



在 Amazon EKS 應用程式中設計 HA 和彈性

AWS 方案指引



AWS 方案指引: 在 Amazon EKS 應用程式中設計 HA 和彈性

Copyright © 2026 Amazon Web Services, Inc. and/or its affiliates. All rights reserved.

Amazon 的商標和商業外觀不得用於任何非 Amazon 的產品或服務，也不能以任何可能造成客戶混淆、任何貶低或使 Amazon 名譽受損的方式使用 Amazon 的商標和商業外觀。所有其他非 Amazon 擁有的商標均為其各自擁有者的財產，這些擁有者可能附屬於 Amazon，或與 Amazon 有合作關係，亦或受到 Amazon 贊助。

Table of Contents

| | |
|-------------------------------|----|
| 簡介 | 1 |
| HA 和彈性設計 | 2 |
| 分散工作負載 | 2 |
| 使用 Pod 拓撲分散限制 | 2 |
| Pod 親和性與反親和性 | 6 |
| Pod 中斷預算 | 8 |
| 探查和運作狀態檢查 | 8 |
| 啟動探查 | 9 |
| 活體探查 | 9 |
| 準備度探查 | 9 |
| 傳入資源和負載平衡器運作狀態檢查 | 9 |
| 容器生命週期掛鉤 | 10 |
| 了解區域中斷期間的 Pod 移出 | 12 |
| 實作 Amazon EKS 區域轉移以提高彈性 | 12 |
| 了解區域轉移機制 | 12 |
| 區域轉移啟用方法 | 13 |
| 有效區域轉移的先決條件 | 13 |
| 區域中斷彈性的建議 | 13 |
| 輪班完成和復原 | 14 |
| 結論 | 15 |
| Resources | 16 |
| 文件歷史紀錄 | 17 |
| 詞彙表 | 18 |
| # | 18 |
| A | 18 |
| B | 21 |
| C | 22 |
| D | 25 |
| E | 28 |
| F | 30 |
| G | 31 |
| H | 32 |
| I | 33 |
| L | 35 |

| | |
|-----------|-----|
| M | 36 |
| O | 40 |
| P | 42 |
| Q | 44 |
| R | 44 |
| S | 47 |
| T | 50 |
| U | 51 |
| V | 51 |
| W | 52 |
| Z | 53 |
| liv | liv |

在 Amazon EKS 應用程式中設計高可用性和彈性

Haofei Feng、Frank Fan 和 Rus Kalakutskiy , Amazon Web Services (AWS)

2025 年 10 月 ([文件歷史記錄](#))

確保應用程式設計的高可用性 (HA) 和彈性對於實現近乎零的復原點目標 (RPO) 和復原時間目標 (RTO) 至關重要。隨著組織越來越多地將應用程式遷移和現代化到 Kubernetes 環境，對強大且可擴展解決方案的需求持續增加。Amazon Elastic Kubernetes Service (Amazon EKS) 可協助您大規模有效率地管理容器化應用程式。

本指南深入探討一組公認的建議和最佳實務，用於設計和管理 Amazon EKS 微服務應用程式。根據豐富的經驗和實際部署，這些洞見為架構師和開發人員提供寶貴的指導。針對 Kubernetes 型應用程式的高效能、可靠性和可擴展性實作這些建議，以實現強大的操作。

高可用性和彈性設計考量事項

Kubernetes 的共同責任模型變得更加複雜。Amazon EKS 控制平面可用性和彈性是由 Amazon Web Services () 管理 AWS。您的組織會管理資料平面，這可能會大幅影響微服務應用程式的效能和可用性。

在 Amazon EKS 上設計高可用性和彈性的應用程式時，請考慮下列元件：

- 微服務應用程式：其 Pod 和容器
- 工作負載資料平面：輸入控制器、Pod、[Amazon Virtual Private Cloud \(Amazon VPC\) 容器網路界面 \(CNI\)](#)、服務網格附屬和 kube-proxy 等系統元件
- 工作負載管理層：控制器、許可控制器、網路政策引擎，以及這些元件的持久性資料儲存
- Kubernetes 控制平面
- 基礎設施：節點、網路和網路設備

前三個考量事項是指在 Kubernetes 叢集中執行的元件，本指南涵蓋下列主題：

- [跨節點和可用區域分散工作負載](#)
- [使用 PDB 保護關鍵工作負載](#)
- [設定探查和運作狀態檢查](#)
- [設定容器生命週期掛鉤](#)
- [了解區域中斷期間的 Pod 移出](#)

將工作負載分散到節點和可用區域

將工作負載分散到可用區域和節點等[故障網域](#)，可改善元件可用性，並減少水平擴展應用程式的故障機會。下列各節介紹將工作負載分散到節點和可用區域的方法。

使用 Pod 拓撲分散限制

[Kubernetes Pod 拓撲分散限制](#)，指示 Kubernetes 排程器將由不同故障網域（可用區域、節點和硬體類型）管理 ReplicaSet 的 Pod 分發 StatefulSet。當您使用 Pod 拓撲分散限制時，您可以執行下列動作：

- 根據應用程式需求，將 Pod 分散或集中於不同的故障網域。例如，您可以分配 Pod 以實現彈性，也可以集中 Pod 以實現網路效能。

- 結合不同的條件，例如跨可用區域分佈和跨節點分佈。
- 如果無法滿足條件，請指定偏好的動作：
 - 使用 `whenUnsatisfiable: DoNotSchedule` 搭配 `maxSkew` 和 的組合 `minDomains` 來建立排程器的硬性需求。
 - 使用 `whenUnsatisfiable: ScheduleAnyway` 減少 `maxSkew`。

如果故障區域無法使用，該區域中的 Pod 會變得運作狀態不佳。Kubernetes 會重新排程 Pod，同時盡可能遵守分散限制條件。

下列程式碼顯示使用 Pod 拓撲跨可用區域或跨節點分散限制的範例：

```
...
spec:
  selector:
    matchLabels:
      app: <your-app-label>
  replicas: 3
  template:
    metadata:
      labels: <your-app-label>
    spec:
      serviceAccountName: <ServiceAccountName>
  ...
    topologySpreadConstraints:
      - labelSelector:
          matchLabels:
            app: <your-app-label>
        maxSkew: 1
        topologyKey: topology.kubernetes.io/zone # <---spread those pods evenly over
all availability zones
        whenUnsatisfiable: ScheduleAnyway
      - labelSelector:
          matchLabels:
            app: <your-app-label>
        maxSkew: 1
        topologyKey: kubernetes.io/hostname # <---spread those pods evenly over all
nodes
        whenUnsatisfiable: ScheduleAnyway
```

預設的叢集整體拓撲分散限制

根據預設，Kubernetes 提供一組拓撲分散限制，用於跨節點和可用區域分佈 Pod：

```
defaultConstraints:
- maxSkew: 3
  topologyKey: "kubernetes.io/hostname"
  whenUnsatisfiable: ScheduleAnyway
- maxSkew: 5
  topologyKey: "topology.kubernetes.io/zone"
  whenUnsatisfiable: ScheduleAnyway
```

Note

需要不同類型的拓撲限制的應用程式可以覆寫叢集層級政策。

預設限制條件會設定較高的 maxSkew，這不適用於具有少量 Pod 的部署。截至目前為止，
KubeSchedulerConfiguration無法在 Amazon EKS 中變更。<https://github.com/aws/containers-roadmap/issues/1468>如果您需要強制執行其他一組拓撲分散限制，請考慮使用變動許可控制器，如下節所示。如果您執行替代排程器，也可以控制預設拓撲分散限制條件。不過，管理自訂排程器會增加複雜性，並可能影響叢集彈性和 HA。基於這些原因，我們不建議僅針對拓撲分散限制使用替代排程器。

拓撲分散限制的 Gatekeeper 政策

強制執行拓撲分散限制的另一個選項是使用 Gatekeeper 專案的政策。Gatekeeper 政策是在應用程式層級定義。

下列程式碼範例顯示使用Gatekeeper OPA政策進行部署。您可以根據您的需求修改政策。例如，僅將政策套用到具有標籤的部署HA=true，或使用不同的政策控制器撰寫類似的政策。

第一個範例顯示與 ConstraintTemplate 搭配使用 k8stopologyspreadrequired_template.yaml：

```
apiVersion: templates.gatekeeper.sh/v1
kind: ConstraintTemplate
metadata:
  name: k8stopologyspreadrequired
spec:
  crd:
```

```
spec:
  names:
    kind: K8sTopologySpreadRequired
  validation:
    openAPIV3Schema:
      type: object
      properties:
        message:
          type: string
targets:
  - target: admission.k8s.gatekeeper.sh
    rego: |
      package k8stopologyspreadrequired

      get_message(parameters, _default) =3D msg {
        not parameters.message
        msg :=_default
      }

      get_message(parameters, _default) =3D msg {
        msg := parameters.message
      }

      violation[{"msg": msg}] {
        input.review.kind.kind ="Deployment"
        not input.review.object.spec.template.spec.topologySpreadConstraint
        def_msg :"Pod Topology Spread Constraints are required for Deployments"
        msg :get_message(input.parameters, def_msg)
      }
    }
```

下列程式碼顯示 constraints YAML 資訊清單

k8stopologyspreadrequired_constraint.yml :

```
apiVersion: constraints.gatekeeper.sh/v1beta1
kind: K8sTopologySpreadRequired
metadata:
  name: require-topologyspread-for-deployments
spec:
  match:
    kinds:
      - apiGroups: ["apps"]
```

```
  kinds: ["Deployment"]
  namespaces: ## Without these two lines will apply to the whole cluster
    - "example"
```

