

NYCPS TMS: Prescriptive

Operational Excellence,

Resilience & Readiness

Strategy

I. Introduction: Mandate for Operational Superiority

This document mandates the comprehensive, hyper-detailed strategy for achieving and maintaining **Operational Excellence, Resilience, and Organizational Readiness** for the NYCPS Transportation Management System (TMS). Given the system's critical role in student transportation safety and logistics, its operation within the sensitive AWS GovCloud environment, and the high expectations of a major public sector deployment, operational failures, security incidents, or inability to adapt to disruptions are unacceptable risks. This

strategy provides the prescriptive framework, processes, technical implementations, and organizational alignment necessary to ensure the TMS is not only launched successfully but operates reliably, predictably, securely, and efficiently throughout its lifecycle, capable of weathering diverse disruptions.

Our approach integrates principles from **Site Reliability Engineering (SRE)**, rigorous **Business Continuity Planning (BCP)**, proactive **Disaster Recovery (DR)**, robust **Security Operations (SecOps)**, comprehensive **Organizational Change Management (OCM)**, and **Continuous Improvement** methodologies. It assumes zero tolerance for operational negligence and demands meticulous planning, automation, testing, monitoring, and documented procedures.

Failure to implement and adhere to this operational strategy poses existential risks to the project, potentially leading to service outages impacting student safety, data breaches, compliance violations, reputational damage, and significant financial consequences for NYCPS.

II. Core Pillars of Operational Excellence & Readiness

Mandatory Operational Pillars:

- **Reliability by Design (SRE Focused):** Engineer reliability into the system from the start. Define Service Level Objectives (SLOs) based on user needs, manage via Error Budgets, automate operational tasks (toil), perform blameless post-mortems, and plan capacity proactively.
- **Comprehensive Observability:** Implement deep monitoring across Metrics, Logs, and Traces to understand system behavior, detect anomalies early, and enable rapid troubleshooting (as detailed in the Observability Strategy).
- **Automated & Actionable Alerting:** Configure intelligent alerts focused on user-impacting symptoms and SLO breaches, directly linked to actionable runbooks and routed effectively to on-call personnel.
- **Robust Incident Management:** Implement a structured, practiced incident response process with clear roles, communication protocols, and escalation paths to minimize Mean Time To Detect (MTTD) and Mean Time To Repair (MTTR).
- **Multi-Layered Resilience (BCP/DR):** Implement both technical Disaster Recovery (DR) for IT system continuity and a comprehensive Business Continuity Plan (BCP) for maintaining essential OPT *business operations* during disruptions, covering a wide range of scenarios.
- **Proactive Security Operations (SecOps):****
Continuously monitor for threats, manage vulnerabilities, respond to security incidents rapidly, and maintain a strong security posture aligned with GovCloud and NYCPS standards.

- **Organizational Readiness & Change Management:****
Prepare NYCPSS staff (OPT, Schools, etc.) and SBC personnel for the new system and processes through effective communication, comprehensive training, and ongoing support to ensure adoption and minimize disruption.
- **Continuous Testing & Validation:**** Rigorously test operational procedures, DR/BCP plans, security controls, and system resilience through automated checks, tabletop exercises, DR drills, and chaos engineering (where appropriate).
- **Documentation & Knowledge Management:****
Maintain up-to-date, accessible documentation for operational procedures, runbooks, system configurations, BCP/DR plans, and incident post-mortems.
- **Culture of Continuous Improvement:**** Foster a culture where operational data, incident learnings, and user feedback constantly drive improvements to the system, processes, and documentation.

III. Site Reliability Engineering (SRE)

Implementation

We will establish a dedicated SRE function (or embed SRE responsibilities within DevOps/Ops teams) to focus specifically on the reliability, performance, and scalability of the production TMS environment.

Implementation How-To:

1. Define & Track SLOs/SLIs/Error Budgets:**

- Collaboratively define specific, user-centric SLOs for critical TMS journeys (e.g., Parent App ETA Check < 500ms P95, Ridership Scan Processed & Confirmed < 5s P99, OPT Admin Map Load < 3s P90, Route Generation Job Completion Success > 99.9%).
- Identify corresponding SLIs (metrics like latency, error rate, availability) and implement robust monitoring using CloudWatch/APM tools.

