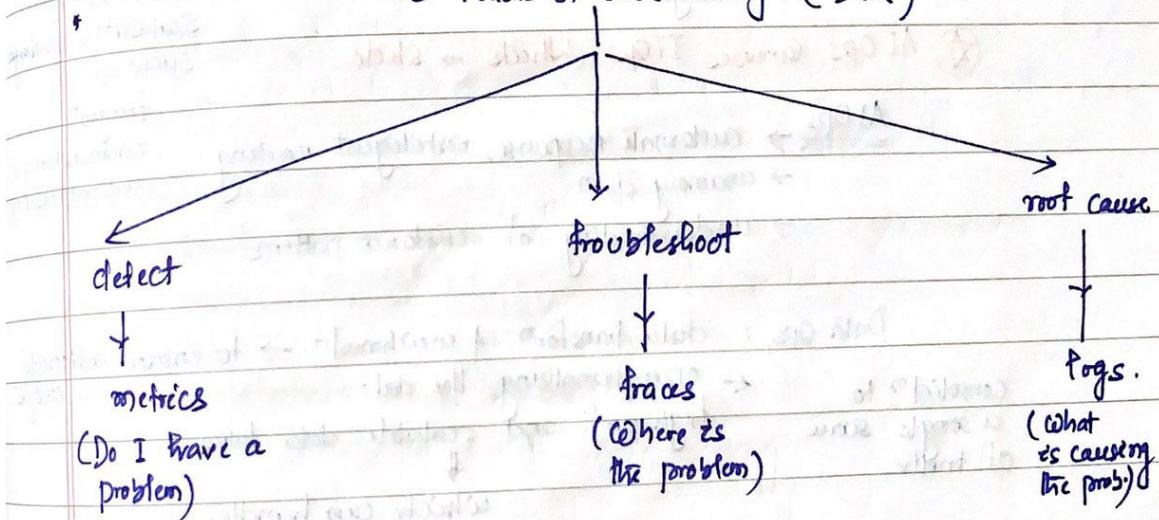
Complexity, elasticity and transience are the key considerations.

→ good observability makes detecting, investigating and responding to problems easier.

Observability: granular visibility into the unexpected behavior of complex sys. along with rich context



03 Pillars of Observability (DTR)



* saturation: degree to which the resource has extra ~~load~~
 Course: SRE - Google
 which it can't service, often queued

Synthetic and End-user monitoring:

Synthetic

measures thru AppDynamics
 Dynatrace

Real User

- * event programmatically generated from user traffic
- distributed tracing: including trace libraries and configure spans that will allow code instrumentation.

trace spans generated asynchronously

and dispatched to a persistence store hooked to the backend analytics eng.

Observability driven devops:

Developers
 developers well have the ability to debug in prod.

You build

you run

you monitor

Merge

only proper observability

are baked in the code

- * Tech. Should enable distributed tracing and tracing the breadcrumbs built in the sys.

Platform Engg. & AI Ops

* Platform SRE: approach to solve the org's scalability challenges such as fragm. inconsistency, unpredictability.

★ AI Ops → anomaly detection

★ AI Ops Greases ITOps wheels → slide

AI Ops → automatic response, intelligent routing
 → anomaly detection
 → understanding of customer patterns

Data Ops: data transform & enrichment → to ensure actionable insights
 consolidated to a single source of truth

← operationalizing the data
 high-perf. and scalable data lakes.

which can handle mixed workloads,
 diff. data types
 from sensors

* **Prescriptive**: 30 mins before the incident occurs.

08/15:

Incident Mgmt:

Time between incidents \rightarrow (MTBF) $\xrightarrow{\text{mean time between failures}}$

④ **Telemetry** - process of recording the behavior of your sys.

* **OODA Loop**

observe - orient - decide - act

* **Incident metrics to Biz. impact**

CLR - closed loop automated response

SRE - improve the first call resolⁿ rate (FCR)

Break the Glass: $\xrightarrow{\text{should be highly restricted}}$ available to SRE Team.

for zero trust networking should be available only from specific locⁿs. e.g.: panic rooms

mech. for BCP/DR.

* ~~Test~~ unplanned

should be tested regularly by the team(s) responsible for prodⁿ services, to make sure it functⁿ when you need it.