何時使用拓撲分散限制條件

考慮在下列案例中使用拓撲分散限制：

- 任何水平擴展的應用程式（例如，無狀態 Web 服務）
- 具有主動-主動或主動-被動複本的應用程式（例如 NoSQL 資料庫或快取）
- 具有待命複本的應用程式（例如，控制器）

可用於水平擴展案例的系統元件，例如，包括下列項目：

- [Cluster Autoscaler](#) 和 [Karpenter](#)（使用 `replicaCount > 1` 和 `leader-elect = true`）
- [AWS Load Balancer](#) 控制器
- [CoreDNS](#)

Pod 親和性與反親和性

在某些情況下，確保節點上執行的特定類型不超過一個 Pod 是有益的。例如，若要避免在同一個節點上排程多個網路密集型 Pod，您可以使用反親和性規則搭配 標籤 `Ingress` 或 `Network-heavy`。使用時 `anti-affinity`，您也可以使用下列的組合：

- 網路最佳化節點上的污點
- 網路密集型 Pod 上的對應公差
- 節點親和性或節點選擇器，以確保網路密集型 Pod 使用網路最佳化執行個體

網路密集型 Pod 做為範例使用。您可能有不同的需求，例如 GPU、記憶體或本機儲存。如需其他使用範例和組態選項，請參閱 [Kubernetes 文件](#)。

重新平衡 Pod

本節討論在 Kubernetes 叢集中重新平衡 Pod 的兩種方法。第一個 使用適用於 Kubernetes 的 `Descheduler`。`Descheduler` 透過強制執行策略來移除違反拓撲分散限制或反親和性規則的 Pod，協助維護 Pod 分佈。第二個方法使用 `Karpenter` 整合和 bin-packing 功能。合併透過將工作負載合併到較少、更有效率的封裝節點上，持續評估和最佳化資源用量。

如果您不使用 Karpenter，建議使用 Descheduler。如果您同時使用 Karpenter 和 Cluster Autoscaler，則可以將 Descheduler 與 Cluster Autoscaler 用於節點群組。

無群組節點的排程器

移除 Pod 時，無法保證拓撲限制仍滿足。例如，縮減部署規模可能會導致 Pod 分佈不平衡。不過，由於 Kubernetes 僅在排程階段使用 Pod 拓撲分散限制，因此整個故障網域的 Pod 會保持不平衡。

若要在這類情況下維持平衡的 Pod 分佈，您可以使用 [Descheduler for Kubernetes](#)。Descheduler 是適用於多種用途的實用工具，例如強制執行最大 Pod 存留期或存留時間 (TTL)，或改善基礎設施的使用。在彈性和高可用性 (HA) 的情況下，請考慮下列排程器策略：

- [RemovePodsViolatingTopologySpreadConstraint](#)
- [RemovePodsViolatingInterPodAntiAffinity](#)
- [RemoveDuplicates](#)

Karpenter 整合和 bin-packing 功能

對於使用 Karpenter 的工作負載，您可以使用整合和 bin-packing 功能來最佳化資源使用率，並降低 Kubernetes 叢集中的成本。Karpenter 會持續評估 Pod 置放和節點使用率，並盡可能嘗試將工作負載合併到較少、更有效率的封裝節點。此程序涉及分析資源需求、考慮 Pod 親和性規則等限制，以及可能在節點之間移動 Pod 以提高整體叢集效率。下列程式碼提供範例：

```
apiVersion: karpenter.sh/v1beta1
kind: NodePool
metadata:
  name: default
spec:
  disruption:
    consolidationPolicy: WhenUnderutilized
    expireAfter: 720h
```

對於 `consolidationPolicy`，您可以使用 `WhenUnderutilized` 或 `WhenEmpty`：

- 當 `consolidationPolicy` 設定為 `WhenUnderutilized`，Karpenter 會將所有節點視為合併。當 Karpenter 探索到空節點或未使用的節點時，Karpenter 會嘗試移除或取代節點以降低成本。
- 當 `consolidationPolicy` 設為 `WhenEmpty`，Karpenter 會考慮僅整合不包含工作負載 Pod 的節點。

Karpenter 整合決策並非僅根據您在監控工具中可能看到的 CPU 或記憶體使用率百分比。反之，Karpenter 會根據 Pod 資源請求和潛在的成本最佳化，使用更複雜的演算法。如需詳細資訊，請參閱 [Karpenter](#) 文件。

使用 PDB 保護關鍵工作負載

Pod 中斷預算 (PDB) 是維護叢集中應用程式高可用性的重要功能。PDB 會指定目標大小，這是特定類型 Pod 的最低可用性。這表示特定 Pod 類型的複本數目下限必須在任何指定時間執行。如果執行中的複本數量低於目標大小，Kubernetes 會防止對剩餘的複本進一步中斷，直到達到目標大小為止。PDBs 有助於確保工作負載不受這些事件影響，並可繼續不間斷地執行。發生中斷時，Kubernetes 會嘗試從受影響的節點正常移出 Pod，同時維持 PDB 中指定的複本數量。

您可以使用 PDB 告知複本的 `minAvailable` 和 `maxUnavailable` 數目。例如，如果您希望應用程式至少有三個副本可供使用，請建立類似下列範例的 PDB：

```
apiVersion: policy/v1beta1
kind: PodDisruptionBudget
metadata:
  name: my-svc-pdb
spec:
  minAvailable: 3
  selector:
    matchLabels:
      app: my-svc
```

為您的應用程式正確設定 PDBs，有助於將計劃或非計劃事件期間的中斷降至最低。您可以使用反親和性規則來排程不同節點上的部署 Pod，並避免節點升級期間的 PDB 延遲。

設定探查和負載平衡器運作狀態檢查

除了負載平衡器運作狀態檢查之外，Kubernetes 還提供數種執行應用程式運作狀態檢查的方式。您可以執行下列 Kubernetes 內建探查與負載平衡器運作狀態檢查，做為 Pod 內容中的命令，或做為 `kubelet` 或主機 IP 地址的 HTTP/TCP 探查。

活體探查和整備探查應該不同且獨立（或至少具有不同的逾時值）。如果應用程式發生暫時問題，整備探查會將 Pod 標記為未就緒，直到問題解決為止。如果活體探查設定不正確，活體探查可能會終止 Pod。

啟動探查

使用啟動探查來保護初始化週期較長的應用程式。在啟動探查成功之前，其他探查都會停用。

您可以定義 Kubernetes 應等待應用程式啟動的時間上限。如果在設定時間上限之後，Pod 仍然失敗啟動探查，應用程式會終止，並建立新的 Pod。

當應用程式的啟動時間無法預測時，請使用啟動探查。如果您知道您的應用程式需要 10 秒才能啟動，請 `initialDelaySeconds` 改用活體探查或整備探查。

活體探查

使用即時性探查來偵測應用程式問題，或程序是否在沒有問題的情況下執行。活體探查可以偵測程序繼續執行但應用程式變得沒有回應的死結條件。使用活體探查時，請執行下列動作：

- 使用 `initialDelaySeconds` 延遲第一個探查。
- 請勿為活體和整備探查設定相同的規格。
- 請勿將活體探查設定為取決於 Pod 外部的因素（例如資料庫）。
- 設定特定的活體探查 `terminationGracePeriodSeconds`。如需詳細資訊，請參閱 [Kubernetes 文件](#)。

準備度探查

使用整備探查來偵測下列項目：

- 應用程式是否已準備好接受流量
- 部分可用性，其中應用程式可能暫時無法使用，但預期在特定操作完成後會再次正常運作

準備度探查有助於確保應用程式組態和相依性在執行時沒有問題或錯誤，以便應用程式可以提供流量。不過，設定不佳的準備度探查可能會導致中斷，而不是防止中斷。取決於資料庫連線等外部因素的整備探查可能會導致所有 Pod 失敗探查。此類故障可能會導致中斷，而且可能會導致從後端服務到使用故障 Pod 的其他 服務發生串聯故障。

傳入資源和負載平衡器運作狀態檢查

Application Load Balancer 和 Kubernetes ingress 提供運作狀態檢查功能。針對 Application Load Balancer 運作狀態檢查，指定目標連接埠和路徑。

Note

對於 Kubernetes ingress，將會有取消註冊延遲。Application Load Balancer 的預設值為 300 秒。請考慮使用您用於整備探查的相同值來設定輸入資源或負載平衡器運作狀態檢查。

NGINX 也提供運作狀態檢查。如需詳細資訊，請參閱 [NGINX 文件](#)。

Istio 輸入和輸出閘道沒有與 NGINX 的 HTTP 運作狀態檢查相當的運作狀態檢查機制。不過，您可以使用 [Istio 斷路器](#) 或 DestinationRule 極端值偵測來實現類似功能。

如需詳細資訊，請參閱《Amazon EKS 最佳實務指南》中的 [可用性和 Pod 生命週期](#)。

設定容器生命週期掛鉤

在正常的容器關閉期間，您的應用程式應該透過啟動關閉來回應SIGTERM訊號，以便用戶端不會經歷任何停機時間。您的應用程式應執行清除程序，如下所示：

- 儲存資料
- 關閉檔案描述項
- 關閉資料庫連線
- 正常完成傳輸中請求
- 及時結束，以滿足 Pod 終止請求

設定足夠長的寬限期，以便清除完成。若要了解如何回應SIGTERM訊號，請參閱您用於應用程式的程式設計語言文件。

[容器生命週期掛鉤](#) 可讓容器了解其管理生命週期中的事件。在執行對應的生命週期關聯時，容器可以執行在處理常式中實作的程式碼。容器生命週期掛鉤為 Kubernetes 和雲端的非同步性質提供了解決方法。這種方法可以防止遺失在輸入資源之前轉送到終止 Pod 的連線，並 iptables 更新為不將新流量轉送到 Pod。