- Calculate Error Budgets based on SLO targets (e.g., 99.9% availability = 43 minutes of downtime allowed per month).
- Create dashboards visualizing SLI performance against SLOs and error budget consumption.
- Implement alerts triggered when error budget burn rate exceeds predefined thresholds, prompting proactive investigation *before* the SLO is breached.
- Use error budget status as a key input for prioritizing reliability work vs. new feature development during sprint planning.

Responsibility: SRE Team, Product Owner, Tech Leads.

2. Automate Operational Toil:**

- SRE team actively identifies manual, repetitive tasks performed by Ops/Support (e.g., common alert responses, manual data corrections, report generation, environment

provisioning steps, access requests).

- **Prioritize automation of high-frequency or high-risk toil using scripting (Python/Bash), IaC (Terraform), AWS Systems Manager Automation documents, Lambda functions, or internal tooling.**
- **Track time spent on toil vs. engineering work, aiming for <50% toil.**

Responsibility: SRE Team, DevOps Team.

3. Capacity Planning & Performance Engineering:**

- **Continuously monitor resource utilization trends and application performance metrics.**
- **Develop capacity models based on historical data and projected growth (student numbers, feature usage).**
- **Conduct regular performance/load tests (as per Test Strategy) to validate capacity limits and identify bottlenecks.**

- **Proactively recommend and implement infrastructure scaling (adjusting ASG/Fargate counts, RDS/ElastiCache instance types) or application performance optimizations based on monitoring and testing results.**

Responsibility: SRE Team,
Performance Engineer, Dev Leads,
Cloud Ops.

4. Production Readiness Reviews (PRR):**

- **SRE team defines and enforces a mandatory PRR checklist for any new service or significant feature change before production deployment.**
- **Checklist covers: Defined SLOs, monitoring/alerting configured, runbooks created, capacity testing completed, dependencies documented, rollback plan validated, security review passed, logging implemented.**
- **SRE provides formal Go/No-Go input during the Release Readiness**

Review based on operational readiness.

Responsibility: SRE Team, Dev Leads, QA Lead, Security Lead.

Quality Gate: Successful completion of PRR checklist required before production deployment approval.

5. Blameless Post-Mortems:** Lead and facilitate blameless post-mortems for all significant production incidents, focusing on systemic improvements and ensuring actionable follow-up items are tracked to completion.

IV. Business Continuity & Disaster Recovery (BCP/DR) Strategy

This section details the integrated strategy for ensuring both IT system resilience (DR) and the continuation of essential OPT business functions (BCP) during various disruptive events. It incorporates the concept of a Minimum Viable Operational State (MVOS) as the primary initial recovery target during severe disruptions.

A. Formal Business Continuity Plan (BCP) - Maintaining Essential Operations

1. What (Scope & Objectives):

- **Objective:** Maintain critical OPT operational functions (prioritizing student safety and essential communication) during significant disruptions affecting IT systems, personnel availability, or facilities, even if full TMS functionality is unavailable.
- **Scope:** Covers scenarios beyond standard IT DR, including: Prolonged AWS GovCloud region outage, major cybersecurity incident (ransomware, destructive attack, data exfiltration), physical facility disruption (OPT offices, key data centers), pandemic/staffing shortages, significant third-party vendor failure (e.g., primary telecom provider, hardware supplier), major geo-political events impacting operations or cloud access, insider threats leading to system compromise, large-scale social media reputation attacks impacting operations.
- **Key Elements (Aligned with RFP 3.27 & Enhanced):**
 - **Business Impact Analysis (BIA):** Identify critical OPT business processes reliant on TMS (e.g., daily

dispatch, exception management, emergency communication, basic ridership safety checks). Determine **Maximum Tolerable Downtime (MTD)** for each process. Assess qualitative and quantitative impacts of disruption (financial, safety, reputational, legal/compliance). This BIA directly informs the MVOS definition.

- ****Risk Assessment:**** Analyze likelihood and impact of various BCP scenarios (as listed above, including cyber threats like ransomware, data exfiltration, insider threats).
- ****Continuity Strategies:**** Define specific strategies for maintaining critical functions, prioritized by BIA results, including:
 - ****Minimum Viable Operational State (MVOS) Definition:**** (See Section IV.A.3 below) - The primary technical/procedural target state for initial recovery.

- ****Manual Workarounds:****

Detailed, documented, *and regularly tested* procedures for performing essential tasks identified in the BIA *without* relying on full TMS functionality (e.g., using pre-printed static route manifests, alternate communication devices like radios/satellite phones, paper-based check-in/exception logging, activating pre-defined call trees).

