

DevOps Interview Q&A – Do Cao Hieu

DevOps Interview Q&A – Do Cao Hieu

Bộ câu hỏi phỏng vấn dựa trên CV Portfolio (docaohieu.com)

Chuẩn bị cho các vị trí: DevOps Engineer / AI Infrastructure Engineer

PHẦN 1: INFRASTRUCTURE & ARCHITECTURE

Q1: Em ghi quản lý 10-server self-hosted infrastructure với 70+ Docker containers. Hãy mô tả kiến trúc tổng thể – server nào chạy gì, tại sao chia như vậy?

Dạ, em xin trình bày kiến trúc tổng thể của hệ thống. Em chia 10 servers ra theo vai trò rõ ràng để tách biệt concerns và tối ưu resource. Cụ thể thì em có 1 server làm load balancer chạy Nginx reverse proxy, server này đứng trước nhận tất cả traffic từ internet rồi route tới các backend servers phía sau. Tiếp theo em có 2 app servers chạy các ứng dụng chính như web apps, API services – mỗi server chạy khoảng 15-20 containers. Rồi em có 1 server dành riêng cho CI/CD chạy Jenkins và GitLab Runner, Harbor container registry – server này cần nhiều CPU và disk I/O nên em tách riêng để không ảnh hưởng tới production workloads. 1 server monitoring chạy toàn bộ observability stack gồm Prometheus, Grafana, Loki, Tempo, Alertmanager – em tách riêng vì monitoring phải independent, nếu chạy chung với app mà app server down thì mất luôn monitoring. 1 server database chạy PostgreSQL primary với Patroni, Redis. 1 server database standby cho PostgreSQL replication. 1 server cho AI services chạy OpenClaw, Engram, Qdrant vector database – nhóm này tốn RAM nhiều nên cần server riêng. Và 2 servers còn lại cho các services phụ trợ và staging environment.

Tất cả 10 servers kết nối với nhau qua WireGuard full mesh VPN, nên internal traffic đều encrypted. Lý do em chọn 10 servers thay vì ít hơn là vì em muốn tách biệt rõ ràng – monitoring không chạy chung app, database không chạy chung CI/CD. Nếu gộp lại 5 servers thì khi 1 server có vấn đề sẽ ảnh hưởng nhiều services cùng lúc. Còn nếu nhiều hơn 10 thì chi phí vận hành và quản lý tăng mà chưa cần thiết với scale hiện tại.

Q2: Con số tiết kiệm ~\$13,800/năm so với cloud — em tính ra sao? Breakdown cụ thể?

Dạ, em tính bằng cách so sánh chi phí thực tế hàng tháng của VPS với chi phí tương đương nếu dùng AWS managed services. Cụ thể, 10 VPS servers của em tốn khoảng \$150-200/tháng tổng cộng trên các provider như Hetzner. Nếu chuyển sang AWS thì tương đương khoảng 10 EC2 instances với specs tương tự sẽ tốn khoảng \$500-600/tháng. RDS PostgreSQL Multi-AZ thay cho self-hosted Patroni tốn thêm khoảng \$200-300/tháng. ElastiCache Redis khoảng \$100/tháng. Managed Grafana + CloudWatch khoảng \$100-150/tháng. ALB load balancer khoảng \$50-80/tháng. ECR container registry thay Harbor khoảng \$30-50/tháng. Cộng thêm data transfer costs, EBS volumes, và các chi phí phụ khác.

Tổng cộng trên AWS ước tính khoảng \$1,200-1,350/tháng, trong khi self-hosted của em chỉ khoảng \$150-200/tháng. Vậy tiết kiệm khoảng \$1,000-1,150/tháng, tương đương \$12,000-\$13,800/năm. Tuy nhiên em cũng thừa nhận rằng con số này chưa tính chi phí thời gian em bỏ ra để quản lý infrastructure — đó là trade-off lớn nhất. Nếu tính giờ công quản lý thì khoảng cách sẽ nhỏ hơn. Nhưng với em thì việc tự quản lý cũng là cơ hội học hỏi rất nhiều, nên em coi đó vừa là tiết kiệm chi phí vừa là đầu tư vào kiến thức.

Q3: 5-node WireGuard full mesh VPN — em triển khai cụ thể như thế nào?

Dạ, full mesh có nghĩa là mỗi node đều có peer config tới tất cả các node còn lại. Với 5 nodes thì em có $5 \times 4 / 2 = 10$ tunnels. Em generate key pair cho từng node bằng `wg genkey | tee privatekey | wg pubkey > publickey`. Mỗi node có 1 WireGuard interface `wg0` với 1 IP riêng trong subnet 10.0.0.0/24, ví dụ node 1 là 10.0.0.1, node 2 là 10.0.0.2, và cứ thế.

Em chọn full mesh thay vì hub-spoke vì lý do chính là không có single point of failure. Nếu dùng hub-spoke mà hub server down thì tất cả nodes mất kết nối với nhau. Full mesh thì mỗi node nói chuyện trực tiếp, latency thấp hơn, và nếu 1 node down thì các node còn lại vẫn giao tiếp bình thường. Trade-off là config complexity tăng theo $O(n^2)$ — mỗi khi thêm 1 node mới phải update config trên tất cả nodes hiện tại.

Để giải quyết vấn đề config complexity, em viết Ansible role để tự động generate WireGuard config. Khi thêm node mới, em chỉ cần thêm thông tin vào inventory file rồi chạy playbook — nó sẽ tự generate config cho node mới và update config trên tất cả existing nodes. Em cũng set `PersistentKeepalive = 25 seconds` cho NAT traversal, vì một số VPS nằm sau NAT. Firewall rules thì em configure để internal services chỉ listen trên WireGuard interface, không expose ra public internet.

Q4: Khi 1 trong 10 server bị down lúc 2AM, quy trình incident response của em là gì?

Dạ, đầu tiên về detection thì em có Alertmanager cấu hình gửi critical alerts qua Telegram. Khi 1 server down, Prometheus detect `node_exporter` không response trong vòng 2 phút, trigger alert `InstanceDown` và em nhận notification trên Telegram ngay lập tức, kể cả lúc 2 giờ sáng.

Bước tiếp theo là triage — em xác định server đó chạy services gì. Nếu là database primary server hay load balancer thì đó là P0, phải xử lý ngay. Nếu là staging server hay server phụ trợ thì có thể để sáng hôm sau. Với P0, em SSH vào server ngay, nếu SSH được thì check `journalctl -xe` xem system logs, `docker ps -a` xem containers nào bị crash, `dmesg | tail` xem kernel errors, `df -h` check disk space, `free -m` check memory. Thường thì nguyên nhân rơi vào mấy nhóm chính: disk full, OOM killer kill process, hoặc Docker daemon crash.

Sau khi xác định root cause, em fix — ví dụ nếu disk full thì cleanup Docker images cũ và logs, nếu OOM thì restart service và adjust resource limits. Nếu server hoàn toàn không SSH được thì em reboot qua provider console. Sau khi resolve xong, em viết post-mortem ngắn ghi lại root cause, timeline, và preventive measures — ví dụ thêm alert cho disk usage threshold thấp hơn, hoặc tăng resource limits. Nếu là hardware failure và server không recover được thì em dùng Ansible provision server mới và migrate containers sang.

Q5: Em có disaster recovery plan không? Nếu toàn bộ cluster mất thì rebuild mất bao lâu?

Dạ có. Backup strategy của em gồm 3 lớp. Thứ nhất là database: em chạy `pg_dump` daily lúc 2 giờ sáng qua cron job, backup file được rsync tới 1 server khác ở datacenter khác. Thứ hai là Docker volumes cho các services quan trọng cũng được rsync daily. Thứ ba là tất cả configuration — Ansible playbooks, Docker Compose files, Nginx configs — đều nằm trong git repositories.

Nếu toàn bộ cluster mất, quy trình rebuild sẽ là: provision VPS mới trên Hetzner hoặc DigitalOcean, chạy Ansible playbooks để setup base system, install Docker, configure WireGuard, deploy monitoring stack, rồi deploy applications. Về thời gian thì em ước tính RTO khoảng 4-6 giờ cho critical services, và khoảng 1-2 ngày để restore hoàn toàn 100% với tất cả data. RPO là khoảng 24 giờ vì backup chạy daily — nghĩa là worst case mất data của 1 ngày gần nhất.

Em thành thật nói là em chưa test full disaster recovery từ scratch — em đã test restore database từ backup thành công, đã test rebuild individual server từ Ansible, nhưng chưa simulate mất toàn bộ cluster. Đây là điểm em cần cải thiện. Em đang plan quarterly DR drill để test toàn bộ quy trình.

PHẦN 2: DOCKER & CONTAINER ORCHESTRATION

Q6: 70+ containers — em tổ chức và quản lý chúng bằng gì? Docker Compose hay Kubernetes?

Dạ, em dùng Docker Compose là chính. Em chưa dùng Kubernetes trong production vì với scale 10 servers và 70+ containers thì Docker Compose vẫn đáp ứng tốt, và complexity của K8s chưa justify cho quy mô này. Em biết K8s concepts — pods, deployments, services, ingress, ConfigMaps, Secrets — và đã lab trên local cluster, nhưng production thì vẫn là Compose.

Cách em tổ chức là mỗi server có một thư mục chứa các Docker Compose files, chia theo service group. Ví dụ trên monitoring server có `docker-compose.monitoring.yml` chứa Prometheus, Grafana, Loki, Alertmanager. Trên CI/CD server có `docker-compose.cicd.yml` chứa Jenkins, GitLab Runner, Harbor. Em dùng naming convention rõ ràng cho containers, networks, và volumes — prefix theo project hoặc service group, ví dụ `monitoring_prometheus`, `monitoring_grafana`. Networks thì em tách frontend, backend, database — containers chỉ join network cần thiết. Resource limits thì em set memory limit và CPU limit trong Compose file cho từng container, ví dụ PostgreSQL được allocate 4GB RAM, Prometheus 2GB, các app containers thường 512MB-1GB.

Logging driver em configure là json-file với max-size và max-file rotation, rồi Promtail đọc Docker logs và ship về Loki. Em cũng đang evaluate chuyển sang K8s khi scale lớn hơn, nhưng hiện tại Docker Compose vẫn là lựa chọn phù hợp nhất.

Q7: Khi cần update 1 container trong production mà không downtime, em làm như thế nào?

Dạ, quy trình zero-downtime update của em phụ thuộc vào service đó có chạy nhiều instance hay không. Đối với services có thể chạy multiple instances, em dùng kiểu blue-green: pull image mới, start container mới với port khác, verify health check endpoint `/health` trả về OK, rồi update Nginx upstream config để route traffic sang container mới, reload Nginx bằng `nginx -s reload` (không downtime), rồi stop container cũ sau khi confirm container mới hoạt động ổn.

Đối với services single instance như database thì em plan maintenance window ngắn, thường ngoài giờ peak. Đối với stateless app containers thì em viết script automate quá trình: pull new image, create new container, health check, switch traffic, remove old container. Rollback plan là em luôn giữ old image tag — nếu container mới có vấn đề thì revert lại image cũ trong vòng 1-2 phút.

Nếu có database migration kèm theo thì em apply migration trước, đảm bảo migration backward-compatible — nghĩa là schema mới vẫn hoạt động với code cũ. Deploy code mới sau khi migration xong. Đây là pattern em học được và áp dụng nhất quán.