容器生命週期、Endpoint 和 EndpointSlice 是不同 APIs 的一部分。請務必協調這些 APIs。不過，當 Pod 終止時，Kubernetes API 會同時通知 kubelet（容器生命週期）和 EndpointSlice 控制器。如需包括圖表的詳細資訊，請參閱《Amazon EKS 最佳實務指南》中的 [Gracefully 處理用戶端請求](#)。

當 kubelet SIGTERM 傳送至 Pod 時，EndpointSlice 控制器正在終止 EndpointSlice 物件。該終止會通知 Kubernetes API 伺服器，通知每個節點 kube-proxy 的更新 iptables。雖然這些動作

同時發生，但它們之間沒有相依性或序列。容器收到SIGKILL訊號的機率很高，遠比每個節點kube-proxy上的更新本機iptables規則還早。在這種情況下，可能的情況包括下列項目：

- 如果您的應用程式在收到時立即並暗中捨棄處理中的請求和連線SIGTERM，用戶端會看到500錯誤。
- 如果您的應用程式確保所有處理中的請求和連線在收到時完全處理SIGTERM，則在寬限期內，新的用戶端請求仍會傳送到應用程式容器，因為iptables規則可能尚未更新。在清除程序關閉容器上的伺服器通訊端之前，這些新請求將產生新的連線。當寬限期結束時，在SIGTERM傳送之後建立的新連線會無條件捨棄。

若要解決先前的案例，您可以實作應用程式內整合或 PreStop 生命週期掛鉤。如需包括圖表的詳細資訊，請參閱《Amazon EKS 最佳實務指南》中的優雅關機應用程式。

Note

無論應用程式是否正常關閉，或是preStop勾點的結果，應用程式容器最終都會在寬限期結束時透過終止SIGKILL。

使用preStop勾點搭配sleep命令來延遲傳送 SIGTERM。這將有助於在輸入物件將新連線路由至Pod時，繼續接受新連線。測試sleep命令的時間值，以確保將Kubernetes和其他應用程式相依性的任何延遲納入考量，如下列範例所示：

```
apiVersion: apps/v1
kind: Deployment
metadata:
  name: nginx
spec:
  containers:
    - name: nginx
      lifecycle:
        # This "sleep" preStop hook delays the Pod shutdown until
        # after the Ingress Controller removes the matching Endpoint or EndpointSlice
      preStop:
        exec:
          command:
            - /bin/sleep
            - "20"
            # This period should be turned to Ingress/Service Mesh update latency
```

如需詳細資訊，請參閱 Kubernetes 文件中的 [容器掛鉤](#)，以及《Amazon EKS 最佳實務指南》中的 [逐步關閉應用程式](#)。

了解區域中斷期間的 Pod 移出

當完全可用區域中斷時，也就是當該可用區域中的所有節點失去與 Kubernetes 控制平面的連線時，Kubernetes 中的 [節點生命週期控制器](#) 會偵測情況，並從受影響的區域移出 Pod。無法連線節點上的 Pod 會標記為 `Terminating` 並在可用可用區域中的運作狀態良好的節點上排程新的 Pod。在此期間，受影響的節點會顯示 `NotReady` 狀態、排程器防止新的 Pod 放置在這些節點上，而 EndpointSlice 控制器會從服務路由中移除與受損可用區域相關聯的端點，直到恢復連線為止。

對於涉及區域內部分節點故障的案例，其中只有一部分節點變得無法連線，節點生命週期控制器會套用不同的移出行為。如果中斷持續超過設定的容忍期（預設為五分鐘），中斷連線節點上的 Pod 會標記為 `Terminating` 並在可用可用區域中運作狀態良好的節點上排程新的 Pod。

實作 Amazon EKS 區域轉移以提高彈性

[Amazon EKS 區域轉移](#) 與 Amazon Application Recovery Controller (ARC) 整合，提供一種機制，可在可用區域受損期間主動管理流量。此功能可讓網路流量暫時從運作狀態不佳的可用區域重新導向至運作狀態良好的區域內 AWS 區域，以將服務中斷降至最低。

了解區域轉移機制

Amazon EKS 區域轉移會處理東西流量（叢集內的內部 Pod 通訊）。使用 Application Load Balancer 或 Network Load Balancer 設定區域轉移時，也支援傳入流量路由。此機制透過協調多個 Kubernetes 和 AWS 控制平面元件來安全地重新導向流量，而不會中斷執行中的工作負載。在作用中區域轉移期間，Amazon EKS 會自動執行下列協同動作：

- **節點封鎖**：受損可用區域中的所有節點都會進行封鎖。這可防止 Kubernetes 排程器在維護現有工作負載時，在節點上放置新的 Pod。
- **可用區域重新平衡暫停**：對於受管節點群組，可用區域重新平衡操作會暫停，而且 Auto Scaling 群組會更新，僅在運作狀態良好的可用區域中啟動新的資料平面節點。這可確保新的容量不會佈建在受損區域中。
- **端點移除**：EndpointSlice 控制器會從所有相關 EndpointSlices 移除受損可用區域中的 Pod 端點。這可確保服務探索和負載平衡機制只會將流量路由到在運作狀態良好的可用區域中執行的 Pod。
- **工作負載保留**：Amazon EKS 不會終止節點或移出受影響可用區域中的 Pod。它會在受損區域中維持完整容量，因此當區域轉移過期或取消時，流量可以安全地傳回，而不需要額外的擴展操作。

區域轉移啟用方法

根據您的操作模型，您可以選擇兩種方法來啟動區域轉移：

- 透過監控、警報或客戶報告偵測到特定可用區域問題時，[手動區域轉移](#)可提供運算子驅動的控制。此方法需要透過 ARC 主控台、AWS Command Line Interface (AWS CLI) 或區域轉移 APIs 明確動作，其中運算子會指定受損的可用區域並定義轉移的過期時間。當團隊具有專用監控和待命功能，並偏好直接控制流量管理決策時，手動輪班是適當的。
- 當 ARC 根據內部遙測和多個的運作狀態訊號，包括網路指標、Amazon Elastic Compute Cloud (Amazon EC2) 和 Elastic Load Balancing。當指標顯示問題已解決時 AWS 服務，AWS 會自動結束自動轉移，區域[自動轉移](#)會授權 AWS 自動啟動轉移。如果您想要以最少的手動介入來達到最高的可用性狀態，建議您使用此方法，因為它可對偵測到的可用區域受損啟用次分鐘的回應。

有效區域轉移的先決條件

若要在可用區域受損期間成功保護應用程式的區域轉移，您必須先架構叢集以實現異地同步備份彈性，才能啟用區域轉移功能：

- 多可用區域節點分佈：在至少三個可用區域佈建工作者節點，以確保一個區域無法使用時有足夠的備援。
- 容量規劃：跨運作狀態良好的可用區域預先佈建足夠的運算容量，以便在從服務中移除一個可用區域時容納完整工作負載，因為在主動中斷期間擴展操作可能會遇到容量不足的情況。
- Pod 分佈和預先擴展：在所有可用區域部署每個應用程式的多個複本，並預先擴展每個區域中的關鍵系統元件，例如 [CoreDNS](#)。這有助於確保在轉移區域後仍有足夠的容量。

區域中斷彈性的建議

- 在叢集建立時啟用區域轉移：對於新的 EKS 叢集，請在透過 Amazon EKS 主控台的初始佈建期間啟用與 ARC 的區域轉移整合 AWS CLI，或基礎設施即程式碼 (IaC) 工具，例如 AWS CloudFormation。使用快速組態建立 [的 EKS Auto Mode 叢集](#)預設會啟用區域轉移。
- 選取適當的啟用方法：針對需要自動化回應的最大可用性的生產環境選擇區域自動轉移，尤其是在可用區域受損期間停機時間可能帶來重大業務影響的面向客戶應用程式。對於營運團隊偏好在流量轉移之前提供明確核准的環境，或仍在進行應用程式測試和驗證的環境，請使用手動區域轉移。
- 在生產部署之前測試彈性：透過手動啟動測試區域轉移或啟用區域自動轉移實務執行來驗證單一可用區損失下的叢集行為，以確認應用程式在減少可用區域計數的情況下操作時維持可用性、效能仍然可接受且容量足夠。我們強烈建議您進行此測試，以便在實際可用區域受損之前識別組態差距。

- **與負載平衡器組態協調**：對於接收外部流量的應用程式，在相關聯的 Application Load Balancer 和 Network Load Balancer 上啟用 ARC 區域轉移，以確保在可用區域受損期間傳入流量和叢集內東西流量一起轉移。此協調可防止外部請求達到運作狀態良好的 Pod，但這些 Pod 無法與轉移區域中的相依性通訊的情況。
- **監控輪班操作**：啟用區域輪班後，請設定輪班事件的監控和提醒，包括自動輪班啟用、手動輪班啟動和輪班過期，以保持對流量管理動作及其對應用程式行為的影響的操作可見性。

輪班完成和復原

當區域轉移根據其設定的持續時間過期，或在可用區域受損解決後手動取消時，EndpointSlice 控制器會自動更新所有 EndpointSlices，以重新整合還原可用區域中的端點。當用戶端重新整理端點資訊和建立新的連線時，流量會逐漸返回先前受影響的區域。這可啟用完整的叢集容量使用率，而不需要手動介入或 Pod 重新排程。