Specify required offline data/forms.
- ****Alternate Site Operations:**** Plans and logistics for relocating essential OPT staff (identified by role) to

designated alternate work locations (primary/secondary DOE sites, pre-arranged facilities) or enabling secure, effective mandated remote work for critical functions.
Include IT/comms requirements for alternate sites.

- ****Staffing Plans (BCP Focus):** Identify critical personnel/roles needed for both MVOS and manual operations. Define clear backup assignments, implement mandatory cross-training programs on critical manual/MVOS procedures, maintain updated emergency contact lists. Address potential**

pandemic/absenteeism scenarios.

- ****Critical Communication Plan (BCP Focus):** Define communication trees, primary *and* backup methods (e.g., mass notification system, phone trees, satellite phones, runner system if necessary), pre-approved templates for notifying staff, SBCs, schools, parents, emergency services, and NYCPS leadership during diverse BCP events (including scenarios where standard email/VoIP are down).
- ****Third-Party Coordination (BCP**

Focus): Document specific plans for coordinating with critical external parties (SBCs - ensuring they have compatible contingency plans, emergency services, key suppliers, potentially other City agencies) during different disruption scenarios.**

Maintain emergency contact lists for vendors/partners.

- ****Cyber Incident Response Integration:** Explicitly integrate with the Security Incident Response Plan (SIRP) for scenarios involving ransomware, data breaches, or destructive malware, defining**

**business-side actions
(e.g., activating manual
processes while systems
are
investigated/restored).
Include communication
strategies for managing
reputational impact
from cyber events,
coordinating with
NYCPS communications.**

- ****Plan Activation Triggers & Authority:** Define clear, unambiguous criteria for activating different levels of the BCP (ranging from partial activation of manual processes for a specific function to full BCP activation for major disasters). Designate specific roles (e.g., OPT Incident Commander, Senior OPT Leadership) authorized to activate the plan.**
- ****Coordination with NYCPS Overall BCP:** Ensure alignment and clear delineation of responsibilities with the broader NYC Department of Education**

BCP and emergency management frameworks.

2. How (BCP Implementation):

- 1. **Develop BCP Document:**** Create a comprehensive, formal BCP document incorporating all elements above. Store in Confluence, version controlled, and distribute securely to key personnel. Ensure offline copies are available at designated locations/to key roles.
- 2. **Document & Validate Manual Workarounds:**** Work meticulously with OPT SMEs to document *and validate* (through walkthroughs/simulations) step-by-step manual procedures. Ensure required offline data/forms (e.g., emergency route binders, contact lists, paper forms) are generated and distributed regularly (e.g., weekly/monthly updates).
- 3. **Establish & Test Communication Systems:**** Implement and regularly test primary and backup communication methods/call trees (including out-of-band options). Store pre-approved communication templates.
- 4. **Identify Critical Personnel & Cross-Train:**** Maintain and regularly update the list of critical roles/personnel/backups. Schedule and track mandatory cross-training on manual/MVOS procedures.

5. **Tabletop Exercises & Drills (Mandatory &

Frequent):**

- Conduct **quarterly** targeted tabletop exercises focusing on specific scenarios (e.g., key system module unavailable, key vendor failure, specific cyber threat like ransomware). Involve relevant OPT business users, IT/SRE, Security, Comms.**
- Conduct **annual** comprehensive drills simulating large-scale disruptions (e.g., regional power outage, major cyberattack). Test communication plans, manual workarounds, decision-making under pressure, and coordination between BCP and technical DR.**
- Document all exercise outcomes meticulously, identify gaps/weaknesses, and assign actionable improvement items tracked in Jira/Confluence. Update BCP based on lessons learned.**

6. **Integrate with DR Testing: Coordinate BCP exercises (especially those involving IT system**

unavailability) with technical DR tests for holistic validation.

7. ****Review & Update Cycle:**** The BCP ***must*** be formally reviewed, updated, and approved by OPT leadership, CISO, CPO, and Steering Committee at least annually, and after any significant incident or exercise revealing deficiencies.

Responsibility: Business Continuity Manager (Lead), OPT Leadership/SMEs (Process Definition/Validation), PMs, SRE/Ops, Security Team, Comms Lead, HR (Staffing), Training Lead.

The BCP requires significant input and validation from OPT business stakeholders and alignment with overall NYCPSC emergency preparedness. Regular testing is non-negotiable.