Swarming: collaborⁿ based process

no tiered support groups

AI/ML use cases: increase first-touch resolution

* AI-ML Use Cases

CLASSMATE
Date _____
Page _____

Atomic Habits	- Book
Lateral Thinking	→ Edward

Chaos Engg:

ID based zero trust arch.
adversity focussed

most common adopt DevSec - Cloud transform

not abt breaking things to prod

failure modes should be derived from the real world



Game Day → What
Who
When
Where
How

max visibility

Game Days → IP → C-E

- Pre-requisites:

Instrumentation
Social awareness

- DERT (Disaster Recovery Testing) → wheel of misfortune

Module-B

(SRE is the purist form of DevOps)

Table-top exercises

VAMOM

vision/values/methods/measures
objectives/measures

SRE and Incident Response Management

- Ability to respond and resolve incidents
 - ↳ true indicator of operational capabilities

- Key Responsibilities:

Manage incidents

Debugging the sys.

right mitigation

Org. the incident response if it req. broader coordination.

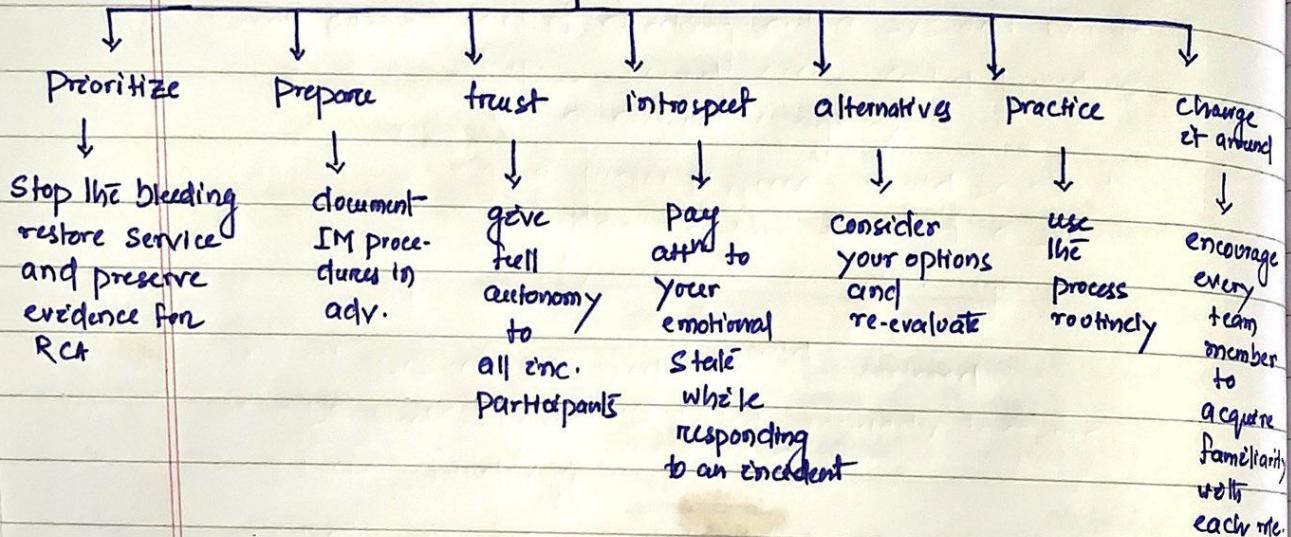
↑ detect prev^o

Strategies for deploy / roll back / roll forward (feature flags, blue green, canary)

auto remediation

reverting sys. to safe state

- Best practices for inc. mgmt:



- Key principles of incident response:

- maintain a clear line of command (recusive separation of responsibilities)
- designate clearly defined roles (inc commander, operational work, comm, planning)
- Keep a working record of debugging and mitigation
- close incidents early and often

3 C's:

- Coordinate → response effort
- Communicate → between incident responders within the org and outside also.
- Control → maintain control over incident response