結論

當您設計具有高可用性和彈性的架構時，請考慮下列元件：

- 微服務應用程式（其 Pod 和容器）
- 工作負載資料平面（輸入控制器、Pod、[Amazon VPC CNI](#)、服務網格附屬和 kube-proxy 等系統元件）
- 工作負載管理層（控制器、許可控制器、網路政策引擎，以及這些元件的持久性資料儲存）
- Kubernetes 控制平面
- 基礎設施（節點、網路和網路設備）

若要解決這些元件考量，請使用下列關鍵策略：

- 為了協助確保高可用性和容錯能力，請將工作負載分散到節點和可用區域。
- 為了保護關鍵工作負載，請使用 Pod 中斷預算 (PDBs) 在中斷期間維持應用程式穩定性。
- 為了協助確保 Pod 正確執行和提供流量，請設定啟動探查、活體探查、整備探查和負載平衡器運作狀態檢查。
- 若要有效率地管理容器狀態轉換，請設定容器生命週期關聯。
- 若要在節點故障或維護期間控制移出程序，請設定 Pod 移出時間。

透過實作這些實務，您可以大幅提升在 Amazon EKS 上執行之應用程式的可靠性和彈性，確保強大的效能和高可用性。

Resources

- [Kubernetes Pod 拓撲分散限制條件](#) (Kubernetes 文件)
- [Karpenter FAQs](#) (Karpenter 文件)
- [Descheduler for Kubernetes](#) (GitHub 儲存庫)
- [可用性和 Pod 生命週期](#) Amazon EKS 最佳實務指南
- [正常關閉應用程式](#) Amazon EKS 最佳實務指南
- [【EKS】 【request】：能夠設定 pod-eviction-timeout 和解決方法](#) (Containers Roadmap 儲存庫)

文件歷史紀錄

下表描述了本指南的重大變更。如果您想收到有關未來更新的通知，可以訂閱 [RSS 摘要](#)。

| 變更 | 描述 | 日期 |
|----------------------|---------------------------------|------------------|
| 更新 | 已修訂有關 <u>區域中斷期間 Pod 移出的</u> 章節。 | 2025 年 10 月 29 日 |
| 更新 | 修訂 <u>使用 Pod 拓撲分散限制區段</u> 。 | 2025 年 1 月 27 日 |
| 初次出版 | — | 2024 年 10 月 23 日 |

AWS 規範性指引詞彙表

以下是 AWS Prescriptive Guidance 提供的策略、指南和模式中常用的術語。若要建議項目，請使用詞彙表末尾的提供意見回饋連結。

數字

7 R

將應用程式移至雲端的七種常見遷移策略。這些策略以 Gartner 在 2011 年確定的 5 R 為基礎，包括以下內容：

- 重構/重新架構 – 充分利用雲端原生功能來移動應用程式並修改其架構，以提高敏捷性、效能和可擴展性。這通常涉及移植作業系統和資料庫。範例：將您的現場部署 Oracle 資料庫遷移至 Amazon Aurora PostgreSQL 相容版本。
- 平台轉換 (隨即重塑) – 將應用程式移至雲端，並引入一定程度的優化以利用雲端功能。範例：將內部部署 Oracle 資料庫遷移至 中的 Amazon Relational Database Service (Amazon RDS) for Oracle AWS 雲端。
- 重新購買 (捨棄再購買) – 切換至不同的產品，通常從傳統授權移至 SaaS 模型。範例：將您的客戶關係管理 (CRM) 系統遷移至 Salesforce.com。
- 主機轉換 (隨即轉移) – 將應用程式移至雲端，而不進行任何變更以利用雲端功能。範例：將您的現場部署 Oracle 資料庫遷移至 中 EC2 執行個體上的 Oracle AWS 雲端。
- 重新放置 (虛擬機器監視器等級隨即轉移) – 將基礎設施移至雲端，無需購買新硬體、重寫應用程式或修改現有操作。您可以將伺服器從內部部署平台遷移到相同平台的雲端服務。範例：將 Microsoft Hyper-V 應用程式遷移至 AWS。
- 保留 (重新檢視) – 將應用程式保留在來源環境中。其中可能包括需要重要重構的應用程式，且您希望將該工作延遲到以後，以及您想要保留的舊版應用程式，因為沒有業務理由來進行遷移。
- 淘汰 – 解除委任或移除來源環境中不再需要的應用程式。

A

ABAC

請參閱[屬性型存取控制](#)。

抽象服務

請參閱 [受管服務](#)。

ACID

請參閱 [原子性、一致性、隔離性、耐久性](#)。

主動-主動式遷移

一種資料庫遷移方法，其中來源和目標資料庫保持同步（透過使用雙向複寫工具或雙重寫入操作），且兩個資料庫都在遷移期間處理來自連接應用程式的交易。此方法支援小型、受控制批次的遷移，而不需要一次性切換。它更靈活，但比 [主動-被動遷移](#) 需要更多的工作。

主動-被動式遷移

一種資料庫遷移方法，其中來源和目標資料庫保持同步，但只有來源資料庫會在資料複寫至目標資料庫時處理來自連線應用程式的交易。目標資料庫在遷移期間不接受任何交易。

彙總函數

在一組資料列上運作的 SQL 函數，會計算群組的單一傳回值。彙總函數的範例包括 SUM 和 MAX。

AI

請參閱 [人工智慧](#)。

AIOps

請參閱 [人工智慧操作](#)。

匿名化

永久刪除資料集中個人資訊的程序。匿名化有助於保護個人隱私權。匿名資料不再被視為個人資料。

反模式

經常用於重複性問題的解決方案，其中解決方案具有反生產力、無效或比替代解決方案更有效。

應用程式控制

一種安全方法，僅允許使用核准的應用程式，以協助保護系統免受惡意軟體攻擊。

應用程式組合

有關組織使用的每個應用程式的詳細資訊的集合，包括建置和維護應用程式的成本及其商業價值。此資訊是 [產品組合探索和分析程序](#) 的關鍵，有助於識別要遷移、現代化和優化的應用程式並排定其優先順序。

人工智慧 (AI)

電腦科學領域，致力於使用運算技術來執行通常與人類相關的認知功能，例如學習、解決問題和識別模式。如需詳細資訊，請參閱[什麼是人工智慧？](#)

人工智慧操作 (AIOps)

使用機器學習技術解決操作問題、減少操作事件和人工干預以及提高服務品質的程序。如需有關如何在 AWS 遷移策略中使用 AIOps 的詳細資訊，請參閱[操作整合指南](#)。

非對稱加密

一種加密演算法，它使用一對金鑰：一個用於加密的公有金鑰和一個用於解密的私有金鑰。您可以共用公有金鑰，因為它不用於解密，但對私有金鑰存取應受到高度限制。

原子性、一致性、隔離性、耐久性 (ACID)

一組軟體屬性，即使在出現錯誤、電源故障或其他問題的情況下，也能確保資料庫的資料有效性和操作可靠性。

屬性型存取控制 (ABAC)

根據使用者屬性 (例如部門、工作職責和團隊名稱) 建立精細許可的實務。如需詳細資訊，請參閱《AWS Identity and Access Management (IAM) 文件》中的[ABAC for AWS](#)。

授權資料來源

存放主要版本資料的位置，被視為最可靠的資訊來源。您可以將授權資料來源中的資料複製到其他位置，以處理或修改資料，例如匿名、修訂或假名化資料。

可用區域

中的不同位置 AWS 區域，可隔離其他可用區域中的故障，並提供相同區域中其他可用區域的低成本、低延遲網路連線能力。

AWS 雲端採用架構 (AWS CAF)

的指導方針和最佳實務架構 AWS，可協助組織制定高效且有效的計劃，以成功地移至雲端。AWS CAF 將指導方針組織到六個重點領域：業務、人員、治理、平台、安全和營運。業務、人員和控管層面著重於業務技能和程序；平台、安全和操作層面著重於技術技能和程序。例如，人員層面針對處理人力資源 (HR)、人員配備功能和人員管理的利害關係人。因此，AWS CAF 為人員開發、訓練和通訊提供指引，協助組織做好成功採用雲端的準備。如需詳細資訊，請參閱[AWS CAF 網站](#)和[AWS CAF 白皮書](#)。

AWS 工作負載資格架構 (AWS WQF)

評估資料庫遷移工作負載、建議遷移策略並提供工作預估值的工具。 AWS WQF 隨附於 AWS Schema Conversion Tool (AWS SCT)。它會分析資料庫結構描述和程式碼物件、應用程式程式碼、相依性和效能特性，並提供評估報告。

B

錯誤的機器人

旨在中斷或傷害個人或組織的機器人。

BCP

請參閱業務持續性規劃。

行為圖

資源行為的統一互動式檢視，以及一段時間後的互動。您可以將行為圖與 Amazon Detective 搭配使用來檢查失敗的登入嘗試、可疑的 API 呼叫和類似動作。如需詳細資訊，請參閱偵測文件中的行為圖中的資料。

大端序系統

首先儲存最高有效位元組的系統。另請參閱 [Endianness](#)。

二進制分類

預測二進制結果的過程 (兩個可能的類別之一)。例如，ML 模型可能需要預測諸如「此電子郵件是否是垃圾郵件？」等問題 或「產品是書還是汽車？」

Bloom 篩選條件

一種機率性、記憶體高效的資料結構，用於測試元素是否為集的成員。

藍/綠部署

一種部署策略，您可以在其中建立兩個不同但相同的環境。您可以在一個環境（藍色）中執行目前的應用程式版本，並在另一個環境（綠色）中執行新的應用程式版本。此策略可協助您快速復原，並將影響降至最低。

機器人

透過網際網路執行自動化任務並模擬人類活動或互動的軟體應用程式。有些機器人有用或有益，例如在網際網路上編製資訊索引的 Web 爬蟲程式。某些其他機器人稱為惡意機器人，旨在中斷或傷害個人或組織。