3. Minimum Viable Operational State (MVOS) Definition

Purpose:

The MVOS defines the absolute essential subset of TMS functionalities, data, infrastructure, and personnel required to maintain critical student safety oversight and core operational communication during a severe disruption.

Achieving MVOS is the ****primary, immediate objective**** of technical recovery efforts during a SEV1 incident or DR activation, allowing essential functions to resume while full system restoration continues.

How (Implementation & Content):

- 1. Identify MVOS Functions:** Based on the BIA, formally list the minimum required functions documented in the BCP:**
 - ***Example:*** Core User Authentication (Essential OPT Admins, Dispatchers, potentially Drivers via simplified/cached mechanism).
 - ***Example:*** Basic, stable GPS Location Display (potentially with increased latency tolerance, e.g., 1-2 minutes) for active buses, accessible via a simplified view in the OPT Admin Console (or dedicated status page).
 - ***Example:*** Ability to view/access pre-loaded static route assignments for the current day via Driver App (offline mode) and Admin Console.

- *Example:* Core Emergency Broadcast functionality (e.g., OPT Admin direct SMS/Push to active Drivers via a dedicated, resilient channel).
- *Example:* Mechanism to view/input critical, safety-related student exceptions *manually* via a dedicated support channel feeding into OPT operational workflow (outside potentially compromised core systems).
- *Example:* Essential Audit Logging capture for MVOS functions, ensuring basic accountability trail.

2. Define MVOS Data Requirements:** Identify the minimum dataset needed, potentially pre-staged or accessible via read-only replicas:

- *Example:* Pre-calculated static route manifests for the current day (generated daily, accessible offline/via simple storage).

- *Example:* Critical student flags (safety alerts, medical needs) linked to routes/drivers (potentially cached on devices or accessible via simple lookup).
- *Example:* Core driver/vehicle assignment data for the day.
- *Example:* Essential OPT/SBC/Emergency contact information (accessible offline).

3. Define MVOS Infrastructure Subset:** Identify the minimal AWS resource footprint in both primary and DR regions required to run MVOS functions (e.g., core authentication service, minimal API Gateway for emergency comms, specific Lambda functions, RDS read replica for lookups, basic S3, core monitoring/logging). This defines the "Hot/Warm Standby" portion of the DR architecture.

4. Define MVOS Monitoring & Alerting:** Specify the critical subset of CloudWatch alarms focused *solely* on the health and availability of MVOS components (e.g., MVOS API endpoint availability, core auth success rate, emergency broadcast queue depth).

5. Define MVOS Personnel:** List the minimum essential roles/individuals (and their backups) needed to manage/operate the system in MVOS state (likely a subset of the full Incident Response team and core OPT dispatch).

6. Document MVOS Restoration Runbook:** Create specific, streamlined steps within the main DR and Incident Management runbooks detailing *how* to prioritize and restore *only* the MVOS components first, including necessary validation checks.

7. Test MVOS Restoration:** Include specific scenarios in DR drills and Chaos Engineering experiments to validate the ability to establish and operate in the defined MVOS state within the target initial RTO (e.g., <60 mins).

Responsibility: BCM, SRE Lead, Architect, OPT SMEs, PM.

The MVOS definition requires formal sign-off from OPT leadership and technical leads as the agreed-upon initial recovery target during major disruptions.

B. Technical Disaster Recovery (DR) Plan (AWS GovCloud - Enhanced Detail)

1. What (Scope & Objectives - Reiteration):

- ****Objective:**** Restore TMS IT services and data within defined RTOs/RPOs in a secondary AWS GovCloud region following a catastrophic failure impacting the primary region, prioritizing MVOS restoration first.
- ****Scope:**** Covers failure of an entire AWS Region or multiple Availability Zones. Assumes secondary region infrastructure is available.
- ****Mandatory RPO/RTO Targets (Prioritized):****
 - ****MVOS RTO:**** Target initial restoration of defined MVOS functions within 30-60 minutes (requires specific design choices like hot/warm standby for MVOS components).
 - **GPS Data Ingestion/Processing:** RPO=0 (no data loss), RTO=MVOS RTO (requires active-active or hot-standby for core ingestion/location path).
 - **Full Route Planning/Admin Functions:** RPO <= 1 hour, RTO <= 4 hours (post-MVOS restoration).