Q8: Container nào tốn resource nhất trong 70+ containers? Em optimize ra sao?

Dạ, em monitor resource usage qua cAdvisor metrics trên Prometheus và visualize trên Grafana. Các containers tốn resource nhất là PostgreSQL (khoảng 3-4GB RAM), Prometheus (khoảng 2GB RAM vì lưu time-series data), Qdrant vector database (khoảng 2GB RAM cho vector indexing), và Jenkins (khoảng 1.5-2GB RAM khi build đồng thời). Loki cũng tốn khá nhiều disk I/O vì ingest logs từ tất cả containers.

Về optimization, em áp dụng mấy cách chính. Một là multi-stage Docker builds để giảm image size — ví dụ Go application build stage dùng `golang` image rồi copy binary sang `alpine` image, giảm từ 800MB xuống còn 20-30MB. Hai là set resource limits rõ ràng trong Compose file để tránh 1 container chiếm hết resource của server. Ba là connection pooling cho database — em dùng PgBouncer đứng trước PostgreSQL để giảm số connections thực tế. Bốn là tune Prometheus retention và scrape interval — không cần scrape mỗi 10 giây cho tất

cả targets, một số non-critical metrics em set 60 giây. Năm là Docker image cleanup tự động bằng cron job chạy docker system prune hàng tuần để xóa unused images và volumes.

Q9: Docker networking — em setup networking giữa containers trên nhiều server như thế nào?

Dạ, em không dùng Docker Swarm overlay network. Thay vào đó, em dùng WireGuard VPN làm network layer giữa các servers. Containers trên cùng 1 server giao tiếp qua Docker bridge networks — em tạo các custom networks riêng cho từng nhóm, ví dụ frontend-net, backend-net, db-net. Containers chỉ join network mà chúng cần, ví dụ web app join cả frontend-net và backend-net, nhưng database chỉ join db-net.

Khi containers trên server A cần nói chuyện với containers trên server B, traffic đi qua WireGuard tunnel. Em expose port của service trên WireGuard interface, ví dụ PostgreSQL listen trên 10.0.0.3:5432 (WireGuard IP), app containers trên server khác kết nối tới IP đó. Nginx reverse proxy đứng trước tất cả web-facing services, terminate SSL và route traffic dựa trên domain name.

DNS resolution thì em dùng combination — Docker internal DNS cho container-to-container trên cùng server, và config file /etc/hosts hoặc hostname resolution qua WireGuard IPs cho cross-server communication. Em cũng set up internal DNS entries trong hosts file được quản lý bởi Ansible để đảm bảo consistency across all servers.

Q10: Em quản lý Docker images như thế nào? Versioning strategy?

Dạ, em self-host Harbor container registry. Tất cả Docker images được build trong CI/CD pipeline rồi push lên Harbor thay vì Docker Hub. Lý do chính là em muốn full control, không bị rate limit, và pull speed nhanh hơn vì registry nằm trong internal network.

Tagging strategy của em là mỗi image có ít nhất 2 tags: git commit SHA (ví dụ harbor.internal/myapp:a1b2c3d) để traceable, và semantic version (ví dụ harbor.internal/myapp:1.2.3) cho production releases. Em không dùng latest tag cho production vì nó không deterministic. CI/CD pipeline tự động build, tag, và push khi code merge vào main branch.

Harbor có tích hợp Trivy vulnerability scanning — mỗi image push lên được scan tự động, nếu có critical vulnerability thì em nhận alert. Em cũng configure retention policy trong Harbor để tự động xóa images cũ hơn 30 ngày (giữ lại 10 versions gần nhất), và chạy garbage collection weekly để reclaim disk space. Multi-stage Dockerfile là standard practice — em enforce qua CI pipeline, check Dockerfile best practices.

PHẦN 3: MONITORING & OBSERVABILITY

Q11: Prometheus, Grafana, Loki, Tempo, Alertmanager — em triển khai full stack này như thế nào? Mô tả data flow.

Dạ, toàn bộ observability stack chạy trên 1 server monitoring riêng. Data flow như sau: Prometheus chạy ở chế độ pull — cứ mỗi 15-30 giây nó scrape metrics từ các targets. Trên mỗi server em cài `node_exporter` để expose system metrics (CPU, RAM, disk, network), `cAdvisor` để expose container metrics (per-container CPU, memory, network I/O). Các applications cũng expose custom metrics endpoints. Prometheus lưu data vào local TSDB, retention em set 30 ngày.

Loki nhận logs từ tất cả containers. Em cài Promtail trên mỗi server — Promtail đọc Docker container logs (JSON log files) và ship về Loki central server. Loki index labels (container name, server, log level) nhưng không index full text — đây là thiết kế của Loki giúp tiết kiệm storage so với Elasticsearch.

Tempo nhận distributed traces — hiện tại em mới integrate cho một số services chính, gửi traces qua OpenTelemetry SDK. Alertmanager nhận alerts từ Prometheus recording rules — khi metric thỏa điều kiện alert (ví dụ CPU > 90% trong 10 phút), Prometheus gửi alert tới Alertmanager, Alertmanager group alerts theo labels rồi route tới Telegram bot của em.

Grafana là visualization layer — nó query Prometheus bằng PromQL, query Loki bằng LogQL, query Tempo cho traces. Em configure data sources cho cả 3 và tạo dashboards tập trung. Đây là single pane of glass cho toàn bộ infrastructure.

Q12: Em viết alert rules như thế nào? Cho ví dụ 3-5 critical alerts quan trọng nhất.

Dạ, alert rules em viết trong Prometheus config dưới dạng YAML. Những critical alerts quan trọng nhất của em gồm:

Thứ nhất là `InstanceDown` — khi `node_exporter` của bất kỳ server nào không response quá 3 phút thì fire alert. Đây là alert quan trọng nhất vì nó báo hiệu server có thể đã crash hoặc mất network.

Thứ hai là `DiskSpaceLow` — khi disk usage vượt 85%. Em từng bị incident disk full làm PostgreSQL crash, nên alert này rất critical. Em set 2 mức: warning ở 80%, critical ở 90%.

Thứ ba là `HighMemoryUsage` — khi container memory usage vượt 90% limit. Đây thường là dấu hiệu memory leak, và nếu không xử lý kịp thì OOM killer sẽ kill container.

Thứ tư là `PostgreSQLReplicationLag` — khi replication lag giữa primary và standby vượt 30 giây. Điều này có thể ảnh hưởng data consistency nếu failover xảy ra.

Thứ năm là `ContainerRestarting` — khi 1 container restart nhiều hơn 3 lần trong 15 phút, báo hiệu crash loop cần investigate.

Để tránh alert fatigue, em set severity levels rõ ràng — critical alerts gửi Telegram ngay, warning alerts group lại gửi mỗi 30 phút. Em cũng dùng inhibition rules — ví dụ khi server down thì suppress tất cả alerts từ containers trên server đó, chỉ nhận 1 alert `InstanceDown` thôi.

Q13: Khi một service chậm lại (high latency), em debug bằng observability stack như thế nào? Mô tả step-by-step.

Dạ, quy trình debug của em theo RED method — Rate, Errors, Duration. Bước đầu tiên em mở Grafana dashboard, check request rate có spike đột ngột không — nếu traffic tăng gấp 3-4 lần so với bình thường thì đó có thể là nguyên nhân. Tiếp theo check error rate — nếu error rate tăng cùng lúc với latency thì khả năng cao là downstream dependency có vấn đề.

Bước hai em drill down vào resource metrics của server chạy service đó — CPU, memory, disk I/O, network. Nếu CPU saturated thì service đang bị compute-bound. Nếu memory gần limit thì có thể đang swap. Nếu disk I/O cao thì có thể database đang full scan.

Bước ba em vào Loki xem logs của service đó, filter theo time range lúc latency tăng. Em tìm error messages, slow query warnings, timeout errors. PromQL query em hay dùng là `rate(http_request_duration_seconds_bucket[5m])` để xem latency distribution.

Bước bốn em check downstream dependencies — database query time qua PostgreSQL metrics (`pg_stat_statements`), Redis latency, external API response time. Rất nhiều trường hợp latency cao ở application layer nhưng root cause lại ở database — ví dụ missing index làm query chạy full table scan.

Bước năm nếu có Tempo traces thì em trace 1 slow request end-to-end, xem thời gian bị nghẽn ở đâu trong request flow. Cuối cùng em correlate tất cả data — thời điểm latency tăng có trùng với deployment mới, traffic spike, hay cron job nặng chạy cùng lúc không.

Q14: Grafana dashboard em custom những gì? Metrics quan trọng nhất?

Dạ, em có mấy nhóm dashboards chính. Dashboard đầu tiên là Infrastructure Overview — hiển thị tất cả 10 servers trên 1 trang, mỗi server có panels cho CPU usage, memory usage, disk usage, network I/O. Em dùng Grafana variables để filter theo server. Dashboard này cho em cái nhìn tổng quan trong 5 giây.

Dashboard thứ hai là Container Dashboard — liệt kê tất cả 70+ containers, sorted theo resource usage. Panels gồm CPU per container, memory per container, network I/O, restart count. Em dùng cAdvisor metrics cho dashboard này. Khi 1 container bất thường thì nổi lên ngay trên top.

Dashboard thứ ba là Database Dashboard cho PostgreSQL — connections active/idle, transactions per second, cache hit ratio (target > 99%), replication lag, dead tuples (cần vacuum), query execution time từ `pg_stat_statements`. Dashboard này giúp em proactively detect database issues.

Dashboard thứ tư là Application Dashboard theo RED method — request rate, error rate percentage, và latency percentiles p50, p95, p99. Em track p99 vì đó là tail latency ảnh hưởng user experience.

Ngoài ra em có CI/CD Dashboard tracking build duration trung bình, success/failure rate, và deployment frequency. Em coi monitoring dashboard là living documents — liên tục update và improve dựa trên incidents và lessons learned.

PHẦN 4: CI/CD & GITOPS

Q15: Em dùng Jenkins, GitLab CI, GitHub Actions — tại sao dùng nhiều tool? Use case cho từng cái?

Dạ, lý do em dùng nhiều CI/CD tools là vì mỗi tool phù hợp với context khác nhau. Jenkins em dùng cho các internal projects phức tạp cần nhiều customization — Jenkins chạy self-hosted trên server riêng, em có full control, cài plugins theo nhu cầu, build complex pipelines với Groovy scripting. Ví dụ pipeline build và deploy Engram qua Jenkins vì cần nhiều stages custom.

GitLab CI em dùng cho projects hosted trên GitLab instance — ưu điểm là integration chặt chẽ với GitLab, config bằng `.gitlab-ci.yml` nằm ngay trong repo, runners em cũng self-host. GitHub Actions em dùng cho open source projects và projects trên GitHub — ưu điểm là community actions marketplace rất phong phú, ví dụ GitPocket project dùng GitHub Actions cho React Native build.

Em thừa nhận trade-off là phải maintain knowledge và config cho 3 tools khác nhau, tăng operational overhead. Về lâu dài em đang có xu hướng converge về GitHub Actions cho majority of projects vì ecosystem mạnh hơn và maintenance ít hơn Jenkins. Jenkins em giữ cho legacy pipelines và cases cần deep customization.