殭屍網路

受到惡意軟體感染且受單一方控制之機器人的網路，稱為機器人繼承器或機器人運算子。殭屍網路是擴展機器人及其影響的最佳已知機制。

分支

程式碼儲存庫包含的區域。儲存庫中建立的第一個分支是主要分支。您可以從現有分支建立新分支，然後在新分支中開發功能或修正錯誤。您建立用來建立功能的分支通常稱為功能分支。當準備好發佈功能時，可以將功能分支合併回主要分支。如需詳細資訊，請參閱[關於分支](#) (GitHub 文件)。

碎片存取

在特殊情況下，以及透過核准的程序，讓使用者快速取得他們通常無權存取 AWS 帳戶之的存取權。如需詳細資訊，請參閱 Well-Architected 指南中的 AWS [實作打破玻璃程序](#) 指標。

棕地策略

環境中的現有基礎設施。對系統架構採用棕地策略時，可以根據目前系統和基礎設施的限制來設計架構。如果正在擴展現有基礎設施，則可能會混合棕地和綠地策略。

緩衝快取

儲存最常存取資料的記憶體區域。

業務能力

業務如何創造價值 (例如，銷售、客戶服務或營銷)。業務能力可驅動微服務架構和開發決策。如需詳細資訊，請參閱[在 AWS 上執行容器化微服務](#)白皮書的圍繞業務能力進行組織部分。

業務連續性規劃 (BCP)

一種解決破壞性事件 (如大規模遷移) 對營運的潛在影響並使業務能夠快速恢復營運的計畫。

C

CAF

請參閱[AWS 雲端採用架構](#)。

Canary 部署

版本對最終使用者的緩慢和增量版本。當您有信心時，您可以部署新版本並完全取代目前的版本。

CCoE

請參閱[Cloud Center of Excellence](#)。

CDC

請參閱[變更資料擷取](#)。

變更資料擷取 (CDC)

追蹤對資料來源 (例如資料庫表格) 的變更並記錄有關變更的中繼資料的程序。您可以將 CDC 用於各種用途，例如稽核或複寫目標系統中的變更以保持同步。

混沌工程

故意引入故障或破壞性事件，以測試系統的彈性。您可以使用[AWS Fault Injection Service \(AWS FIS\)](#) 執行實驗，為您的 AWS 工作負載帶來壓力，並評估其回應。

CI/CD

請參閱[持續整合和持續交付](#)。

分類

有助於產生預測的分類程序。用於分類問題的 ML 模型可預測離散值。離散值永遠彼此不同。例如，模型可能需要評估影像中是否有汽車。

用戶端加密

在目標 AWS 服務 接收資料之前，在本機加密資料。

雲端卓越中心 (CCoE)

一個多學科團隊，可推動整個組織的雲端採用工作，包括開發雲端最佳實務、調動資源、制定遷移時間表以及領導組織進行大規模轉型。如需詳細資訊，請參閱 AWS 雲端企業策略部落格上的[CCoE 文章](#)。

雲端運算

通常用於遠端資料儲存和 IoT 裝置管理的雲端技術。雲端運算通常連接到[邊緣運算](#)技術。

雲端操作模型

在 IT 組織中，用於建置、成熟和最佳化一或多個雲端環境的操作模型。如需詳細資訊，請參閱[建置您的雲端操作模型](#)。

採用雲端階段

組織在遷移至 時通常會經歷的四個階段 AWS 雲端：

- 專案 – 執行一些與雲端相關的專案以進行概念驗證和學習用途
- 基礎 – 進行基礎投資以擴展雲端採用 (例如，建立登陸區域、定義 CCoE、建立營運模型)

- 遷移 – 遷移個別應用程式
- 重塑 – 優化產品和服務，並在雲端中創新

這些階段由 Stephen Orban 在部落格文章 [The Journey Toward Cloud-First 和企業策略部落格上的採用階段](#) 中定義。 AWS 雲端 如需有關它們如何與 AWS 遷移策略相關的詳細資訊，請參閱[遷移整備指南](#)。

CMDB

請參閱[組態管理資料庫](#)。

程式碼儲存庫

透過版本控制程序來儲存及更新原始程式碼和其他資產 (例如文件、範例和指令碼) 的位置。常見的雲端儲存庫包括 GitHub 或 Bitbucket Cloud。程式碼的每個版本都稱為分支。在微服務結構中，每個儲存庫都專用於單個功能。單一 CI/CD 管道可以使用多個儲存庫。

冷快取

一種緩衝快取，它是空的、未填充的，或者包含過時或不相關的資料。這會影響效能，因為資料庫執行個體必須從主記憶體或磁碟讀取，這比從緩衝快取讀取更慢。

冷資料

很少存取且通常是歷史資料的資料。查詢這類資料時，通常可接受慢查詢。將此資料移至效能較低且成本較低的儲存層或類別，可以降低成本。

電腦視覺 (CV)

AI 欄位[???](#)，使用機器學習從數位影像和影片等視覺化格式分析和擷取資訊。例如，Amazon SageMaker AI 提供 CV 的影像處理演算法。

組態偏離

對於工作負載，組態會從預期狀態變更。這可能會導致工作負載不合規，而且通常是漸進和無意的。

組態管理資料庫 (CMDB)

儲存和管理有關資料庫及其 IT 環境的資訊的儲存庫，同時包括硬體和軟體元件及其組態。您通常在遷移的產品組合探索和分析階段使用 CMDB 中的資料。

一致性套件

您可以組合的 AWS Config 規則和修補動作集合，以自訂您的合規和安全檢查。您可以使用 YAML 範本，將一致性套件部署為 AWS 帳戶 和 區域中或整個組織的單一實體。如需詳細資訊，請參閱 AWS Config 文件中的一致性套件。

持續整合和持續交付 (CI/CD)

自動化軟體發行程序的來源、建置、測試、暫存和生產階段的程序。CI/CD 通常被描述為管道。CI/CD 可協助您將程序自動化、提升生產力、改善程式碼品質以及加快交付速度。如需詳細資訊，請參閱[持續交付的優點](#)。CD 也可表示持續部署。如需詳細資訊，請參閱[持續交付與持續部署](#)。

CV

請參閱[電腦視覺](#)。

D

靜態資料

網路中靜止的資料，例如儲存中的資料。

資料分類

根據重要性和敏感性來識別和分類網路資料的程序。它是所有網路安全風險管理策略的關鍵組成部分，因為它可以協助您確定適當的資料保護和保留控制。資料分類是 AWS Well-Architected Framework 中安全支柱的元件。如需詳細資訊，請參閱[資料分類](#)。

資料偏離

生產資料與用於訓練 ML 模型的資料之間有意義的變化，或輸入資料隨時間有意義的變更。資料偏離可以降低 ML 模型預測的整體品質、準確性和公平性。

傳輸中的資料

在您的網路中主動移動的資料，例如在網路資源之間移動。

資料網格

架構架構，提供分散式、分散式資料擁有權與集中式管理。

資料最小化

僅收集和處理嚴格必要資料的原則。在中實作資料最小化 AWS 雲端可以降低隱私權風險、成本和分析碳足跡。

資料周邊

AWS 環境中的一組預防性防護機制，可協助確保只有信任的身分才能從預期的網路存取信任的資源。如需詳細資訊，請參閱[在上建置資料周邊 AWS](#)。

資料預先處理

將原始資料轉換成 ML 模型可輕鬆剖析的格式。預處理資料可能意味著移除某些欄或列，並解決遺失、不一致或重複的值。

資料來源

在整個資料生命週期中追蹤資料的來源和歷史記錄的程序，例如資料的產生、傳輸和儲存方式。

資料主體

正在收集和處理資料的個人。

資料倉儲

支援商業智慧的資料管理系統，例如 分析。資料倉儲通常包含大量歷史資料，通常用於查詢和分析。

資料庫定義語言 (DDL)

用於建立或修改資料庫中資料表和物件之結構的陳述式或命令。

資料庫處理語言 (DML)

用於修改 (插入、更新和刪除) 資料庫中資訊的陳述式或命令。

DDL

請參閱[資料庫定義語言](#)。

深度整體

結合多個深度學習模型進行預測。可以使用深度整體來獲得更準確的預測或估計預測中的不確定性。

深度學習

一個機器學習子領域，它使用多層人工神經網路來識別感興趣的輸入資料與目標變數之間的對應關係。

深度防禦

這是一種資訊安全方法，其中一系列的安全機制和控制項會在整個電腦網路中精心分層，以保護網路和其中資料的機密性、完整性和可用性。當您上採用此策略時 AWS，您可以在 AWS Organizations 結構的不同層新增多個控制項，以協助保護資源。例如，defense-in-depth 方法可能會結合多重要素驗證、網路分割和加密。

委派的管理員

在 AWS Organizations 中，相容的服務可以註冊 AWS 成員帳戶來管理組織的帳戶，並管理該服務的許可。此帳戶稱為該服務的委派管理員。如需詳細資訊和相容服務清單，請參閱 AWS Organizations 文件中的[可搭配 AWS Organizations 運作的服務](#)。

deployment

在目標環境中提供應用程式、新功能或程式碼修正的程序。部署涉及在程式碼庫中實作變更，然後在應用程式環境中建置和執行該程式碼庫。

開發環境

請參閱 [環境](#)。

偵測性控制

一種安全控制，用於在事件發生後偵測、記錄和提醒。這些控制是第二道防線，提醒您注意繞過現有預防性控制的安全事件。如需詳細資訊，請參閱在 AWS 上實作安全控制中的[偵測性控制](#)。