- **Full Notification System: RPO <= 1 hour, RTO <= 1-2 hours (post-MVOS restoration).**
- **Full Reporting/Analytics: RPO <= 4-24 hours (depending on ETL schedule), RTO <= 8-12 hours (post-MVOS restoration).**

2. How (Implementation - Prescriptive Details):

1. **DR Architecture (Warm Standby + Hot Components for MVOS/RPO=0):**

- **Secondary Region:** Designated secondary AWS GovCloud (US) region.**
- **IaC for DR:** Maintain identical Terraform code structure (`environments/dr/`) mirroring production, parameterized for DR region specifics. Code *must* be tested regularly against DR region.**
- **Core Infrastructure (Warm Standby):** Continuously run foundational infra in DR: VPC, subnets, core security groups, KMS keys (cross-region replicas where possible/needed), ECR repositories**

(with replicated images), potentially a small ECS/EKS control plane, core IAM roles (mirrored).

- ****Critical Data Replication (Hot/Near-Hot for RPO=0/Low RPO):****
 - ****GPS/Ridership Stream:** Implement Kinesis cross-region delivery OR use DynamoDB Global Tables for near real-time replication of critical location/ridership state needed for MVOS/RPO=0.**
Continuously monitor replication lag/status.
 - ****RDS:** Maintain continuously replicating cross-region Read Replicas for critical databases. Monitor**

**replica lag via
CloudWatch alarms.**

- ****DynamoDB (Non-Global):**** Enable PITR and configure regular (e.g., hourly or more frequent) cross-region snapshot copies via AWS Backup.
- ****S3:**** Enable Cross-Region Replication (CRR) for critical buckets (raw data, IaC state, backups, application artifacts). Monitor replication status.
- ****MVOS Compute Components (Warm/Hot Standby):**** Deploy a minimal, scaled-down instance of critical MVOS services (e.g., core AuthN/AuthZ, basic map display API, emergency notification Lambda) continuously running in the DR region

OR implement infrastructure patterns allowing extremely rapid startup (e.g., Lambda functions, pre-pulled images on minimal Fargate capacity provisioned by IaC).

- ****Other Compute (Pilot Light):**** **Keep IaC templates versioned and tested, ready to rapidly deploy the remaining application stack (non-MVOS Lambdas, Fargate services, EC2 routing engines) via the DR CI/CD pipeline during failover activation.**
- ****DNS Failover:**** **Configure Route 53 health checks monitoring primary region MVOS endpoints *and* full application endpoints with aggressive check intervals. Use Route 53 DNS failover routing policies (e.g., Active-Passive Failover with latency or health check routing) to *automatically* redirect traffic to DR region endpoints if primary health checks fail consistently. Set low TTLs (e.g., 60 seconds or less).**

2. **DR Runbook (Hyper-Detailed & Sectioned - Stored in Confluence, Offline Copies Available):**

- ****Section 1: Activation:**** Trigger criteria (Multiple failed Route 53 health checks, manual declaration process/authority defined in BCP), initial communication steps (Automated PagerDuty alert -> SRE on-call acknowledges -> Activate Incident Response Plan), confirmation of primary outage scope/nature.
- ****Section 2: MVOS Restoration (Target: <60 mins):****
 - Verify critical data replication status/lag (RPO check - DynamoDB Global Table status, RDS replica lag via CloudWatch). Log status.
 - Execute script/manual step (if needed) to promote RDS cross-region read replica(s) to standalone instance(s).

Validate promotion success.

- **Verify MVOS compute components (warm standby Fargate/Lambda) are active/scaled up.**
Execute minimal IaC apply (`terraform apply - target=module.mvOS_infra...` via DR pipeline if needed for dependencies.
- **Verify automatic Route 53 DNS failover has occurred for MVOS endpoints. Manually trigger if necessary via pre-approved change.**
- **Run MVOS-specific automated smoke tests (e.g., check login, basic**

map load, emergency broadcast API).

- **Declare MVOS operational via internal/stakeholder comms (Status Page update).**
- ****Section 3: Full System Restoration (Target: <4 hours post-MVOS):****
 - **Execute main IaC apply (`terraform apply` using DR tfvars) via DR CI/CD pipeline (with appropriate approvals if needed) to provision remaining compute (Fargate, Lambda, EC2) and supporting services.**
 - **Restore less critical databases from cross-region snapshots (e.g., reporting DB). Execute necessary data validation post-restore.**

- Run full application integration and E2E smoke tests suites against the DR environment.
- Update DNS records for full application endpoints (if separate from MVOS) or confirm full functionality behind existing DR endpoints.
- Declare full system operational in DR. Communicate status widely.
- **Section 4: Fallback:** Detailed, tested steps to return service to the primary region once verified stable, including: data synchronization back strategy (critical decision point - DMS, application-level sync, potentially read-only period), infrastructure spin-down in DR (via `terraform destroy` or scaling down), DNS changes

(manual or automated), and comprehensive post-failback validation.