Q16: Mô tả CI/CD pipeline hoàn chỉnh nhất em từng build — từ commit tới production.

Dạ, pipeline hoàn chỉnh nhất của em cho một web application service gồm các stages sau. Trigger là khi developer push code hoặc merge merge request vào main branch. Stage 1 là Lint — chạy linter check code style, format check, static analysis. Stage 2 là Unit Test — chạy test suite, report coverage percentage. Stage 3 là Build — build Docker image bằng multi-stage Dockerfile. Stage 4 là Push — tag image với git commit SHA và push lên Harbor registry. Tại đây Trivy tự động scan image cho vulnerabilities.

Stage 5 là Deploy to Staging — pull image mới trên staging server, docker compose up, chạy health checks. Stage 6 là Integration Test — chạy automated tests against staging environment, verify API endpoints, check database migrations. Stage 7 — nếu tất cả pass — là manual approval gate, em review staging hoạt động OK rồi approve. Stage 8 là Deploy to Production — pull image trên production server, thực hiện blue-green switch qua Nginx, verify health check endpoint.

Xuyên suốt pipeline, mỗi stage fail đều gửi notification Telegram cho em biết. Secrets trong pipeline được inject qua environment variables — database passwords, API keys — không hardcode trong code hay pipeline config. Build artifacts như test reports, coverage reports được archive cho reference.

Q17: Zero-downtime deployment — em implement cụ thể như thế nào?

Dạ, cách em implement zero-downtime deployment cho web services là sử dụng blue-green pattern kết hợp Nginx. Cụ thể, application chạy trên port 8080 (blue instance). Khi deploy version mới, em start

container mới trên port 8081 (green instance). Em chờ health check endpoint /health trên port 8081 trả về 200 OK — thường em set timeout 30 giây, retry mỗi 2 giây.

Khi green instance healthy, em update Nginx upstream config để point tới port 8081 thay vì 8080, rồi `nginx -s reload` — Nginx reload gracefully không drop existing connections. Sau khi verify traffic đang chạy qua green instance bình thường (check logs, check metrics trong 2-3 phút), em stop blue instance cũ trên port 8080.

Nếu green instance có vấn đề sau deploy — error rate tăng, latency tăng — em rollback bằng cách update Nginx config trở lại port 8080 (blue instance vẫn đang chạy), reload Nginx. Toàn bộ rollback mất dưới 1 phút.

Đối với database migrations, em áp dụng nguyên tắc backward-compatible migrations — chỉ thêm columns, không xóa hay rename trong cùng deployment. Migration chạy trước khi deploy code mới, đảm bảo code cũ vẫn hoạt động được với schema mới. Column cũ được cleanup trong deployment tiếp theo khi chắc chắn code mới đã stable.

Q18: Harbor container registry — tại sao self-hosted? So sánh với Docker Hub, GHCR?

Dạ, em chọn self-host Harbor vì mấy lý do chính. Thứ nhất là control — em full control over data, không phụ thuộc third-party service uptime. Docker Hub đã nhiều lần có rate limiting issues, ảnh hưởng deployments. Thứ hai là security — Harbor tích hợp Trivy vulnerability scanning, mỗi image push lên được auto-scan. Em configure policy block images có critical CVEs. Thứ ba là performance — registry nằm trong internal network nên pull speed rất nhanh, đặc biệt khi deploy 70+ containers cùng lúc không bị bottleneck bandwidth. Thứ tư là cost — Harbor miễn phí, trong khi Docker Hub paid plan cho private repos, và GHCR cũng có giới hạn storage.

So sánh thì Docker Hub có ecosystem lớn nhất, community images phong phú, nhưng rate limits phiền. GHCR tiện khi dùng GitHub Actions nhưng ít features hơn Harbor. Harbor có thêm replication (replicate images giữa registries), RBAC chi tiết, retention policies, và webhook integrations. Trade-off là em phải maintain Harbor instance — update, backup, monitor. Nhưng vì em đã có observability stack sẵn rồi nên overhead không quá nhiều.

PHẦN 5: DATABASE & DATA MANAGEMENT

Q19: PostgreSQL Patroni HA — em setup như thế nào? Bao nhiêu node? Failover hoạt động ra sao?

Dạ, em setup PostgreSQL HA bằng Patroni với 2 nodes — 1 primary và 1 standby. Patroni quản lý leader election và automatic failover. Em dùng etcd làm DCS (Distributed Configuration Store) để Patroni các nodes đồng bộ state với nhau. Architecture là Patroni wrap PostgreSQL process trên mỗi node, liên tục report health status vào etcd.

Khi primary node down, Patroni trên standby node detect leader key trong etcd expire (thường sau 30 giây), tự động promote standby lên làm primary. Replication em dùng streaming replication asynchronous — em chấp nhận rủi ro mất vài transactions cuối vì synchronous replication ảnh hưởng write performance. Applications connect qua PgBouncer connection pooler, và PgBouncer config trỏ tới Patroni API endpoint để auto-detect primary node mới.

Em đã test manual failover bằng cách stop PostgreSQL trên primary — Patroni promote standby thành công trong khoảng 30-45 giây. Trong thời gian failover thì database unavailable, applications nhận connection errors nhưng retry logic handle được. Sau khi primary cũ recover, nó tự động join lại cluster như standby và sync data từ primary mới. Em monitor replication lag qua Prometheus, alert khi lag vượt 30 giây.

Q20: Redis em dùng cho mục đích gì cụ thể? Setup ra sao?

Dạ, em dùng Redis cho mấy use cases chính. Thứ nhất là caching — cache database query results và API responses để giảm load cho PostgreSQL. Thứ hai là session store cho web applications — sessions lưu trong Redis thay vì disk, fast access và dễ share giữa multiple app instances. Thứ ba là message queue — em dùng Redis pub/sub cho real-time notifications và lightweight message passing giữa services. Thứ tư là rate limiting cho API endpoints.

Setup hiện tại em dùng Redis standalone, chưa cần Sentinel hay Cluster vì scale chưa đòi hỏi. Persistence em configure cả RDB snapshots (mỗi 5 phút nếu có > 100 keys thay đổi) và AOF (append-only file) để minimize data loss khi crash. Memory policy em set allkeys-lru — khi Redis đạt maxmemory limit (em set 1GB) thì tự động evict least recently used keys. Em monitor Redis qua redis_exporter gửi metrics về Prometheus — track memory usage, hit rate (đang đạt khoảng 95%), connected clients, và operations per second. Cache hit rate 95% nghĩa là chỉ 5% requests phải query database, giảm đáng kể database load.

Q21: Backup strategy cho toàn bộ data — database, configs, volumes?

Dạ, backup strategy của em gồm mấy lớp. PostgreSQL em chạy pg_dump daily lúc 2 giờ sáng qua cron job — dump ra file compressed .sql.gz, retain 7 bản daily, 4 bản weekly (mỗi chủ nhật), và 3 bản monthly. Backup files được rsync tới 1 server backup ở location khác. Redis thì RDB snapshot cũng được backup daily cùng schedule.

Docker volumes cho các services quan trọng (Grafana dashboards, Prometheus data, Loki data) em rsync daily tới backup server. Tất cả configuration files — Ansible playbooks, Docker Compose files, Nginx configs, monitoring configs — đều commit trong git repositories, nên luôn có version history.

Em thành thật nói là retention policy chưa hoàn hảo — em chưa có offsite backup ở cloud (S3 hay tương đương), tất cả backup hiện tại trên 2 servers khác nhau nhưng cùng 1 provider. Đây là risk em đang plan mitigate bằng cách thêm S3-compatible storage (ví dụ Backblaze B2) cho offsite backup. Về test restore, em đã test pg_restore từ backup file thành công, verify data integrity bằng cách compare row counts và checksums. Nhưng full DR test từ scratch thì em chưa thực hiện, đang plan quarterly.

Q22: Database performance tuning — em đã tune PostgreSQL những gì?

Dạ, em tune PostgreSQL dựa trên workload và available RAM. Với server có 8GB RAM dành cho PostgreSQL, em set `shared_buffers = 2GB` (25% RAM), `effective_cache_size = 6GB` (75% RAM — giúp query planner biết có bao nhiêu cache available), `work_mem = 64MB` (cho sort operations), `maintenance_work_mem = 512MB` (cho VACUUM, CREATE INDEX). Connection pooling em dùng PgBouncer ở chế độ transaction mode, pool size 100 connections, giảm từ hàng trăm application connections xuống còn 20-30 actual database connections.

Về query optimization, em bật `pg_stat_statements` extension để track slow queries. Có 1 case cụ thể em nhớ rõ: query tìm kiếm entities trong Engram knowledge graph chạy 2-3 giây vì full table scan trên bảng lớn. Em chạy EXPLAIN ANALYZE, thấy Sequential Scan trên 500K rows. Em thêm composite index trên columns thường query, query time giảm xuống còn 5-10ms — improvement khoảng 200-300x.

Autovacuum em tune để chạy aggressive hơn default — `autovacuum_vacuum_scale_factor = 0.05` (thay vì 0.2), vì em thấy dead tuples tích tụ gây bloat. Em monitor vacuum progress qua `pg_stat_user_tables`. Slow query log em set `log_min_duration_statement = 500ms` để capture queries chậm hơn 500ms.

PHẦN 6: INFRASTRUCTURE AS CODE

Q23: Ansible — em dùng như thế nào? Playbook lớn nhất gồm những gì?

Dạ, Ansible là IaC tool chính của em. Em dùng cho server provisioning, configuration management, và application deployment. Organization theo best practices: em có inventory file chia servers thành groups (`app_servers`, `db_servers`, `monitoring`, `ci/cd`), `group_vars` chứa variables chung cho mỗi group, `host_vars` cho variables riêng từng server. Roles em tách rõ ràng: role base cho OS hardening, role docker cho Docker installation, role wireguard cho VPN setup, role monitoring cho `node_exporter` + `promtail`, role nginx cho reverse proxy.

Playbook lớn nhất là “full cluster provision from scratch” — provision 10 servers hoàn chỉnh. Playbook này gồm: install base packages, create users với SSH keys, configure firewall (UFW rules), setup automatic security updates, install Docker và Docker Compose, generate và deploy WireGuard configs cho full mesh VPN, install monitoring agents, deploy application containers, configure Nginx reverse proxy. Chạy từ đầu tới cuối mất khoảng 30-40 phút cho 10 servers.

Tất cả playbooks em viết idempotent — chạy lại bao nhiêu lần cũng cho kết quả giống nhau, không side effects. Em test bằng `--check` --diff trước khi apply để xem changes sẽ là gì. Secrets em encrypt bằng `ansible-vault` — database passwords, API keys, WireGuard private keys đều encrypted trong vault, decrypt at runtime bằng vault password. Em chưa dùng Molecule cho role testing nhưng đang plan integrate.

Q24: Terraform — em dùng cho mục đích gì? State management ra sao?

Dạ, em thành thật nói là Terraform em chủ yếu đang học và lab, chưa dùng nhiều trong production. Infrastructure của em chủ yếu self-hosted trên VPS, và em provision VPS bằng provider dashboard hoặc API rồi configure bằng Ansible. Terraform em đã dùng thử cho việc provision Hetzner VPS — viết Terraform config define server specs, SSH keys, networking — nhưng chưa triển khai rộng.