開發值串流映射 (DVSM)

一種程序，用於識別並優先考慮對軟體開發生命週期中的速度和品質造成負面影響的限制。DVSM 擴展了最初專為精簡製造實務設計的價值串流映射程序。它著重於透過軟體開發程序建立和移動價值所需的步驟和團隊。

數位分身

真實世界系統的虛擬呈現，例如建築物、工廠、工業設備或生產線。數位分身支援預測性維護、遠端監控和生產最佳化。

維度資料表

在[星星結構描述](#)中，較小的資料表包含有關事實資料表中量化資料的資料屬性。維度資料表屬性通常是文字欄位或離散數字，其行為類似於文字。這些屬性通常用於查詢限制、篩選和結果集標記。

災難

防止工作負載或系統在其主要部署位置實現其業務目標的事件。這些事件可能是自然災難、技術故障或人為動作的結果，例如意外設定錯誤或惡意軟體攻擊。

災難復原 (DR)

您用來將[災難](#)造成的停機時間和資料遺失降至最低的策略和程序。如需詳細資訊，請參閱 AWS Well-Architected Framework 中的[上工作負載的災難復原 AWS：雲端中的復原](#)。

DML

請參閱資料庫處理語言。

領域驅動的設計

一種開發複雜軟體系統的方法，它會將其元件與每個元件所服務的不斷發展的領域或核心業務目標相關聯。Eric Evans 在其著作 Domain-Driven Design: Tackling Complexity in the Heart of Software (Boston: Addison-Wesley Professional, 2003) 中介紹了這一概念。如需有關如何將領域驅動的設計與 strangler fig 模式搭配使用的資訊，請參閱使用容器和 Amazon API Gateway 逐步現代化舊版 Microsoft ASP.NET (ASMX) Web 服務。

DR

請參閱災難復原。

偏離偵測

追蹤與基準組態的偏差。例如，您可以使用 AWS CloudFormation 來偵測系統資源中的偏離，也可以使用 AWS Control Tower 來偵測登陸區域中可能影響控管要求合規性的變更。<https://docs.aws.amazon.com/AWSCloudFormation/latest/UserGuide/using-cfn-stack-drift.html>

DVSM

請參閱開發值串流映射。

E

EDA

請參閱探索性資料分析。

EDI

請參閱電子資料交換。

邊緣運算

提升 IoT 網路邊緣智慧型裝置運算能力的技術。與雲端運算相比，邊緣運算可以減少通訊延遲並改善回應時間。

電子資料交換 (EDI)

在組織之間自動交換商業文件。如需詳細資訊，請參閱什麼是電子資料交換。

加密

將人類可讀取的純文字資料轉換為加密文字的運算程序。

加密金鑰

由加密演算法產生的隨機位元的加密字串。金鑰長度可能有所不同，每個金鑰的設計都是不可預測且唯一的。

端序

位元組在電腦記憶體中的儲存順序。大端序系統首先儲存最高有效位元組。小端序系統首先儲存最低有效位元組。

端點

請參閱 [服務端點](#)。

端點服務

您可以在虛擬私有雲端 (VPC) 中託管以與其他使用者共用的服務。您可以使用 [建立端點服務](#)，AWS PrivateLink 並將許可授予其他 AWS 帳戶或 AWS Identity and Access Management (IAM) 委託人。這些帳戶或主體可以透過建立介面 VPC 端點私下連接至您的端點服務。如需詳細資訊，請參閱 Amazon Virtual Private Cloud (Amazon VPC) 文件中的[建立端點服務](#)。

企業資源規劃 (ERP)

一種系統，可自動化和管理企業的關鍵業務流程（例如會計、[MES](#) 和專案管理）。

信封加密

使用另一個加密金鑰對某個加密金鑰進行加密的程序。如需詳細資訊，請參閱 [\(\) 文件中的信封加密](#)。 AWS Key Management Service AWS KMS

環境

執行中應用程式的執行個體。以下是雲端運算中常見的環境類型：

- 開發環境 – 執行中應用程式的執行個體，只有負責維護應用程式的核心團隊才能使用。開發環境用來測試變更，然後再將開發環境提升到較高的環境。此類型的環境有時稱為測試環境。
- 較低的環境 – 應用程式的所有開發環境，例如用於初始建置和測試的開發環境。
- 生產環境 – 最終使用者可以存取的執行中應用程式的執行個體。在 CI/CD 管道中，生產環境是最後一個部署環境。
- 較高的環境 – 核心開發團隊以外的使用者可存取的所有環境。這可能包括生產環境、生產前環境以及用於使用者接受度測試的環境。

epic

在敏捷方法中，有助於組織工作並排定工作優先順序的功能類別。epic 提供要求和實作任務的高層級描述。例如，AWS CAF 安全概念包括身分和存取管理、偵測控制、基礎設施安全、資料保護和事件回應。如需有關 AWS 遷移策略中的 Epic 的詳細資訊，請參閱[計畫實作指南](#)。

ERP

請參閱[企業資源規劃](#)。

探索性資料分析 (EDA)

分析資料集以了解其主要特性的過程。您收集或彙總資料，然後執行初步調查以尋找模式、偵測異常並檢查假設。透過計算摘要統計並建立資料可視化來執行 EDA。

F

事實資料表

[星狀結構描述](#)中的中央資料表。它存放有關業務操作的量化資料。一般而言，事實資料表包含兩種類型的資料欄：包含度量的資料，以及包含維度資料表外部索引鍵的資料欄。

快速失敗

一種使用頻繁且增量測試來縮短開發生命週期的理念。這是敏捷方法的關鍵部分。

故障隔離界限

在 AWS 雲端，像是可用區域 AWS 區域、控制平面或資料平面等界限會限制故障的影響，並有助於改善工作負載的彈性。如需詳細資訊，請參閱[AWS 故障隔離界限](#)。

功能分支

請參閱[分支](#)。

特徵

用來進行預測的輸入資料。例如，在製造環境中，特徵可能是定期從製造生產線擷取的影像。

功能重要性

特徵對於模型的預測有多重要。這通常表示為可以透過各種技術來計算的數值得分，例如 Shapley Additive Explanations (SHAP) 和積分梯度。如需詳細資訊，請參閱[的機器學習模型可解譯性 AWS](#)。

特徵轉換

優化 ML 程序的資料，包括使用其他來源豐富資料、調整值、或從單一資料欄位擷取多組資訊。這可讓 ML 模型從資料中受益。例如，如果將「2021-05-27 00:15:37」日期劃分為「2021」、「五月」、「週四」和「15」，則可以協助學習演算法學習與不同資料元件相關聯的細微模式。

少量擷取提示

在要求 [LLM](#) 執行類似的任務之前，提供少量示範任務和所需輸出的範例。此技術是內容內學習的應用程式，其中模型會從內嵌在提示中的範例（快照）中學習。少量的提示對於需要特定格式、推理或網域知識的任務來說非常有效。另請參閱[零鏡頭提示](#)。

FGAC

請參閱[精細存取控制](#)。

精細存取控制 (FGAC)

使用多個條件來允許或拒絕存取請求。

閃切遷移

一種資料庫遷移方法，透過[變更資料擷取](#)使用連續資料複寫，以盡可能在最短的時間內遷移資料，而不是使用分階段方法。目標是將停機時間降至最低。

FM

請參閱[基礎模型](#)。

基礎模型 (FM)

大型深度學習神經網路，已在廣義和未標記資料的大量資料集上進行訓練。FMs 能夠執行各種一般任務，例如了解語言、產生文字和影像，以及以自然語言交談。如需詳細資訊，請參閱[什麼是基礎模型](#)。

G

生成式 AI

已針對大量資料進行訓練的 [AI](#) 模型子集，可使用簡單的文字提示建立新的內容和成品，例如影像、影片、文字和音訊。如需詳細資訊，請參閱[什麼是生成式 AI](#)。

地理封鎖

請參閱[地理限制](#)。

地理限制 (地理封鎖)

Amazon CloudFront 中的選項，可防止特定國家/地區的使用者存取內容分發。您可以使用允許清單或封鎖清單來指定核准和禁止的國家/地區。如需詳細資訊，請參閱 CloudFront 文件中的[限制內容的地理分佈](#)。

Gitflow 工作流程

這是一種方法，其中較低和較高環境在原始碼儲存庫中使用不同分支。Gitflow 工作流程被視為舊版，而以[幹線為基礎的工作流程](#)是現代、偏好的方法。

黃金影像

系統或軟體的快照，做為部署該系統或軟體新執行個體的範本。例如，在製造中，黃金映像可用於在多個裝置上佈建軟體，並有助於提高裝置製造操作的速度、可擴展性和生產力。

綠地策略

新環境中缺乏現有基礎設施。對系統架構採用綠地策略時，可以選擇所有新技術，而不會限制與現有基礎設施的相容性，也稱為[棕地](#)。如果正在擴展現有基礎設施，則可能會混合棕地和綠地策略。

防護機制

有助於跨組織單位 (OU) 來管控資源、政策和合規的高層級規則。預防性防護機制會強制執行政策，以確保符合合規標準。透過使用服務控制政策和 IAM 許可界限來將其實作。偵測性防護機制可偵測政策違規和合規問題，並產生提醒以便修正。它們是透過使用 AWS Config AWS Security Hub CSPM、Amazon GuardDuty、Amazon Inspector AWS Trusted Advisor 和自訂 AWS Lambda 檢查來實作。

H

HA

請參閱[高可用性](#)。

異質資料庫遷移

將來源資料庫遷移至使用不同資料庫引擎的目標資料庫 (例如，Oracle 至 Amazon Aurora)。異質遷移通常是重新架構工作的一部分，而轉換結構描述可能是一項複雜任務。[AWS 提供有助於結構描述轉換的 AWS SCT](#)。