3. **DR Testing (Mandatory & Rigorous):**

- **Annual Full Simulation:** Execute the *entire* runbook, including failover, MVOS RTO validation, full system RTO validation, running simulated load, verifying RPO, and executing failback. Document meticulously.**
- **Quarterly Component Tests:** Test specific parts: RDS replica promotion/failback, IaC deployment validation in DR, Kinesis/DynamoDB replication verification, specific manual BCP workaround activation alongside DR component test.**
- **Chaos Engineering (Mature Stage):** Use AWS FIS to simulate failures (AZ outage, RDS instance failure, high network latency between AZs) in *non-production* environments to test HA and failover mechanisms resilience *before* needing full DR.**

- ****Documentation:** All test plans, execution records, results (actual RTO/RPO vs. target), issues found, and remediation actions *must* be documented in Confluence and reviewed by leadership/governance bodies. DR plan/runbooks updated immediately based on findings.**

Responsibility: SRE/Ops Team (Lead/Execution), DevOps Team (IaC/Pipelines), DBA (DB failover/restore/sync), Network Team (DNS), Security Team (DR security posture), Application Teams (Validation/Testing).

Meeting defined RPO/RTOs, especially RPO=0 for critical data, requires careful architectural design (e.g., Global Tables vs. Replicas) and robust, regularly tested replication and failover procedures. This is essential for operational continuity and potentially compliance mandates.

V. Proactive Security Operations (SecOps) Integration

Security is integrated into daily operations, focusing on continuous monitoring, rapid threat response, and maintaining a hardened security posture.

Implementation How-To:

1. Continuous Security Monitoring & SIEM:**

- Deploy and configure AWS native security services: GuardDuty (all detectors enabled), Security Hub (integrated with GuardDuty, Inspector, Config, Macie, partner tools), AWS Config (with compliance packs - FedRAMP, NIST), CloudTrail (enabled, encrypted, integrated).
- Forward high-priority findings from these services AND critical application security logs (auth failures, WAF blocks, input validation failures) via Kinesis Firehose or direct integrations to a central SIEM or security data lake (e.g., AWS OpenSearch with Security Analytics plugin, Splunk)

for correlation, analysis, and alerting.

- **Develop specific detection rules and correlation searches within the SIEM/Security Hub for TMS-specific threats (e.g., anomalous OPT admin activity, attempts to access cross-school data, high rate of driver login failures from unusual locations).**
- **Actively monitor AWS Trusted Advisor security checks and AWS Health Dashboard security notifications.**

Responsibility: Security Team (Setup/Tuning), SRE/Ops (Log Forwarding), DevOps (IaC config).

2. Integrated Vulnerability Management:**

- **Automate SCA, SAST, DAST, and Container Scanning within GitLab CI/CD pipelines (as per DevSecOps strategy).**
- **Run authenticated AWS Inspector scans regularly against EC2 instances (if used).**

- Aggregate findings into a central view (e.g., Security Hub, dedicated Jira project, vulnerability management platform).
- Mandate SLAs for patching/remediation based on severity (e.g., Critical: 7 days, High: 30 days, Medium: 90 days). Track progress via Jira/dashboards.
- Regularly review patching status via Systems Manager Patch Manager compliance reports.

Responsibility: Security Team
(Tracking/Oversight/Tooling),
DevOps/Developers (Patching/Remediation).

3. **Security Incident Response Plan (SIRP)

Execution:**

- Maintain a detailed, actionable SIRP tailored for TMS, covering scenarios like malware, ransomware, DoS/DDoS, unauthorized access, insider threat, data breach/exfiltration.
- Integrate security alerts (GuardDuty High/Medium, Critical

SIEM alerts, WAF significant blocks) into PagerDuty/Opsgenie with dedicated escalation policies for the Security Team (often working alongside SRE/Ops).