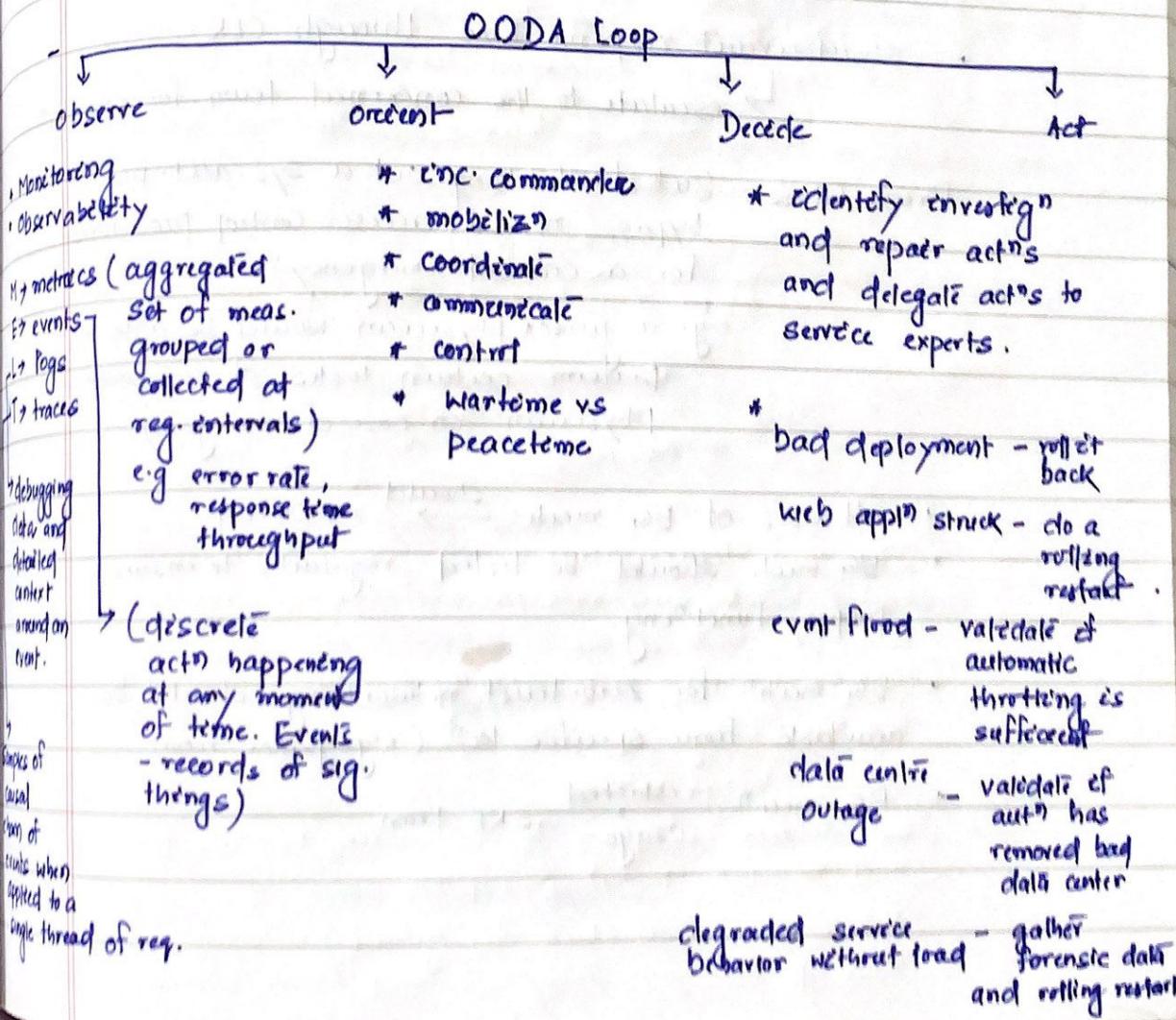
vertical functional aligned teams that operate inside a command and control approval system,

(Traditional ITSM)

Migration

horizontal alignment self-regulatory system

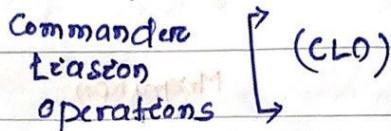
(DevOps + SRE)



- Tie Incident metrics to Biz. impact:

- no. of outbound notⁿ per sec.
 - no. of orders per sec.
 - no. of stream starts per sec.