高可用性 (HA)

在遇到挑戰或災難時，工作負載能夠在不介入的情況下持續運作。HA 系統的設計目的是自動容錯移轉、持續提供高品質的效能，並處理不同的負載和故障，並將效能影響降至最低。

歷史現代化

一種方法，用於現代化和升級操作技術 (OT) 系統，以更好地滿足製造業的需求。歷史資料是一種資料庫，用於從工廠中的各種來源收集和存放資料。

保留資料

從用於訓練機器學習模型的資料集中保留的部分歷史標記資料。您可以使用保留資料，透過比較模型預測與保留資料來評估模型效能。

異質資料庫遷移

將您的來源資料庫遷移至共用相同資料庫引擎的目標資料庫 (例如，Microsoft SQL Server 至 Amazon RDS for SQL Server)。同質遷移通常是主機轉換或平台轉換工作的一部分。您可以使用原生資料庫公用程式來遷移結構描述。

熱資料

經常存取的資料，例如即時資料或最近的轉譯資料。此資料通常需要高效能儲存層或類別，才能提供快速的查詢回應。

修補程序

緊急修正生產環境中的關鍵問題。由於其緊迫性，通常會在典型 DevOps 發行工作流程之外執行修補程式。

超級護理期間

在切換後，遷移團隊在雲端管理和監控遷移的應用程式以解決任何問題的時段。通常，此期間的長度為 1-4 天。在超級護理期間結束時，遷移團隊通常會將應用程式的責任轉移給雲端營運團隊。



IaC

將基礎設施視為程式碼。

身分型政策

連接至一或多個 IAM 主體的政策，可定義其在 AWS 雲端環境中的許可。

閒置應用程式

90 天期間 CPU 和記憶體平均使用率在 5% 至 20% 之間的應用程式。在遷移專案中，通常會淘汰這些應用程式或將其保留在內部部署。



IIoT

請參閱 [工業物聯網](#)。

不可變的基礎設施

為生產工作負載部署新基礎設施的模型，而不是更新、修補或修改現有的基礎設施。不可變基礎設施本質上比[可變基礎設施](#)更一致、可靠且可預測。如需詳細資訊，請參閱 AWS Well-Architected Framework 中的[使用不可變基礎設施部署最佳實務](#)。

傳入 (輸入) VPC

在 AWS 多帳戶架構中，接受、檢查和路由來自應用程式外部之網路連線的 VPC。[AWS 安全參考架構](#)建議您使用傳入、傳出和檢查 VPC 來設定網路帳戶，以保護應用程式與更廣泛的網際網路之間的雙向介面。

增量遷移

一種切換策略，您可以在其中將應用程式分成小部分遷移，而不是執行單一、完整的切換。例如，您最初可能只將一些微服務或使用者移至新系統。確認所有項目都正常運作之後，您可以逐步移動其他微服務或使用者，直到可以解除委任舊式系統。此策略可降低與大型遷移關聯的風險。

工業 4.0

2016 年 [Klaus Schwab](#) 推出的術語，透過連線能力、即時資料、自動化、分析和 AI/ML 的進展，指製造程序的現代化。

基礎設施

應用程式環境中包含的所有資源和資產。

基礎設施即程式碼 (IaC)

透過一組組態檔案來佈建和管理應用程式基礎設施的程序。IaC 旨在協助您集中管理基礎設施，標準化資源並快速擴展，以便新環境可重複、可靠且一致。

工業物聯網 (IIoT)

在製造業、能源、汽車、醫療保健、生命科學和農業等產業領域使用網際網路連線的感測器和裝置。如需詳細資訊，請參閱[建立工業物聯網 \(IIoT\) 數位轉型策略](#)。

檢查 VPC

在 AWS 多帳戶架構中，集中式 VPC，可管理 VPCs (在相同或不同的 AWS 區域)、網際網路和內部部署網路之間的網路流量檢查。[AWS 安全參考架構](#)建議您使用傳入、傳出和檢查 VPC 來設定網路帳戶，以保護應用程式與更廣泛的網際網路之間的雙向介面。

物聯網 (IoT)

具有內嵌式感測器或處理器的相連實體物體網路，其透過網際網路或本地通訊網路與其他裝置和系統進行通訊。如需詳細資訊，請參閱[什麼是 IoT？](#)

可解釋性

機器學習模型的一個特徵，描述了人類能夠理解模型的預測如何依賴於其輸入的程度。如需詳細資訊，請參閱[的機器學習模型可解釋性 AWS。](#)

IoT

請參閱[物聯網。](#)

IT 資訊庫 (ITIL)

一組用於交付 IT 服務並使這些服務與業務需求保持一致的最佳實務。ITIL 為 ITSM 提供了基礎。

IT 服務管理 (ITSM)

與組織的設計、實作、管理和支援 IT 服務關聯的活動。如需有關將雲端操作與 ITSM 工具整合的資訊，請參閱[操作整合指南。](#)

ITIL

請參閱[IT 資訊庫。](#)

ITSM

請參閱[IT 服務管理。](#)

L

標籤型存取控制 (LBAC)

強制存取控制 (MAC) 的實作，其中使用者和資料本身都會獲得明確指派的安全標籤值。使用者安全標籤和資料安全標籤之間的交集會決定使用者可以看到哪些資料列和資料欄。

登陸區域

登陸區域是架構良好的多帳戶 AWS 環境，可擴展且安全。這是一個起點，您的組織可以從此起點快速啟動和部署工作負載與應用程式，並對其安全和基礎設施環境充滿信心。如需有關登陸區域的詳細資訊，請參閱[設定安全且可擴展的多帳戶 AWS 環境。](#)

大型語言模型 (LLM)

預先訓練大量資料的深度學習 [AI](#) 模型。LLM 可以執行多個任務，例如回答問題、摘要文件、將文字翻譯成其他語言，以及完成句子。如需詳細資訊，請參閱[什麼是 LLMs](#)。

大型遷移

遷移 300 部或更多伺服器。

LBAC

請參閱[標籤型存取控制](#)。

最低權限

授予執行任務所需之最低許可的安全最佳實務。如需詳細資訊，請參閱 IAM 文件中的[套用最低權限許可](#)。

隨即轉移

請參閱[7 個 R](#)。

小端序系統

首先儲存最低有效位元組的系統。另請參閱[Endianness](#)。

LLM

請參閱[大型語言模型](#)。

較低的環境

請參閱[環境](#)。

M

機器學習 (ML)

一種使用演算法和技術進行模式識別和學習的人工智慧。機器學習會進行分析並從記錄的資料 (例如物聯網 (IoT) 資料) 中學習，以根據模式產生統計模型。如需詳細資訊，請參閱[機器學習](#)。

主要分支

請參閱[分支](#)。

惡意軟體

旨在危及電腦安全或隱私權的軟體。惡意軟體可能會中斷電腦系統、洩露敏感資訊，或取得未經授權的存取。惡意軟體的範例包括病毒、蠕蟲、勒索軟體、特洛伊木馬、間諜軟體和鍵盤記錄器。

受管服務

AWS 服務會 AWS 操作基礎設施層、作業系統和平台，而您會存取端點來存放和擷取資料。Amazon Simple Storage Service (Amazon S3) 和 Amazon DynamoDB 是受管服務的範例。這些也稱為抽象服務。

製造執行系統 (MES)

一種軟體系統，用於追蹤、監控、記錄和控制生產程序，將原物料轉換為現場成品。

MAP

請參閱遷移加速計劃。

機制

建立工具、推動工具採用，然後檢查結果以進行調整的完整程序。機制是在操作時強化和改善自身的循環。如需詳細資訊，請參閱 AWS Well-Architected Framework 中的建置機制。

成員帳戶

屬於組織一部分的管理帳戶 AWS 帳戶以外的所有 AWS Organizations。帳戶一次只能是一個組織的成員。

製造執行系統

請參閱製造執行系統。

訊息佇列遙測傳輸 (MQTT)

根據發佈/訂閱模式的輕量型 machine-to-machine (M2M) 通訊協定，適用於資源受限的 IoT 裝置。

微服務

一種小型的獨立服務，它可透過定義明確的 API 進行通訊，通常由小型獨立團隊擁有。例如，保險系統可能包含對應至業務能力（例如銷售或行銷）或子領域（例如購買、索賠或分析）的微服務。微服務的優點包括靈活性、彈性擴展、輕鬆部署、可重複使用的程式碼和適應力。如需詳細資訊，請參閱使用無 AWS 伺服器服務整合微服務。

微服務架構

一種使用獨立元件來建置應用程式的方法，這些元件會以微服務形式執行每個應用程式程序。這些微服務會使用輕量型 API，透過明確定義的介面進行通訊。此架構中的每個微服務都可以進行

更新、部署和擴展，以滿足應用程式特定功能的需求。如需詳細資訊，請參閱在上實作微服務 AWS。

Migration Acceleration Program (MAP)

一種 AWS 計畫，提供諮詢支援、訓練和服務，協助組織建立強大的營運基礎，以移至雲端，並協助抵銷遷移的初始成本。MAP 包括用於有條不紊地執行舊式遷移的遷移方法以及一組用於自動化和加速常見遷移案例的工具。

大規模遷移

將大部分應用程式組合依波次移至雲端的程序，在每個波次中，都會以更快的速度移動更多應用程式。此階段使用從早期階段學到的最佳實務和經驗教訓來實作團隊、工具和流程的遷移工廠，以透過自動化和敏捷交付簡化工作負載的遷移。這是 AWS 遷移策略的第三階段。