- **Follow defined SIRP steps: Triage -> Containment -> Eradication -> Recovery -> Post-Incident Analysis (coordinate with general Incident Management).**
- **Mandate coordination with NYCPS CISO, Legal, Privacy Officer, and potentially NYC3/Law Enforcement for significant incidents, especially data breaches, following established NYCPS protocols.**
- **Conduct quarterly SIRP tabletop exercises focusing on specific threat scenarios.**

Responsibility: Security Team (Lead/Execution), SRE/Ops, Legal, Comms, CISO Liaison.

4. **Proactive Defense & Hardening:**

- **Regularly review and tune AWS WAF rules (Managed + Custom)**

based on application traffic and threat intelligence.

- **Implement rigorous egress filtering (Security Groups, potentially Network Firewall) to limit outbound connections from application instances to only necessary endpoints.**
- **Perform periodic reviews of IAM policies and Security Group rules using automated tools (e.g., IAM Access Analyzer) and manual inspection to ensure least privilege.**
- **Conduct regular security architecture reviews.**

Responsibility: Security Team, Cloud Architect, DevOps/SRE.

VI. Organizational Readiness & Change

Management (OCM) - Operational Focus

Ensuring OPT Staff, School Admins, SBC

Dispatchers/Drivers/Attendants, and other users are prepared, trained, and supported is crucial for successful adoption and operation, especially during disruptions.

A. Formal Organizational Change Management (OCM) Plan Integration

Implementation How-To:

- 1. The dedicated **OCM Plan** *must* explicitly incorporate operational readiness, BCP, DR, and MVOS implications for all user groups.**
- 2. **Impact Assessments:** OCM impact assessments *must* detail how system outages, use of manual workarounds, or activation of DR/MVOS affect specific job roles, workflows, communication needs, and required skills.**
- 3. **Targeted Communications:** Develop communication materials (FAQs, email templates, website updates, meeting briefings) specifically explaining BCP/DR procedures, MVOS operations,**

and emergency communication protocols *in clear, non-technical language* for each stakeholder group (OPT, SBCs, Schools, Parents). Coordinate timing via BCP Comms Plan.

- 4. **Mandatory Training Modules:** Incorporate specific, mandatory modules into role-based training covering:**
 - How to execute relevant manual workaround procedures.**
 - How their role functions during MVOS state (what tools work, what doesn't, alternative communication).**
 - Emergency communication procedures (how to receive alerts, who to contact).**
 - How to report operational issues during disruptions.**

Track completion of this critical training.

- 5. **Readiness Assessments:** Conduct pre-go-live and pre-BCP/DR drill readiness assessments for key operational groups (OPT Dispatch, SBC Dispatch) using checklists to verify understanding of**

procedures, access to offline materials/tools, and communication protocols.

- 6. **Feedback Mechanisms:** Use post-incident surveys and exercise debriefs to gather feedback specifically on the effectiveness of OCM efforts (communication, training, documentation) related to operational resilience.**

Responsibility: OCM Lead, Training Lead, Communications Lead, OPT Leadership/SMEs, PM.

B. Operational Runbook & Knowledge Base Management

Implementation How-To:

- 1. Maintain the central, version-controlled runbook repository (Confluence) covering Incident Response, DR, *and* BCP Manual Workarounds.**
- 2. Runbooks *must* be written clearly, step-by-step, testable, and regularly validated/updated (especially after incidents/exercises or system changes). Assign owners to each runbook.**
- 3. Develop a comprehensive, easily searchable Knowledge Base (KB) within the Support Ticketing**

system (Jira Service Management) or Confluence, specifically including articles for:

- **Common troubleshooting steps for TMS features (accessible by L1/L2/End Users).**
- **Known Error Database (KEDB) with documented workarounds for recurring issues.**
- **Simplified guides on performing key tasks during MVOS or using manual procedures (tailored for different roles).**

4. Establish a process for regularly reviewing and updating KB articles based on support ticket trends and user feedback.

Responsibility: SRE/Ops Team (Technical Runbooks), BCM/OPT SMEs (BCP Workarounds), Support Manager (KB Curation), Technical Writer.

VII. Continuous Improvement Cycle for Operations

Operational excellence requires a relentless focus on learning from experience and data to refine systems and processes.