Em hiểu Terraform concepts: HCL language, providers, resources, data sources, modules, state management. State em lưu local file khi lab, nhưng em biết production nên dùng remote backend như S3 + DynamoDB locking hoặc Terraform Cloud. Plan → Apply workflow em nắm rõ — luôn review plan trước khi apply, không bao giờ apply blind.

Gap thực sự của em là chưa có kinh nghiệm Terraform production-grade — managing complex infrastructure with modules, handling state drift, managing multiple environments (dev/staging/prod) bằng workspaces. Đây là area em đang tích cực học và plan integrate vào workflow — ví dụ dùng Terraform provision VPS rồi Ansible configure, tạo complete IaC pipeline.

Q25: Khi thêm 1 server mới vào cluster, quy trình từ A-Z là gì?

Dạ, quy trình từ A-Z khi em thêm 1 server mới gồm các bước sau. Bước 1: Provision VPS trên Hetzner dashboard — chọn specs (CPU, RAM, disk), chọn OS image (Ubuntu 22.04), add SSH key. Mất khoảng 2 phút.

Bước 2: Thêm server vào Ansible inventory file — khai báo IP, hostname, group membership, host-specific variables. Bước 3: Chạy Ansible base playbook trên server mới — playbook này setup OS hardening (disable root SSH, configure UFW firewall, install fail2ban, setup unattended-upgrades), tạo admin user với SSH key, install base packages.

Bước 4: Chạy Docker role — install Docker Engine, Docker Compose, configure logging driver (json-file với rotation). Bước 5: Chạy WireGuard role — đây là phần hay nhất. Ansible role generate WireGuard keypair cho node mới, generate config file với peer entries cho tất cả existing nodes, đồng thời update config trên tất cả existing nodes để thêm peer entry cho node mới, rồi restart WireGuard trên tất cả nodes. Sau bước này, node mới đã có VPN connectivity tới toàn bộ cluster.

Bước 6: Chạy monitoring role — install node_exporter và Promtail, update Prometheus config thêm target mới. Bước 7: Deploy services cần thiết lên server mới bằng Docker Compose. Bước 8: Verify — check Grafana dashboard thấy server mới xuất hiện với metrics, check Prometheus targets page thấy status UP, test SSH qua WireGuard IP, test container connectivity.

Toàn bộ quy trình mất khoảng 15-20 phút, phần lớn là automated qua Ansible.

PHẦN 7: NETWORKING & SECURITY

Q26: Nginx / HAProxy — em dùng cho load balancing như thế nào? Config ra sao?

Dạ, em dùng Nginx làm reverse proxy và load balancer chính. Nginx chạy trên dedicated load balancer server, terminate SSL/TLS cho tất cả domains. SSL certificates em dùng Let's Encrypt với certbot auto-renewal — cron job chạy certbot renew 2 lần/ngày, certificates auto-renew trước khi expire 30 ngày.

Config Nginx em tổ chức mỗi domain/service 1 file trong sites-available/, symlink tới sites-enabled/. Mỗi site config gồm: SSL termination, upstream definition pointing tới backend containers (qua WireGuard IP hoặc localhost port), proxy headers (X-Real-IP, X-Forwarded-For, X-Forwarded-Proto), và security headers (HSTS, X-Frame-Options, X-Content-Type-Options, CSP).

Về load balancing, em dùng round-robin cho stateless services khi có multiple instances. Health checks em configure bằng proxy_next_upstream directive — nếu upstream trả error thì Nginx tự route sang instance khác. Rate limiting em set bằng limit_req_zone để chống abuse — ví dụ API endpoints limit 10 requests/second per IP. Static assets em cache tại Nginx level với expires headers.

HAProxy em biết và đã lab nhưng chủ yếu dùng Nginx vì quen thuộc hơn và Nginx vừa làm reverse proxy vừa serve static files. HAProxy có advantage ở layer 4 load balancing và health checking granular hơn, nhưng với scale hiện tại Nginx đáp ứng đủ.

Q27: DNS management — em quản lý DNS records như thế nào?

Dạ, DNS management của em chia 2 phần. External DNS em dùng Cloudflare — quản lý qua Cloudflare dashboard và API. Tất cả public domains trỏ về load balancer server IP. Em dùng Cloudflare proxy mode cho DDoS protection và CDN caching cho static assets. TTL em set 300 giây (5 phút) cho dynamic records và 3600 giây (1 giờ) cho stable records.

Internal DNS thì em dùng cách đơn giản là Ansible-managed /etc/hosts file trên tất cả servers. Mỗi server có entries cho tất cả servers khác với WireGuard IP, ví dụ 10.0.0.1 db-primary, 10.0.0.2 db-standby, 10.0.0.3 monitoring. Khi thêm server mới, Ansible update hosts file trên tất cả nodes. Cách này đơn giản, reliable, không cần maintain DNS server riêng.

Em chưa dùng CoreDNS hay consul DNS cho service discovery — đây là improvement em plan khi scale lớn hơn hoặc chuyển sang Kubernetes. Với Docker Compose thì container name resolution trong cùng network là đủ, cross-server thì dùng WireGuard IPs trong hosts file.

Q28: Security hardening — em áp dụng những gì cho servers?

Dạ, security hardening em áp dụng nhiều layers. SSH security: em disable password authentication, chỉ allow SSH key, disable root login, change SSH port từ 22 sang port custom. Fail2ban em cài trên tất cả servers — ban IP sau 5 failed SSH attempts trong 10 phút, ban duration 1 giờ.

Firewall: em dùng UFW trên mỗi server, default policy deny incoming. Chỉ open ports cần thiết — SSH custom port, WireGuard UDP port. Tất cả application ports chỉ listen trên WireGuard interface, không expose ra public internet. Load balancer server là server duy nhất open port 80 và 443 ra public.

System updates: em configure unattended-upgrades cho security patches tự động. Docker containers em chạy với non-root user khi possible, set read_only: true cho filesystem khi container không cần write, và drop unnecessary Linux capabilities.

Network security: toàn bộ internal traffic encrypt qua WireGuard, services không expose publicly trừ khi cần thiết. Secret management hiện tại em dùng ansible-vault cho infrastructure secrets và environment variables cho application secrets. Em đang work towards HashiCorp Vault cho centralized secret management — đã lab nhưng chưa deploy production.

Em cũng chạy Docker Bench Security periodically để audit Docker configuration theo CIS benchmarks, và review Trivy scan results cho container images trên Harbor.

Q29: Secret management — em quản lý passwords, API keys, certificates như thế nào?

Dạ, hiện tại secret management của em dùng combination of tools. Ansible-vault cho infrastructure secrets — database passwords, WireGuard private keys, API keys cho monitoring — tất cả encrypted trong Ansible vault files, decrypt bằng vault password khi chạy playbooks. Environment variables cho application secrets — mỗi service có .env file trên server, file này không commit vào git, permissions set 600 (chỉ owner đọc được).

Trong CI/CD pipelines, secrets được inject qua pipeline secret variables — Jenkins credentials store, GitLab CI variables, GitHub Actions secrets. Em không bao giờ hardcode secrets trong code hay Dockerfiles.

Certificate management: Let's Encrypt certificates auto-renew bằng certbot, WireGuard keys generate 1 lần và rotate khi cần. Rotation policy em thừa nhận chưa strict — em chưa có automated rotation schedule cho database passwords hay API keys. Đây là gap em đang address bằng cách plan deploy HashiCorp Vault — Vault sẽ centralize secret storage, provide dynamic secrets cho database, và handle automatic rotation.

Access control thì em follow principle of least privilege — mỗi service chỉ có access tới secrets nó cần, không share secrets giữa services. Audit trail hiện tại limited — em biết ai có access gì nhưng chưa có automated logging khi secret được access.

Q30: Có bao giờ bị security incident chưa? Mô tả và cách xử lý.

Dạ, em chưa bị security incident nghiêm trọng nào, nhưng có mấy tình huống đáng kể. Phổ biến nhất là SSH brute force attacks — em thấy trong logs hàng nghìn failed login attempts mỗi ngày từ các IP addresses khác nhau. Fail2ban handle tự động, block IPs ngay. Sau khi thấy volume brute force lớn, em quyết định thay đổi SSH port từ 22 sang custom port — số lượng brute force attempts giảm khoảng 99%.

Có 1 lần em phát hiện 1 container image trên Harbor bị flag critical vulnerability bởi Trivy — CVE liên quan tới OpenSSL trong base image. Em xử lý bằng cách rebuild image với updated base image, re-scan confirm vulnerability fixed, rồi redeploy. Từ đó em set policy trong Harbor block push images có critical CVEs.

Một incident khác là em vô tình expose 1 internal service ra public internet do misconfigure Nginx — dù service không chứa sensitive data nhưng em detect qua access logs thấy external IPs access internal endpoint. Em fix ngay bằng cách update Nginx config để bind service chỉ trên WireGuard interface. Post-incident em review tất cả Nginx configs, thêm alert cho unexpected access patterns. Em cũng thêm vào Ansible playbook 1 task verify rằng internal services không listen trên public interface.

PHẦN 8: AI INFRASTRUCTURE

Q31: OpenClaw agent orchestration — em tự build? Kiến trúc như thế nào?

Đạ đúng, OpenClaw là hệ thống AI agent orchestration em tự build. Kiến trúc gồm mấy components chính. Core là agent runtime — manage agent lifecycle, context window, và tool execution. Mỗi agent có khả năng sử dụng tools qua tool execution engine — gọi bash commands, đọc/ghi files, browse web, gửi messages.

Communication protocol dùng system events — em có thể inject events vào running agent session bằng lệnh `openclaw system event --text "message"`. Đây là cách agents coordinate với nhau — ví dụ Claude Code (em) gửi task cho OpenClaw, OpenClaw thực hiện rồi report kết quả lại.

Capabilities đặc biệt của OpenClaw mà Claude Code không có: web browsing và browser automation, messaging qua Telegram/WhatsApp/Discord, scheduled tasks (cron jobs, reminders), và text-to-speech. Deployment thì OpenClaw chạy containerized với multiple services — gateway service, agent runtime, các MCP (Model Context Protocol) servers. Monitoring qua custom metrics gửi về Prometheus, logs ship về Loki.

Em cũng build shared workspace giữa Claude Code và OpenClaw tại `~/openclaw/workspace/` — inbox cho messages, proposals cho collaboration, research cho research results, và shared memory file. Đây là multi-agent collaboration architecture mà em document trong sách "AI Agents in Production".

Q32: Engram dual-memory system — giải thích chi tiết kiến trúc. Tại sao cần episodic + semantic?

Đạ, Engram là dual-memory system em build cho AI agents, lấy cảm hứng từ cách bộ nhớ con người hoạt động. Kiến trúc gồm 2 memory types. Episodic memory dùng Qdrant vector database — lưu các sự kiện cụ thể, conversations, actions dưới dạng text embeddings. Khi agent cần recall "chuyện gì đã xảy ra", nó query Qdrant bằng semantic search, tìm memories tương tự về ý nghĩa.

Semantic memory dùng PostgreSQL knowledge graph — lưu entities (người, project, concept) và relations giữa chúng (A works-on B, C depends-on D). Khi agent cần answer "tôi biết gì về X", nó query

knowledge graph để lấy structured facts.