- Incident Response Roles:



- CLR: Closed loop/automated response

- * SRE Key responsibility → automate the response processes
manual response: expensive on time and resources

- * Short circuit response time through CLR:

↳ escalate to the concerned team for investigation

- Break Glass: (act of checking out a sys. acc't. p/w to bypass normal access control procedures for a critical emergency.)

e.g. a junior physician would be able to perform certain tasks of a senior physician in case of emergency.

- * all uses of BG mech. → monitored

- * BG mech should be tested regularly to ensure proper functioning.

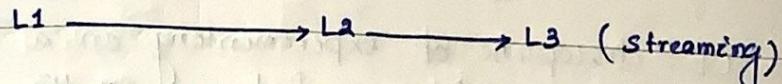
- * BG mech for zero-trust networking should be available from specific locⁿ (e.g. panic room)

- * BG mech restricted usage SRE team

Swarming: (apply organizational context)

- * collaborator based process
- * no tiered support groups
- * case should move directly to the person most likely to resolve it (direct alert to the SRE of the product)
- * no escalations from one group to another

Escalation based process:



AI/ML use cases:

- * auto categorization using chatbot
- * remove blind spots (ML algo. to detect how changes impact dynamic IT landscape)
- * intelligent routing (enables alerting with context and notifying teams with actionable alerts) → minimizing MTTD
- * automatically fulfill basic requests (NLP can help chat agents handle categories of req. and enc.)
- * increase first-touch resolutions
- * anomaly detection by flagging unusual repeat incidents (control the blast radius by clustering related incidents)
- * use predictive analytics to flag req's that could violate SLAs
- * predict impact
- * RCA
- * automate remediation (semi-autonomy remote with little human interv'n)
- * event resolution guidance

CHAOS ENGG:

- complexity of sys.
messy sys. engg.
driven for CE
Changing security landscape (ID based zero-trust arch.,
adversity focussed)

- * 50% of incidents → due to internal issues
remaining 50% → malicious or criminal attack

- CE: discipline of experimenting on a dev. sys. in order to build confidence on the sys. ability to withstand turbulent cond'n's

- * most common adopter - cloud transfr
- * not an advanced tech.
- * not about breaking things in prod
- * Failure modes - derived from the real world
- * minimize the blast radius
- * Brave an emergency stop.
- * not about injecting random ^{chaos} exp. into the sys. and seeing what happens
- * Start - in a controlled env.
Stop - immediately if things go badly

- Stop your docker service.
add delay to your network and chk for resiliency
burn CPU with stress

disaster recovery → periodically switch over sets
fill up storage and see sys. behavior.

detach storage volumes

use traffic shaping to introduce Network packet loss and network packet corruption

Kill services and check for graceful restore

mess up with your /etc/host and check behavior

Gamedays:

- * get relevant people in a room who are responsible for a sys. or a set of sys.
- * shut off a component that the sys. should be robust enough to tolerate without it
- * record the learning outcomes and bring the comp. outline

UAT
Env

Chaos Monkey Story:

- * born from closed transfⁿ (megⁿ) from DC to cloud)
- * put well defined prob.. in front of engg.
- * Terminate VMs or random VPC instances

CE Myths:

- * advanced capability (Certain level of sophistication within your sys. in order to undertake CE)
- * adds more chaos into the sys.
- * Only valuable in Prodⁿ

Gremions: destructⁿ as a serviceProperties of a Chaos Exp:

Game Days allows you to perform exp. with max^m visibility.

- * define steady state
- * formulate hypotheses
- * outline methodology
- * identify blast radius
- * observability - key
- * readily abortable

- Develop a learning culture about failure:

- * sys. safety
- * building safety margin into sys.
- * replace blame culture with learning culture
- * telemetry, expⁿ and instrumⁿ

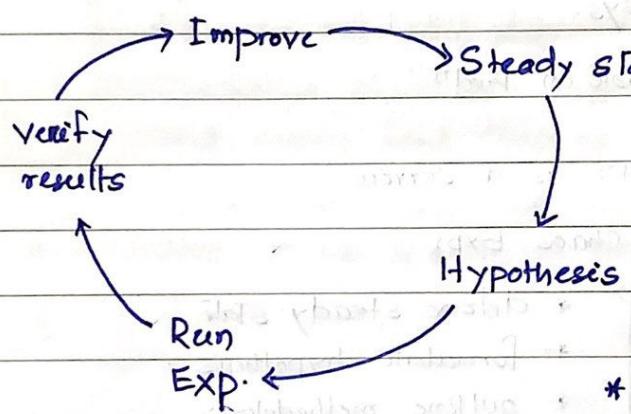
- GameDay Exercises:

- * 2-4 hrs in length
- * cross-functional teams
- * focussed on increasing resilience
- * used for manual CE
- * great intro. to CE

Recommendations

- i) use GameDays for new C. Exp
- ii) instead Exp. Deployment on new targets
- iii) Proving new C.E Tools
- iv) get everyone in the same locⁿ

- State of Chaos:



Quantifying Resilience:

- * time to detect
- * time for notⁿ/time for escalⁿ
- * time when graceful degⁿ Kicked-in
- * time for self-healing to happen
- * time to recovery
- * time to all-clear and back to normalcy

Experiment LifeCycle:

- 1) perform a GameDay exercise
- 2) validate exp. hypotheses
- 3) remediate findings and repeat exp.
- 4) once successful: automate exp.

GameDays: (The Basics)

- * plan and org. gameday exercise
- * dev. and evaluate C. Exp
- * execute live GameDay opns
- * conduct Pre-incident review
- * automate & evangelize results & actn taken

Pre-requisites for CE:

- * instⁿ to be able to detect degraded state in your sys.
- * social awareness (buy-in)
- * expectⁿ that hypotheses should be upheld (fix what is known to be broken)
- * alignment to respond
- * cultural infrastructure

before C.exp starts

- able to detect
- diff. in degradation between a control grp and an experimental grp
- pro-grade observability (before CE)
it's a myth.

- CE experiments focusses on security use cases
security inc. response
security observability
- ChaosStinger: open source framework for security focussed C. Exp

- CE @ Google: (improving dist. sys. reliability)

* DiRT (Disaster Recovery Testing)

↓
developed to

find critical vulnerabilities

determine tech. robustness

* commnd protocol

* revert/rollback for every step

* processes

humans

commnd

test everything that is involved
in the MTTR processes.

SRE - Purest Form of DevOps

SRE: making machines work for humans
↳ implement sys. that fix themselves.

metrics for success:

- * deployment freq.
- * lead time for changes
- * time to restore service
- * change failure rate

Impin' Pattern:

- * orgⁿ current state, maturity and priorities
- * SRE: sys. thinking + s/w Engg.
- * SRE - create pragmatic and cost-effective decisions across the entire product life cycle
- * SRES: bring singularity
 - ↓
(collaborators)
: create env's and autⁿ that will self-enable sys. to run meeting SLAs
 - : work on the entire continuum between s/w design, dev, testing, deployment, oper's and cust. exp.

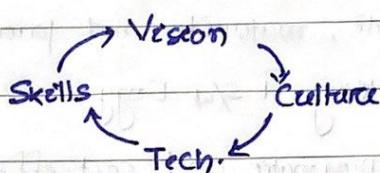
- 1. understand cust. expectations
- a. start small on obvious areas and scale out
- b. a) identify new areas of sys. engg and apply innovⁿ as needed
- b) Target areas: cloud migrn

C/I automⁿ
Observability
opsⁿ automⁿ

strategies

- ImplTM Roadmap:

- i) Blueprint for SRE
- ii) Assess / Maturity / Roadmap : identify bottlenecks in flow of value
- iii) Best Practices
- iv) Enablement → training strategy defined
- v) Governance → enterprise perf. dashboards
- vi) Evangelism → cross pollinate learnings
- vii) Rinse and repeat : refine & review (continuous improvement)
- viii) Ready SRE transfⁿ



- Build Stage Observability:

- * measure and optimize DevOps perf. to maximize ROI
 - optimize s/w delivery with goal based KPIs
 - support dev's with data driven recommendations
 - identify trends and predict release risk

- Run Stage observability:

- * Service
- infrastructure
- appTM ecosystem
- appTM
- biz. parameters
- security analytics

- SRE Execution:

SREs are

- * ^A owners and know the toolset of the product thoroughly

CLASSMATE

Date _____
Page _____

- * ensure consistency of tooling
- * error budgets remove conflict and help self-regulate flow.