Implementation How-To:

- 1. **Post-Mortem Actions:**** All action items from incident RCAs (technical, process, communication improvements) *must* be tracked as Jira tickets, prioritized, and implemented. Regularly review status of open post-mortem actions.
- 2. **Metrics & SLO Review:**** Conduct monthly operational reviews analyzing trends in monitoring metrics, SLO adherence, error budget consumption, alert frequency/noise, MTTR/MTTD, and support ticket volumes. Identify areas needing investigation or improvement.
- 3. **BCP/DR Exercise Lessons Learned:**** Findings and recommendations from tabletop exercises and DR drills *must* lead to concrete updates to plans, runbooks, training materials, or technical configurations, tracked via Jira tickets.
- 4. **Runbook Refinement:**** Regularly update runbooks based on operational experience – simplifying steps, adding diagnostic tips, correcting inaccuracies identified during incidents or exercises.

- 5. **Automation Backlog:**** Maintain a backlog of potential automation opportunities (toil reduction) identified by SRE/Ops/Support teams, prioritize based on effort vs. impact.
- 6. **Feedback Integration:**** Systematically review feedback from user channels, support interactions, and readiness assessments to identify recurring operational pain points or areas for system/process enhancement. Feed into product backlog or operational improvement initiatives.
- 7. **Chaos Engineering Program (Mature Stage):**** Use results from Chaos Engineering experiments to proactively identify and fix resilience weaknesses before they cause real incidents.

Responsibility: SRE/Ops Lead (Driving process), All Teams (Contributing feedback/actions), PM (Prioritization interface).

VIII. Resource Implications for Operational Excellence

Implementing and sustaining this level of operational rigor requires dedicated, skilled personnel beyond the core development teams. Under-resourcing these functions is a direct path to operational instability and project failure.

- **SRE Team (~4-8 FTEs):**** Focus on reliability, SLOs, automation, incident management leadership, performance, capacity, DR. Skills: Deep AWS Infra, IaC, Monitoring/Observability Tools, Scripting (Python/Go), Performance Analysis, Incident Command.
- **DevOps Team (~2-4 FTEs):**** Focus on CI/CD pipelines, IaC development/support, developer tooling, environment management automation. Skills: GitLab CI, Terraform/CloudFormation, Docker, Kubernetes/ECS, Scripting.
- **Security Operations (SecOps) Team (~2-3+ FTEs):**** Focus on security monitoring (SIEM), vulnerability management, security incident response, compliance checks, tool tuning. Skills: SIEM tools, AWS Security Services, Incident Response Frameworks, Vulnerability Analysis, Compliance Standards.
- **Database Administrator (DBA - ~1-2 FTEs):**** Focus on production database performance tuning, backup/restore validation, security hardening, patching. Skills: PostgreSQL/PostGIS (deep),

DynamoDB (optional), Performance Tuning, Backup Strategies.

- **Business Continuity Manager (BCM - ~1 FTE - can be dedicated or senior OPT role):**** Owns BCP development, maintenance, training coordination, exercise planning/facilitation. Skills: BCP Standards (ISO 22301), Risk Assessment, Process Analysis, Workshop Facilitation.
- **Organizational Change Management (OCM) Lead (~1-2 FTEs):**** Develops/executes OCM plan, stakeholder analysis, communications, training coordination (process focus). Skills: OCM methodologies (Prosci etc.), Communication, Training Design, Stakeholder Management.
- **Tiered Support Staff (L1/L2):**** Requires dedicated staffing sized based on anticipated user volume and complexity, trained on KB/runbooks/escalation. Skills: Customer Service, Basic Troubleshooting, Ticketing System Usage.
- **Dedicated Training Team:**** Resources to develop and deliver comprehensive training for *all* user groups on system use *and* operational/BCP procedures.

Note: Effective execution requires strong collaboration and clearly defined interfaces between these specialized operational teams and the core development teams.

IX. Conclusion: Building an Anti-Fragile System & Organization

This Operational Excellence, Resilience, and Readiness Strategy provides the mandatory, hyper-detailed blueprint required to build, operate, and maintain the NYCPS TMS as a truly robust, secure, and dependable system. By prescriptively integrating SRE principles, comprehensive BCP/DR (including MVOS), proactive SecOps, rigorous testing (including Chaos Engineering), meticulous OCM, and a culture of continuous learning, we move beyond basic stability towards creating an anti-fragile ecosystem – one that not only withstands disruptions but learns and improves from them.

The exhaustive detail regarding processes, technical implementation (especially within AWS GovCloud), roles, tooling (GitLab, AWS native services), automation, documentation, and governance ensures the required level of control and accountability for this critical public infrastructure. Unwavering adherence to this strategy is paramount for safeguarding student transportation, meeting compliance mandates, managing risks effectively, and delivering sustained operational superiority for the NYCPS.