Tại sao cần cả hai? Vì chúng phục vụ mục đích khác nhau. Episodic memory cho context — “lần trước user yêu cầu gì, em đã làm gì, kết quả ra sao”. Semantic memory cho knowledge — “server A chạy PostgreSQL, port 5432, connection string là XYZ”. Khi recall, Engram query cả 2 sources và merge results, cho agent cả context lẫn knowledge.

Ingestion pipeline: raw text vào → memory classifier (heuristic, không cần LLM, tiết kiệm cost) phân loại thành todo, decision, preference, error, workflow, lesson, fact → entity extractor rút entities và relations → store vào cả Qdrant (episodic) và PostgreSQL (semantic). Interfaces gồm CLI, MCP server (để Claude Code, Cursor dùng), HTTP API, và WebSocket.

Technical challenges lớn nhất là asyncpg connection pool hang — em fix bằng cách đảm bảo graph query và pool close nằm trong cùng 1 `asyncio.run()` call. Và IPv6 issue với litellm khi connect tới oracle server — em phải force IPv4 trong code.

Q33: Vector database (Qdrant) trong production — deployment, scaling, performance tuning?

Dạ, Qdrant em deploy bằng Docker container với persistent volume mount cho data storage. Collection design — em dùng embedding model generate vectors, mỗi memory entry là 1 point trong collection với vector embedding và metadata payload (timestamp, source, category, agent_id).

HNSW indexing parameters em tune: $m = 16$ (connections per node, balance giữa recall accuracy và memory), $ef_construct = 100$ (build time accuracy). Query performance em đạt latency dưới 50ms cho similarity search trên collection vài chục nghìn points — đủ nhanh cho real-time agent memory recall.

Backup em dùng Qdrant snapshot API — cron job tạo snapshot daily, copy snapshot file tới backup server. Monitoring em expose Qdrant metrics endpoint, scrape bằng Prometheus, track collection size, query latency, memory usage trên Grafana dashboard.

Scaling hiện tại chưa cần vì data volume còn nhỏ (vài GB). Nếu scale lên em sẽ consider Qdrant cluster mode với sharding. Thực tế với use case AI agent memory, data growth rate không cao lắm — mỗi ngày thêm vài trăm memory entries, mỗi entry vài KB, nên single node Qdrant handle tốt.

Q34: MCP (Model Context Protocol) — em hiểu và implement như thế nào?

Dạ, MCP là protocol chuẩn hóa cách AI tools giao tiếp với external services, do Anthropic phát triển. Em implement Engram như 1 MCP server — expose memory operations (remember, recall, ingest, query_graph, etc.) dưới dạng MCP tools. Bất kỳ MCP-compatible client nào — Claude Code, Cursor, hay custom client — đều có thể connect tới Engram MCP server và sử dụng memory operations.

Implementation technical: em dùng stdio transport — MCP client spawn Engram process, communicate qua stdin/stdout bằng JSON-RPC protocol. Mỗi tool có schema definition (name, description, parameters), client discover tools qua `tools/list` method, invoke tool qua `tools/call` method.

Use case thực tế: Claude Code connect tới Engram MCP server, đầu mỗi session gọi `engram_recall` để load context từ sessions trước, trong quá trình làm việc gọi `engram_remember` để lưu thông tin quan trọng. Điều này tạo continuity across sessions — agent “nhớ” được những gì đã làm trước đó.

Em cũng implement MCP resources (read-only data sources) — ví dụ expose memory status, session context dưới dạng MCP resources. Đây là pattern mới và rất promising cho AI agent ecosystem — standardized protocol cho tool integration.

Q35: So sánh self-hosted AI infrastructure vs cloud AI services (AWS Bedrock, Azure OpenAI). Trade-offs?

Dạ, em có perspective thực tế từ cả 2 phía. Self-hosted AI infrastructure (mà em đang làm) có ưu điểm chính là control và customization. Em tự build orchestration layer (OpenClaw), memory system (Engram), và tool integrations theo đúng nhu cầu. Cost model là pay-per-API-call tới LLM providers (Anthropic, OpenAI) + fixed VPS cost cho infrastructure — không có markup từ managed service intermediary.

Cloud AI services như AWS Bedrock hay Azure OpenAI có ưu điểm là managed, scalable, có compliance certifications (SOC2, HIPAA), enterprise support, và tích hợp sẵn với cloud ecosystem. Nhược điểm là vendor lock-in, pricing markup, và ít customization cho orchestration logic.

Approach của em là hybrid: LLM inference dùng cloud API (Anthropic Claude API, không self-host models vì GPU cost quá cao), nhưng orchestration, memory, và tool execution self-hosted. Đây là sweet spot — em get best LLM quality từ cloud providers, nhưng giữ full control over agent behavior, memory, và data.

Khi nào nên chọn gì? Startup/individual: self-hosted orchestration + cloud LLM API (như em đang làm). Enterprise cần compliance: managed cloud services. Large scale cần custom models: self-hosted GPU cluster. Em thấy trend hiện tại là orchestration layer ngày càng quan trọng hơn model layer, vì models commoditize nhanh nhưng orchestration logic là competitive advantage.

PHẦN 9: OUTLIER AI EXPERIENCE

Q36: Top 5% trong 300,000+ contributors — em đạt được cụ thể bằng cách nào?

Dạ, em đạt top 5% trên Outlier AI platform chủ yếu nhờ focus vào chất lượng thay vì số lượng. Platform có accuracy rating system, em maintain 4.5/5 consistently. Chiến lược của em là chọn projects phù hợp với strengths — em focus vào math và STEM projects nơi em có domain knowledge sâu, thay vì spread thin qua nhiều categories.

Cụ thể, em dành thời gian để hiểu evaluation criteria thật kỹ trước khi bắt đầu mỗi project. Mỗi task em review kỹ lưỡng, double-check logic và reasoning trước khi submit. Em cũng consistent — không có gaps

lớn giữa các submissions, maintain steady output. Khi nhận feedback từ reviewers, em analyze patterns — “tại sao task này bị rate thấp?” — rồi adjust approach cho tasks sau.

Kết quả là em được nominated Team Lead position ở rate \$40/hour, có trách nhiệm review work của contributors khác, mentor newcomers, và suggest process improvements. Kinh nghiệm này tuy không phải DevOps trực tiếp nhưng cho em skills quan trọng: attention to detail, quality assurance mindset, working in distributed team (300K+ contributors globally), và communication qua text-based platforms.

Q37: AI training tasks — RLHF, SFT, Evals, Adversarial Training — em đã làm những gì?

Dạ, trong Outlier AI platform em đã tham gia nhiều types of AI training tasks. RLHF (Reinforcement Learning from Human Feedback) — em rank multiple AI responses theo quality, chọn response nào tốt hơn và giải thích tại sao. Data này được dùng để train reward model cho LLMs. SFT (Supervised Fine-Tuning) — em tạo high-quality training examples, viết cả prompt và ideal response, đảm bảo response accurate, helpful, và well-structured.

Evals (Evaluations) — em evaluate model outputs theo multiple dimensions: correctness, helpfulness, safety, formatting. Cho điểm từng dimension và provide detailed justification. Adversarial Training — em cố tìm weaknesses trong model responses, craft prompts khiến model trả lời sai hoặc unsafe, giúp identify failure modes.

Đặc biệt, thành tựu chính của em ở math/STEM projects — em verify mathematical proofs, check computation correctness, evaluate step-by-step reasoning. Đây là nơi engineering background giúp nhiều — systematic thinking, attention to precision. Em cũng contribute vào coding tasks — review AI-generated code, check logic errors, verify test coverage. Kinh nghiệm này cho em hiểu sâu hơn về cách LLMs hoạt động, strengths và limitations — knowledge rất useful khi build AI infrastructure.

Q38: Nominated Team Lead \$40/h — leadership experience gì từ đây?

Dạ, việc được nominated Team Lead cho em leadership experience trong mấy areas. Thứ nhất là quality review — em review work submissions từ contributors khác, provide constructive feedback giúp họ cải thiện. Em phải balance giữa maintain high quality standards và không quá harsh làm discourage contributors.

Thứ hai là mentoring — em guide newcomers qua evaluation criteria, share tips và best practices, help họ ramp up nhanh hơn. Communication chủ yếu qua text-based channels (platform messaging, documentation), nên em develop kỹ năng viết clear, concise instructions.

Thứ ba là process improvement — em observe patterns trong common mistakes, suggest changes to guidelines và evaluation rubrics giúp improve overall quality. Ví dụ em propose thêm check-steps cho math tasks, giảm error rate.

Kinh nghiệm này transferable sang DevOps team environment: review code/infrastructure changes, mentor junior team members, improve processes và runbooks. Tuy scale và context khác — distributed AI training platform vs DevOps team — nhưng core leadership skills tương tự.

PHẦN 10: CAREER TRANSITION & EDUCATION

Q39: Chuyển từ CNC Machine Operator sang DevOps — tại sao? Lộ trình tự học?

Dạ, em bắt đầu career là CNC Machine Operator tại Mộc Phát Furniture từ 2019 tới 2023. Trong thời gian đó em kiêm luôn IT support cho công ty — setup computers, troubleshoot network, maintain internal systems. Từ đó em phát hiện passion thực sự là technology, đặc biệt là automation. CNC machine thực chất cũng là automation — em viết G-code programs để máy tự động cắt theo specifications chính xác. Mindset “automate repetitive tasks with precision” đó transfer trực tiếp sang DevOps.

Lộ trình tự học của em từ 2019-2023: em bắt đầu với Linux basics (Ubuntu on personal laptop), rồi Docker (containerize personal projects), rồi networking fundamentals, rồi CI/CD concepts. Em hoàn thành AWS DevOps Specialization trên Coursera — học EC2, S3, CloudFormation, CodePipeline. Tiếp theo IBM DevOps and Software Engineering Professional Certificate — học Docker deep dive, Kubernetes, Agile methodology. Cùng lúc em build homelab — setup VPS, deploy services, break things, fix them, learn from mistakes.

2024 em quyết định transition full-time. Em build 10-server infrastructure from scratch, áp dụng tất cả kiến thức tích lũy 4 năm. Kết quả là infrastructure production-grade với full observability, CI/CD, HA database, VPN mesh — tất cả self-taught. Key insight mà em rút ra: CNC Machine Operator và DevOps Engineer cùng core philosophy — automation, precision, process optimization, và continuous improvement. Em không bắt đầu từ zero, em transfer engineering mindset sang domain mới.

Q40: Bằng Automotive Engineering Technology (HCMUTE) — có gì transferable sang DevOps?

Dạ, bằng Công nghệ Kỹ thuật Ô tô tại HCMUTE (2015-2019) cho em nền tảng engineering mindset mà em thấy rất transferable. Systematic thinking — khi diagnose sự cố xe hơi, em phải follow systematic process: check symptoms, narrow down possible causes, test hypotheses, find root cause. Đây chính xác là cách em debug infrastructure issues — check metrics, check logs, isolate components, find root cause.

Precision và attention to detail — CNC machine operation đòi hỏi chính xác tới 0.01mm, 1 sai sót nhỏ có thể hỏng cả batch sản phẩm. Trong DevOps, 1 typo trong config file có thể down production. Mindset “double-check everything” từ CNC chuyển sang DevOps rất tự nhiên.

Process optimization — trong manufacturing em tối ưu workflow sản xuất, giảm waste, tăng throughput. Trong DevOps em tối ưu CI/CD pipeline, giảm build time, tăng deployment frequency. Documentation — trong automotive em đọc và viết technical drawings, specifications. Trong DevOps em viết runbooks, architecture docs, operational procedures.

Tuy bằng không directly related tới IT, nhưng engineering foundation cho em analytical thinking, problem-solving methodology, và discipline mà many self-taught developers không có.

Q41: AWS DevOps Specialization + IBM DevOps Certificate — học được gì? Áp dụng thực tế?

Dạ, AWS DevOps Specialization trên Coursera cho em foundational knowledge về AWS ecosystem — EC2 instances, S3 storage, VPC networking, CloudFormation IaC, CodePipeline CI/CD, CloudWatch monitoring. Em hiểu concepts như Auto Scaling Groups, Load Balancers, RDS managed databases, IAM roles và policies. Kiến thức này giúp em khi compare self-hosted vs cloud costs (Q2), và khi discuss cloud architecture design (Q60).

IBM DevOps and Software Engineering Professional Certificate đi sâu hơn vào Docker containerization, Kubernetes orchestration, CI/CD pipelines, Agile và Scrum methodology, microservices architecture. Hands-on labs cho em practice trên real environments.

Áp dụng thực tế: concepts từ courses em translate sang self-hosted infrastructure. VPC networking → WireGuard VPN design. Auto Scaling → manual but planned capacity management. CloudWatch → Prometheus/Grafana stack. CodePipeline → Jenkins/GitHub Actions pipelines. RDS Multi-AZ → PostgreSQL Patroni HA. Containers and orchestration concepts → Docker Compose trong production.

Gap thành thật: em có theoretical knowledge về AWS services nhưng chưa có hands-on production experience quản lý AWS infrastructure ở scale lớn. Em biết navigate AWS Console, understand pricing models, design architectures on paper — nhưng chưa manage production workloads trên AWS. Đây là area em đang actively bridge bằng cách build side projects trên AWS Free Tier.

PHẦN 11: CODING & SCRIPTING

Q42: Python, Go, Bash — em dùng từng ngôn ngữ cho mục đích gì trong DevOps?

Dạ, mỗi ngôn ngữ em dùng cho mục đích khác nhau dựa trên strengths của chúng. Bash là ngôn ngữ em dùng nhiều nhất cho daily DevOps tasks — automation scripts, cron jobs, system administration. Ví dụ backup scripts (pg_dump + rsync + rotation + Telegram notification khi fail), server health check scripts, Docker cleanup scripts, deployment scripts. Bash phù hợp vì nó native trên Linux, không cần install dependencies, chạy nhanh cho simple tasks.

Python em dùng cho complex automation và application development. Engram memory system được viết bằng Python — asyncio cho async operations, asyncpg cho PostgreSQL, aiohttp cho HTTP API, Click cho CLI framework. Em cũng viết Python scripts cho monitoring custom metrics, log analysis, và data processing tasks. Python phù hợp khi cần xử lý complex logic, work với APIs, hoặc build tools có structure.

Go em dùng cho performance-critical CLI tools và services. Go compile ra single binary, no dependencies, rất phù hợp cho DevOps tools cần deploy trên nhiều servers mà không cần install runtime.

OpenClaw có components viết bằng Go cho performance reasons. Em cũng viết small utilities bằng Go khi cần concurrency và low resource usage.

Trong practice, khoảng 50% automation work em dùng Bash, 40% Python, 10% Go. Principle là chọn tool phù hợp cho job — simple tasks dùng Bash, complex logic dùng Python, performance-critical dùng Go.

Q43: Cho ví dụ 1 automation script em viết giải quyết vấn đề thực tế?

Dạ, em cho ví dụ script backup automation mà em thực sự dùng daily. Problem: em cần backup PostgreSQL databases trên nhiều servers, rotate old backups, rsync tới backup server, và alert qua Telegram nếu backup fail. Trước khi có script, em phải SSH vào từng server chạy pg_dump manually — tốn thời gian và dễ quên.

Solution: em viết Bash script kết hợp Ansible. Script chạy trên mỗi database server qua cron (2AM daily): pg_dump compress database ra file .sql.gz với timestamp trong filename, check exit code — nếu fail thì gửi Telegram alert ngay bằng curl tới Telegram Bot API. Nếu success thì tiếp tục: xóa backup files cũ hơn 7 ngày (daily retention), giữ lại weekly backups (chủ nhật) trong 4 tuần, rsync backup file tới remote backup server qua WireGuard tunnel. Cuối cùng log result vào file để monitoring.

Impact: backup chạy tự động 100%, em không cần nhớ hay manually trigger. Khi có incident 1 lần database cần restore, em có backup file ready trong vòng 5 phút. Telegram alert đã catch 2 lần backup fail do disk space issue — em fix trước khi ảnh hưởng production. Script chạy ổn định hơn 1 năm, maintenance minimal.

Q44: PHP và TypeScript trong CV — em dùng cho gì?

Dạ, PHP em dùng chủ yếu trong giai đoạn đầu khi em làm web development. Wellington Decorators project backend ban đầu em xem xét PHP-based CMS trước khi quyết định dùng Strapi (Node.js). Em cũng có kinh nghiệm với WordPress và Laravel từ learning phase. PHP skills giúp em khi cần debug hoặc maintain legacy PHP applications — trong DevOps role đôi khi phải understand application code để troubleshoot deployment issues.

TypeScript em dùng nhiều hơn. GitPocket mobile app (React Native) viết hoàn toàn bằng TypeScript — đây là GitHub client app cho mobile với AI agent integration. Web projects em cũng dùng TypeScript cho frontend (React, Astro). TypeScript cho em type safety và better code organization cho larger projects.

Trong context DevOps, coding skills này rất valuable vì em understand application code — khi developer report bug hay performance issue, em có thể đọc code, hiểu logic, identify bottlenecks, không chỉ look at infrastructure layer. Nhiều DevOps engineers chỉ biết infrastructure tools mà không đọc được application code — em thấy đây là advantage của mình. Em có thể viết custom monitoring scripts, build internal tools, contribute fixes cho application code khi needed.

PHẦN 12: PROJECTS

Q45: Engram project – technical challenges lớn nhất khi build? Cách giải quyết?

Dạ, Engram có mấy technical challenges lớn mà em phải solve. Challenge lớn nhất là asyncpg connection pool hang. Khi em query PostgreSQL knowledge graph rồi close pool ở 2 asyncio.run() calls khác nhau, pool bị hang indefinitely. Root cause là asyncpg pool phải được create, use, và close trong cùng 1 event loop. Solution: em restructure code để graph query và pool close nằm trong cùng 1 asyncio.run() call. Debug mất 2 ngày nhưng learn rất sâu về Python asyncio internals.

Challenge thứ hai là entity-first ingestion. Ban đầu em ingest raw text vào vector store, nhưng retrieval quality thấp vì text chunks không meaningful. Em redesign pipeline: extract entities và relations trước (structured data), rồi mới store vào both vector store (episodic) và knowledge graph (semantic). Dùng /api/v1/ingest endpoint mới thay vì direct storage.

Challenge thứ ba là LLM cost cho memory operations. Ban đầu em dùng LLM để classify memories (todo, decision, preference, etc.) — mỗi ingest operation tốn API call. Solution: em viết heuristic classifier bằng regex patterns — detect keywords và patterns trong text để classify mà không cần LLM. Accuracy không perfect bằng LLM nhưng đủ tốt cho 90%+ cases, và cost giảm về 0 cho classification step.

Challenge thứ tư là IPv6 issue — litellm library default connect qua IPv6 tới oracle server, nhưng IPv6 routing bị hang. Em phải force IPv4 bằng cách patch trong code. Đây là production issue chỉ xuất hiện khi deploy, không thấy khi develop locally.

Q46: GitPocket (React Native) – tại sao build? Tech stack? Challenges?

Dạ, GitPocket là GitHub mobile client em build bằng React Native + TypeScript. Motivation là em muốn access GitHub repositories, review code, manage issues trên mobile phone khi không ở trước laptop. GitHub official mobile app tốt nhưng em muốn integrate AI agent capabilities — ví dụ ask AI to explain code, summarize PR changes, suggest fixes.

Tech stack: React Native cho cross-platform (Android released, iOS coming), TypeScript cho type safety, GitHub REST và GraphQL APIs cho data. UI dùng React Native components, navigation dùng React Navigation. Authentication qua GitHub OAuth.

Challenges chính: GitHub API rate limiting — authenticated users get 5,000 requests/hour, em implement caching strategy để minimize API calls. Rich code rendering trên mobile — display syntax-highlighted code với horizontal scroll, line numbers, diff views trên small screen không dễ. Em dùng WebView component cho code rendering. Offline support em implement bằng local caching — user có thể browse recently viewed repos offline.

Project này tuy không phải DevOps trực tiếp nhưng demonstrate em có full-stack development capability — mobile app development, API integration, CI/CD cho mobile builds (GitHub Actions), app store deployment process.

Q47: OmniPocket AI Marketing Platform – architecture overview?

Dạ, OmniPocket là AI Marketing Platform em đang develop. Concept là giúp businesses automate marketing tasks bằng AI agents — content generation, campaign management, analytics, multi-channel publishing. Architecture gồm backend API (Python/FastAPI), frontend dashboard (React/TypeScript), AI agent layer integrate với LLMs cho content generation và analysis.

Deployment trên self-hosted infrastructure của em — Docker Compose, behind Nginx reverse proxy, monitored bằng full observability stack. Database PostgreSQL cho structured data, Redis cho caching và job queues. AI features connect tới cloud LLM APIs (Anthropic, OpenAI) cho content generation, summarization, sentiment analysis.

Project này demonstrate em apply DevOps skills vào real product development — infrastructure design, CI/CD pipeline, monitoring, database management, security — tất cả cho 1 production application. Đồng thời nó cũng showcase AI infrastructure skills — integrate LLMs, manage API costs, handle rate limiting, build agent workflows.

Q48: Cuốn sách “AI Agents in Production” — 16 chapters cover những gì? Tại sao viết?

Dạ, em viết cuốn “AI Agents in Production” vì em thấy gap lớn giữa AI agent tutorials online (thường là toy examples) và reality khi deploy agents trong production. Hầu hết resources nói về building demo agents, nhưng không cover production concerns như reliability, cost management, memory persistence, multi-agent coordination, security.

16 chapters cover: agent architecture fundamentals, context engineering (manage token limits, context window optimization), memory systems (episodic vs semantic, vector stores, knowledge graphs), tool use and function calling, multi-agent orchestration patterns, error handling and reliability, cost management strategies, security considerations (prompt injection, data leaks), monitoring and observability cho AI systems, và deployment best practices.

Content dựa hoàn toàn trên hands-on experience building OpenClaw và Engram — real production issues, real solutions, real war stories. Ví dụ chapter về memory systems detail asyncpg pool hang fix, heuristic classifier design, entity-first ingestion pipeline — tất cả từ actual Engram development. Chapter về multi-agent orchestration describe Claude Code + OpenClaw collaboration patterns em use daily.

Available trên Leanpub bằng cả English và Vietnamese. Differentiation so với sách khác: em không nói lý thuyết, em share actual production experience với code examples từ real systems. Em coi việc viết sách là cách systematize knowledge, và cũng build professional credibility.

PHẦN 13: SERVICES (từ CV Services section)

Q49: Server Configuration service — em setup một server từ đầu quy trình như thế nào?

Dạ, quy trình em setup server mới từ đầu gồm mấy phases. Phase 1 — OS Installation: chọn Ubuntu 22.04 LTS (em prefer Ubuntu vì ecosystem lớn, LTS cho stability). Provision VPS trên Hetzner, chọn specs phù hợp workload.

Phase 2 — Base Hardening: em chạy Ansible base role. Tạo admin user (không dùng root), copy SSH public key, disable password authentication và root login trong `sshd_config`, change SSH port. Install và configure UFW firewall — default deny incoming, allow SSH custom port, allow WireGuard port. Install fail2ban cho SSH brute force protection. Configure `unattended-upgrades` cho automatic security patches. Set timezone, configure NTP time sync.

Phase 3 — Docker Setup: install Docker Engine và Docker Compose từ official repository (không dùng OS package vì outdated). Configure Docker daemon — logging driver `json-file` với `max-size 50MB` và `max-file 3` để prevent disk fill. Configure Docker to listen only on localhost và WireGuard interface.

Phase 4 — Monitoring: install `node_exporter` (system metrics), `cAdvisor` (container metrics), `Promtail` (log shipping). Configure `Promtail` gửi logs về Loki central server. Update Prometheus config add new target.

Phase 5 — Network: generate WireGuard config, setup VPN tunnel tới cluster. Update `/etc/hosts` với internal DNS entries.

Phase 6 — Application: deploy application containers bằng Docker Compose. Configure Nginx reverse proxy cho web-facing services.

Phase 7 — Verification: check tất cả metrics hiện trên Grafana, verify alerts configured, test SSH access qua VPN, test application endpoints. Document server specs, role, và deployed services.

Q50: System Administration — em quản lý users, access control trên multi-server environment ra sao?

Dạ, user management trên 10 servers em quản lý centralized qua Ansible. Em có 1 Ansible role `users` define tất cả user accounts, SSH keys, sudo policies, và shell configurations. Khi cần thêm user mới, em add vào role config rồi chạy playbook — user được tạo consistent trên tất cả servers (hoặc subset servers tùy group).

SSH key management: mỗi user có SSH public key stored trong Ansible vars. Khi user thêm hoặc thay key, update Ansible rồi run playbook. Không ai dùng password authentication — SSH key only. Key rotation khi member rời team thì em remove key, chạy playbook, access bị revoke trên tất cả servers ngay.

Sudo policies: principle of least privilege. Admin users có limited sudo commands — chỉ được restart Docker containers, view logs. Full sudo chỉ cho em (infra admin). Em define `sudoers` config trong Ansible, deploy consistently.

Hiện tại với 10 servers và team nhỏ thì approach Ansible-based này đủ. Nếu scale lên em sẽ consider centralized auth solutions như FreeIPA hoặc LDAP. Em cũng log SSH access qua `syslog`, ship về Loki để audit ai access server nào, khi nào.

PHẦN 14: BEHAVIORAL &

SITUATIONAL

Q51: Em làm freelance/self-hosted — chưa từng làm trong team DevOps lớn. Làm sao adapt vào corporate environment?

Dạ, em hiểu concern này hoàn toàn valid. Em chưa từng làm trong team DevOps lớn, nhưng em tin mình adapt được vì mấy lý do. Thứ nhất, tools và concepts em dùng daily là industry standard — Docker, Prometheus, Grafana, Ansible, CI/CD pipelines, PostgreSQL HA — đây là exact same tools dùng trong corporate environments. Khác biệt chủ yếu ở scale và process, không phải technology.

Thứ hai, em có team collaboration experience từ Outlier AI — làm việc trong community 300,000+ contributors, review work của người khác, receive và give feedback, follow standardized processes. Tuy không phải DevOps team nhưng transferable skills: async communication, quality review, process adherence.

Thứ ba, em có self-direction rất mạnh — 4 năm tự học, tự build production infrastructure, tự troubleshoot, tự document. Trong corporate environment điều này translate thành: không cần micromanagement, proactive identify và solve problems, tự research solutions trước khi escalate.

Những gì em cần learn trong corporate environment: change management processes (change advisory boards, maintenance windows), on-call rotations và escalation procedures, team code review workflows, compliance và audit requirements. Em sẵn sàng learn và adapt — learning velocity của em đã proven qua career transition từ CNC tới DevOps trong 4 năm.

Q52: Weakness lớn nhất của em trong DevOps là gì?

Dạ, em thành thật chia sẻ mấy weaknesses. Weakness lớn nhất là limited cloud production experience. Em hiểu AWS concepts qua Coursera courses và self-study, nhưng chưa manage production workloads trên AWS/GCP/Azure ở enterprise scale. Em chủ yếu self-hosted trên VPS, chưa dùng managed services như EKS, RDS, Lambda trong production. Em đang actively bridge gap này bằng cách build side projects trên cloud platforms.

Weakness thứ hai là English communication. IELTS 4.5 — em comfortable đọc technical documentation và viết code comments, git commits bằng English. Nhưng speaking và listening trong meetings hay technical discussions em cần improve. Em đang practice daily — đọc sách technical bằng English, viết book “AI Agents in Production” bằng English, watch technical talks. Đây là area em committed improve continuously.

Weakness thứ ba là Kubernetes production experience. Em biết concepts, đã lab trên local cluster, nhưng chưa operate K8s cluster trong production. Production em dùng Docker Compose. Nếu role yêu cầu K8s operations em cần ramp-up time.

Em tin rằng nhận biết weaknesses và có plan address chúng quan trọng hơn pretend không có weaknesses. Cloud experience em plan bridge bằng AWS projects, English em improve daily, K8s em plan migrate gradually khi scale justify.

Q53: Mô tả 1 incident phức tạp nhất em từng xử lý. Root cause? Resolution time?

Dạ, incident phức tạp nhất em xử lý là PostgreSQL primary server OOM crash lúc khoảng 11 giờ đêm. Timeline: em nhận Telegram alert InstanceDown cho database server. SSH vào thì thấy server vẫn online nhưng PostgreSQL process bị OOM killer terminate — check dmesg thấy “Out of memory: Kill process postgresql”.

Diagnosis: check memory usage history trên Grafana, thấy memory tăng dần suốt 3 ngày trước incident. Root cause: 1 application mới deploy có bug tạo quá nhiều database connections — mỗi connection tốn ~10MB RAM, và không có connection pooling, tích tụ tới hàng trăm idle connections chiếm hết RAM.

Resolution: em restart PostgreSQL process, nó recover nhưng sau 30 phút lại bị OOM vì application vẫn tạo connections mới. Em phải stop application container gây vấn đề, restart PostgreSQL, set max_connections thấp hơn tạm thời. Hôm sau em deploy PgBouncer connection pooler đứng trước PostgreSQL, config pool_size = 50, application connect qua PgBouncer thay vì direct tới PostgreSQL. Đồng thời notify developer fix connection leak trong application code.

Resolution time: khoảng 2 giờ cho immediate fix (stop app, restart DB), thêm 4 giờ hôm sau cho proper fix (deploy PgBouncer). Post-mortem: em thêm alert cho PostgreSQL connection count (warn > 100, critical > 200), thêm Grafana dashboard panel track connections per application, và document runbook cho OOM incidents. Từ đó em mandate tất cả applications phải connect qua PgBouncer, không direct tới PostgreSQL.

Q54: Khi có quá nhiều tasks cùng lúc — server down, deployment failed, urgent feature request — em prioritize thế nào?

Dạ, em prioritize theo severity và impact. Production down luôn là P0, highest priority — nếu server chạy critical services bị down, mọi thứ khác đều dừng lại để fix. Deployment failed là P1 — quan trọng nhưng không ảnh hưởng current users nếu rollback kịp. Feature request dù urgent cũng chỉ P2 — business can wait, production cannot.

Cụ thể khi gặp situation 3 tasks cùng lúc: em handle server down trước — quick triage, nếu có runbook thì follow, nếu cần 30+ phút thì em communicate status cho stakeholders. Trong khi đó nếu deployment failed, em rollback về version trước (mất 1-2 phút) để stabilize. Feature request em acknowledge receipt và schedule cho sau khi incidents resolved.

Communication là key — em immediately notify relevant people về situation: “DB server down, đang investigate. Deployment rolled back. Feature request sẽ handle sau khi resolve incidents.” Không ai expect em handle 3 tasks simultaneously, nhưng expect em communicate clearly about priorities và timeline.

Nếu có team, em sẽ delegate — 1 người handle server incident, 1 người investigate deployment failure. Hiện tại solo nên em phải sequential, nhưng principle là always stabilize production first, then investigate, then feature work.

Q55: Em có kinh nghiệm on-call không? Maintain SLA như thế nào?

Dạ, vì em self-host toàn bộ infrastructure nên effectively em on-call 24/7. Bất cứ khi nào critical alert fire — dù 2 giờ sáng hay weekend — em nhận Telegram notification và respond. Đây là reality của self-hosted infrastructure khi chưa có team rotate.

Response time cho critical alerts em target dưới 15 phút — từ lúc nhận alert tới lúc bắt đầu investigate. Cho actual resolution thì depend vào issue complexity — simple restart mất 5 phút, complex debugging có thể mất vài giờ. Em track metrics: average response time, MTTR (Mean Time To Resolve), number of incidents per month trên Grafana.

Cách em maintain reliability: preventive monitoring — alerts fire trước khi issue trở thành critical (disk warning ở 80% trước khi full ở 100%). Runbooks — em document procedures cho common issues, nên khi alert fire lúc 2AM mà brain chưa tỉnh, em follow runbook step-by-step thay vì phải think. Automation — self-healing scripts cho simple issues, ví dụ auto-restart container khi health check fail (Docker restart policy unless-stopped).

SLA thực tế infrastructure em đạt khoảng 99.5% uptime — downtime chủ yếu từ planned maintenance và occasional incidents. Em chưa có formal SLA vì chưa có external customers, nhưng em treat own infrastructure với same discipline. Nếu join corporate team, em ready adapt vào structured on-call rotation, escalation procedures, và formal SLA commitments.

PHẦN 15: TECHNICAL DEEP-DIVE

Q56: Service mesh — em ghi Istio trong CV. Em dùng Istio hay chỉ biết concept?

Dạ, em thành thật nói em chủ yếu biết Istio ở mức concept, chưa deploy trong production. Em hiểu Istio architecture: data plane (Envoy sidecar proxies inject vào mỗi pod), control plane (istiod manage config, certificates, service discovery). Features em nắm: mTLS (mutual TLS encryption giữa services), traffic management (canary deployments, traffic splitting, circuit breaking), observability (automatic metrics, distributed tracing, access logging).

Lý do chưa dùng production: Istio cần Kubernetes, và em chưa chạy K8s trong production. Complexity overhead của Istio cũng significant — thêm sidecar proxy cho mỗi pod tốn resource, và operational complexity tăng đáng kể.

Tuy nhiên, WireGuard mesh VPN mà em đang dùng cover một phần chức năng của service mesh — encryption in transit (WireGuard encrypt tất cả traffic giữa nodes), network segmentation. Em thiếu traffic management features (canary, circuit breaking) và automatic mTLS cho container-to-container traffic trên cùng host. Nếu em migrate sang K8s trong tương lai, Istio hoặc alternatives như Linkerd sẽ là natural next step.

Q57: GitOps — em hiểu và áp dụng GitOps principles như thế nào?

Dạ, GitOps principle chính là Git as single source of truth cho infrastructure state. Em áp dụng điều này: tất cả Ansible playbooks, Docker Compose files, Nginx configs, monitoring configs, alert rules

— đều commit trong git repositories. Khi em cần change infrastructure, em update config trong git, commit, rồi apply bằng Ansible hoặc docker compose up.

Tuy em chưa dùng ArgoCD hay Flux (cần K8s), approach của em follow GitOps spirit: infrastructure changes reviewable qua git diff, rollback bằng git revert, audit trail qua git log, reproducible state vì tất cả config declarative trong git.

CI/CD pipeline enforce điều này: code merge vào main → pipeline auto-deploy. Không manual changes trên server — mọi thứ phải through git commit. Có exceptions cho emergency hotfixes, nhưng sau đó em luôn backport changes vào git để maintain source of truth.

Nếu em migrate sang Kubernetes, em sẽ implement full GitOps với ArgoCD: Kubernetes manifests trong git → ArgoCD watch repo → auto-sync cluster state match git state. Reconciliation loop detect drift và auto-fix. Đây là natural evolution từ current approach.

Q58: Incident response & On-call operations — quy trình cụ thể?

Dạ, quy trình incident response của em structured theo mấy phases. Phase 1 — Detection: Prometheus monitoring continuously, Alertmanager evaluate rules mỗi 15-30 giây. Khi condition match (ví dụ node_exporter down > 3 phút), Alertmanager gửi notification qua Telegram bot. Em nhận instant notification trên phone.

Phase 2 — Triage: em assess severity trong 5 phút đầu. P0: production service unavailable, ảnh hưởng users. P1: degraded performance, partial outage. P2: non-critical service issue, no user impact. Severity determine response urgency — P0 drop everything, P2 schedule for business hours.

Phase 3 — Diagnosis: SSH vào affected server, check system state — journalctl -u docker, docker ps -a, dmesg | tail, df -h, free -m, top. Đồng thời check Grafana dashboards cho historical context — khi nào metrics bắt đầu abnormal? Correlate với recent changes — deployment, config change, traffic pattern.

Phase 4 — Resolution: apply fix, verify service restored, monitor 15-30 phút ensure stability. Nếu fix nhanh không khả thi, apply temporary workaround (ví dụ restart, rollback) rồi schedule proper fix.

Phase 5 — Post-mortem: em viết brief document: timeline, root cause, resolution, action items to prevent recurrence. Action items thường là: thêm monitoring/alert, update runbook, fix underlying code bug, improve automation.

Communication xuyên suốt: em update status qua Telegram khi start investigating, khi identify root cause, và khi resolve. Transparency về progress quan trọng hơn promise timeline.

Q59: Backup, restore & disaster recovery — em đã test restore thành công chưa?

Dạ, em đã test restore thành công cho individual components nhưng chưa test full cluster DR.

PostgreSQL restore: em đã test multiple lần. Quy trình: take pg_dump backup file → create new empty database → pg_restore → verify row counts match, run application smoke tests. Last tested khoảng 2 tháng

trước, restore 5GB database mất khoảng 15 phút. Data integrity verified bằng compare row counts per table giữa production và restored database.

Server rebuild từ Ansible: em đã test bằng cách provision VPS mới, chạy full Ansible playbook suite — server setup, Docker install, WireGuard config, monitoring agents. Kết quả: server fully configured trong 30-40 phút, ready deploy applications. Ansible playbooks idempotent nên reliable.

Application recovery từ Docker images: em test bằng cách pull images từ Harbor, docker compose up, verify services healthy. Vì images immutable và configs trong git, application recovery straightforward.

Chưa test: full cluster DR — rebuild tất cả 10 servers simultaneously, restore all data, verify full system functionality. RTO target là 4-6 giờ cho critical services, nhưng chưa validate. RPO là 24 giờ (daily backup cycle). Em plan quarterly DR drill: chọn 1 non-critical server, wipe clean, rebuild from scratch + restore data, measure actual RTO. Đây là gap em acknowledge và actively planning to address.

Q60: Nếu được giao redesign toàn bộ infrastructure lên cloud (AWS), em sẽ thiết kế như thế nào?

Dạ, nếu migrate lên AWS em sẽ thiết kế như sau. Networking: VPC với 3 Availability Zones, mỗi AZ có public subnet (load balancer, NAT gateway) và private subnet (applications, databases). Security Groups thay firewall rules, NACLs cho subnet-level security.

Compute: EKS (Elastic Kubernetes Service) cho container orchestration — tận dụng managed control plane, em chỉ manage worker nodes. Node groups cho different workload types: general purpose cho app containers, memory-optimized cho database, GPU instances cho AI workloads nếu cần.

Database: RDS PostgreSQL Multi-AZ cho automatic failover — thay thế self-hosted Patroni. ElastiCache Redis cho caching — managed, auto-failover, multi-AZ.

Monitoring: em sẽ keep self-hosted Prometheus/Grafana stack trên EKS thay vì dùng CloudWatch — vì em đã có dashboards, alert rules, và expertise. CloudWatch dùng cho AWS-specific metrics (RDS, ALB, etc.) feed vào Grafana.

CI/CD: GitHub Actions cho build + test, ArgoCD trên EKS cho GitOps deployment. ECR thay Harbor cho container registry (native EKS integration).

Security: IAM roles cho service-to-service authentication (không dùng static credentials), KMS cho encryption at rest, Secrets Manager thay ansible-vault, WAF trước ALB.

Cost optimization: Reserved Instances cho baseline workloads (predictable savings 30-40%), Spot Instances cho CI/CD build agents và non-critical workloads, S3 Intelligent-Tiering cho backups.

Migration strategy: phased approach — Phase 1: non-critical services → Phase 2: CI/CD → Phase 3: monitoring → Phase 4: databases (most risky, last). Blue-green migration: run parallel on both old và new infrastructure, gradually shift traffic.

PHẦN 16: WELLINGTON DECORATORS (Freelance Project)

Q61: Wellington Decorators — em triển khai web và infrastructure cho client này cụ thể như thế nào?

Dạ, Wellington Decorators là freelance project em làm cho 1 công ty trang trí nội thất tại New Zealand. Em triển khai full stack: website + hosting + SEO. Website em build bằng Astro (static site generator) cho frontend — fast, SEO-friendly, perfect Lighthouse scores — kết hợp Strapi headless CMS cho backend content management. Client có thể tự update portfolio, services, contact info qua Strapi admin panel mà không cần technical knowledge.

Infrastructure: em host trên 1 VPS riêng, Nginx reverse proxy terminate SSL (Let's Encrypt), Strapi chạy trong Docker container, Astro build output served as static files bởi Nginx. Setup đơn giản nhưng reliable — uptime gần 100% vì ít moving parts. CI/CD: khi em push code lên git, pipeline auto-build Astro static site, deploy lên server.

Monitoring: basic nhưng đủ — node_exporter cho server metrics, uptime check alerts qua Telegram. Backup: Strapi database (SQLite) backup daily, configs trong git. Maintenance: em manage ongoing server administration, security updates, SSL renewal.

Project này demonstrate em deliver end-to-end: từ development, deployment, hosting, tới ongoing maintenance cho real client. Và achieve measurable business result — Google My Business ranking #1.

Q62: Google My Business #1 ranking — em đạt được bằng cách nào?

Dạ, em đạt Google My Business #1 ranking cho Wellington Decorators bằng combination of technical SEO và local SEO strategies. Về technical SEO: website build bằng Astro nên performance rất tốt — Lighthouse scores gần 100 cho Performance, Accessibility, Best Practices, SEO. Page load time dưới 1 giây. Mobile-responsive design vì Google prioritize mobile-first indexing. Proper HTML semantic structure, meta tags, Open Graph tags.

Schema markup em implement LocalBusiness structured data — giúp Google hiểu đây là local business ở Wellington, services nào cung cấp, contact info, operating hours. Sitemap.xml và robots.txt configured properly.

Google My Business optimization: em help client optimize GMB profile — complete tất cả fields, upload high-quality portfolio photos thường xuyên, respond to customer reviews, post regular updates (GMB posts). Encourage satisfied customers để lại reviews — số lượng và quality reviews ảnh hưởng lớn ranking.

Local SEO: keywords targeting “decorator Wellington”, “house painter Wellington”, “interior painting Wellington NZ”. Content pages cho từng service area. NAP (Name, Address, Phone) consistency across web. Link building từ local directories.

Result: website rank #1 trên Google Maps pack cho several keywords liên quan tới painting/decorating services in Wellington. Client nhận thêm leads đáng kể qua Google Search. Đây là case study em proud nhất về measurable business impact từ technical work.

PHẦN 17: IELTS & COMMUNICATION

Q63: IELTS score? English communication level cho technical discussions?

Dạ, IELTS score của em là 4.5, em không giấu con số này. Tuy nhiên em muốn chia sẻ context. Technical English và conversational English khác nhau nhiều. Với technical English, em rất comfortable — em đọc documentation, Stack Overflow, GitHub issues, technical blogs bằng English hàng ngày. Code comments, git commit messages, README files em viết bằng English. Em viết cuốn sách “AI Agents in Production” bằng English, 16 chapters. Technical vocabulary em nắm vững — em có thể discuss Docker, Kubernetes, CI/CD, monitoring, networking bằng English mà không thiếu từ.

Gap chính là conversational English — speaking fluency và listening comprehension trong real-time discussions. Em đọc hiểu nhanh nhưng respond bằng lời nói cần thời gian hơn. Trong meeting em có thể follow along nhưng contribute ý kiến bằng English chưa tự nhiên.

Improvement plan: em practice daily — nghe technical podcasts, watch conference talks, đọc sách technical bằng English. Em cũng tương tác với AI assistants bằng English daily qua work. Em commit rằng IELTS score sẽ improve qua practice — và trong technical role, writing và reading proficiency (em mạnh) quan trọng hơn speaking fluency. Technical discussions có advantage là vocabulary limited và specific, dễ improve hơn general conversation. Em sẵn sàng communicate bằng English trong work environment, chấp nhận rằng ban đầu sẽ chậm hơn nhưng sẽ improve rapidly with immersion.

TIPS CHUNG

1. **Luôn cho ví dụ cụ thể** — không nói chung chung “tôi biết Docker”, mà “tôi quản lý 70+ containers trên 10 servers”
2. **Nói về trade-offs** — không có giải pháp hoàn hảo, show rằng em hiểu pros/cons
3. **Thành thật về gaps** — cloud experience, team experience — nhưng show learning plan
4. **Quantify** — con số cụ thể: 10 servers, 70+ containers, \$13,800 saved, top 5%, 300K+ contributors
5. **STAR method** cho behavioral questions: Situation → Task → Action → Result
6. **Hỏi ngược** — cuối phỏng vấn hỏi về team structure, tech stack, on-call rotation, growth path