遷移工廠

可透過自動化、敏捷的方法簡化工作負載遷移的跨職能團隊。遷移工廠團隊通常包括營運、業務分析師和擁有者、遷移工程師、開發人員以及從事 Sprint 工作的 DevOps 專業人員。20% 至 50% 之間的企業應用程式組合包含可透過工廠方法優化的重複模式。如需詳細資訊，請參閱此內容集中的遷移工廠的討論和雲端遷移工廠指南。

遷移中繼資料

有關完成遷移所需的應用程式和伺服器的資訊。每種遷移模式都需要一組不同的遷移中繼資料。遷移中繼資料的範例包括目標子網路、安全群組和 AWS 帳戶。

遷移模式

可重複的遷移任務，詳細描述遷移策略、遷移目的地以及所使用的遷移應用程式或服務。範例：使用 AWS Application Migration Service 重新託管遷移至 Amazon EC2。

遷移組合評定 (MPA)

線上工具，提供驗證商業案例以遷移至 的資訊 AWS 雲端。MPA 提供詳細的組合評定 (伺服器適當規模、定價、總體擁有成本比較、遷移成本分析) 以及遷移規劃 (應用程式資料分析和資料收集、應用程式分組、遷移優先順序，以及波次規劃)。MPA 工具 (需要登入) 可供所有 AWS 顧問和 APN 合作夥伴顧問免費使用。

遷移準備程度評定 (MRA)

使用 AWS CAF 取得組織雲端整備狀態的洞見、識別優缺點，以及建立行動計劃以消除已識別差距的程序。如需詳細資訊，請參閱遷移準備程度指南。MRA 是 AWS 遷移策略的第一階段。

遷移策略

用來將工作負載遷移至的方法 AWS 雲端。如需詳細資訊，請參閱本詞彙表中的 [7 個 Rs 項目](#)，並請參閱[動員您的組織以加速大規模遷移](#)。

機器學習 (ML)

請參閱[機器學習](#)。

現代化

將過時的 (舊版或單一) 應用程式及其基礎架構轉換為雲端中靈活、富有彈性且高度可用的系統，以降低成本、提高效率並充分利用創新。如需詳細資訊，請參閱[《》中的現代化應用程式的策略 AWS 雲端](#)。

現代化準備程度評定

這項評估可協助判斷組織應用程式的現代化準備程度；識別優點、風險和相依性；並確定組織能夠在多大程度上支援這些應用程式的未來狀態。評定的結果就是目標架構的藍圖、詳細說明現代化程序的開發階段和里程碑的路線圖、以及解決已發現的差距之行動計畫。如需詳細資訊，請參閱[《》中的評估應用程式的現代化準備 AWS 雲端程度](#)。

單一應用程式 (單一)

透過緊密結合的程序作為單一服務執行的應用程式。單一應用程式有幾個缺點。如果一個應用程式功能遇到需求激增，則必須擴展整個架構。當程式碼庫增長時，新增或改進單一應用程式的功能也會變得更加複雜。若要解決這些問題，可以使用微服務架構。如需詳細資訊，請參閱[將單一體系分解為微服務](#)。

MPA

請參閱[遷移產品組合評估](#)。

MQTT

請參閱[訊息併列遙測傳輸](#)。

多類別分類

一個有助於產生多類別預測的過程 (預測兩個以上的結果之一)。例如，機器學習模型可能會詢問「此產品是書籍、汽車還是電話？」或者「這個客戶對哪種產品類別最感興趣？」

可變基礎設施

更新和修改生產工作負載現有基礎設施的模型。為了提高一致性、可靠性和可預測性，AWS Well-Architected Framework 建議使用[不可變的基礎設施](#)作為最佳實務。

O

OAC

請參閱 [原始存取控制](#)。

OAI

請參閱 [原始存取身分](#)。

OCM

請參閱 [組織變更管理](#)。

離線遷移

一種遷移方法，可在遷移過程中刪除來源工作負載。此方法涉及延長停機時間，通常用於小型非關鍵工作負載。

OI

請參閱 [操作整合](#)。

OLA

請參閱 [操作層級協議](#)。

線上遷移

一種遷移方法，無需離線即可將來源工作負載複製到目標系統。連接至工作負載的應用程式可在遷移期間繼續運作。此方法涉及零至最短停機時間，通常用於關鍵的生產工作負載。

OPC-UA

請參閱 [開啟程序通訊 - 統一架構](#)。

開放程序通訊 - 統一架構 (OPC-UA)

用於工業自動化的machine-to-machine(M2M) 通訊協定。OPC-UA 提供資料加密、身分驗證和授權機制的互通性標準。

操作水準協議 (OLA)

一份協議，闡明 IT 職能群組承諾向彼此提供的內容，以支援服務水準協議 (SLA)。

操作整備審查 (ORR)

問題和相關最佳實務的檢查清單，可協助您了解、評估、預防或減少事件和可能失敗的範圍。如需詳細資訊，請參閱 AWS Well-Architected Framework 中的[操作準備審查 \(ORR\)](#)。

操作技術 (OT)

使用實體環境控制工業操作、設備和基礎設施的硬體和軟體系統。在製造中，整合 OT 和資訊技術 (IT) 系統是[工業 4.0](#)轉型的關鍵重點。

操作整合 (OI)

在雲端中將操作現代化的程序，其中包括準備程度規劃、自動化和整合。如需詳細資訊，請參閱[操作整合指南](#)。

組織追蹤

由建立的線索 AWS CloudTrail，會記錄 AWS 帳戶 組織中所有 的所有事件 AWS Organizations。在屬於組織的每個 AWS 帳戶 中建立此追蹤，它會跟蹤每個帳戶中的活動。如需詳細資訊，請參閱 CloudTrail 文件中的[建立組織追蹤](#)。

組織變更管理 (OCM)

用於從人員、文化和領導力層面管理重大、顛覆性業務轉型的架構。OCM 透過加速變更採用、解決過渡問題，以及推動文化和組織變更，協助組織為新系統和策略做好準備，並轉移至新系統和策略。在 AWS 遷移策略中，此架構稱為人員加速，因為雲端採用專案所需的變更速度。如需詳細資訊，請參閱 [OCM 指南](#)。

原始存取控制 (OAC)

CloudFront 中的增強型選項，用於限制存取以保護 Amazon Simple Storage Service (Amazon S3) 內容。OAC 支援所有 S3 儲存貯體中的所有伺服器端加密 AWS KMS (SSE-KMS) AWS 區域，以及對 S3 儲存貯體的動態PUT和DELETE請求。

原始存取身分 (OAI)

CloudFront 中的一個選項，用於限制存取以保護 Amazon S3 內容。當您使用 OAI 時，CloudFront 會建立一個可供 Amazon S3 進行驗證的主體。經驗證的主體只能透過特定 CloudFront 分發來存取 S3 儲存貯體中的內容。另請參閱 [OAC](#)，它可提供更精細且增強的存取控制。

ORR

請參閱[操作整備審核](#)。

OT

請參閱[操作技術](#)。

傳出 (輸出) VPC

在 AWS 多帳戶架構中，處理從應用程式內啟動之網路連線的 VPC。[AWS 安全參考架構](#)建議您使用傳入、傳出和檢查 VPC 來設定網路帳戶，以保護應用程式與更廣泛的網際網路之間的雙向介面。

P

許可界限

附接至 IAM 主體的 IAM 管理政策，可設定使用者或角色擁有的最大許可。如需詳細資訊，請參閱 IAM 文件中的[許可界限](#)。

個人身分識別資訊 (PII)

當直接檢視或與其他相關資料配對時，可用來合理推斷個人身分的資訊。PII 的範例包括名稱、地址和聯絡資訊。

PII

請參閱[個人身分識別資訊](#)。

手冊

一組預先定義的步驟，可擷取與遷移關聯的工作，例如在雲端中提供核心操作功能。手冊可以採用指令碼、自動化執行手冊或操作現代化環境所需的程序或步驟摘要的形式。

PLC

請參閱[可程式設計邏輯控制器](#)。

PLM

請參閱[產品生命週期管理](#)。

政策

可定義許可的物件（請參閱[身分型政策](#)）、指定存取條件（請參閱[資源型政策](#)），或定義組織中所有帳戶的最大許可 AWS Organizations（請參閱[服務控制政策](#)）。

混合持久性

根據資料存取模式和其他需求，獨立選擇微服務的資料儲存技術。如果您的微服務具有相同的資料儲存技術，則其可能會遇到實作挑戰或效能不佳。如果微服務使用最適合其需求的資料儲存，則

可以更輕鬆地實作並達到更好的效能和可擴展性。如需詳細資訊，請參閱[在微服務中啟用資料持久性](#)。

組合評定

探索、分析應用程式組合並排定其優先順序以規劃遷移的程序。如需詳細資訊，請參閱[評估遷移準備程度](#)。

述詞

傳回 true 或 false 的查詢條件，通常位於 WHERE 子句中。

述詞下推

一種資料庫查詢最佳化技術，可在傳輸前篩選查詢中的資料。這可減少必須從關聯式資料庫擷取和處理的資料量，並改善查詢效能。

預防性控制

旨在防止事件發生的安全控制。這些控制是第一道防線，可協助防止對網路的未經授權存取或不必變更。如需詳細資訊，請參閱在 AWS 上實作安全控制中的[預防性控制](#)。

委託人

中可執行動作和存取資源 AWS 的實體。此實體通常是 AWS 帳戶、IAM 角色或使用者的根使用者。如需詳細資訊，請參閱 IAM 文件中[角色術語和概念](#)中的主體。

依設計的隱私權

透過整個開發程序將隱私權納入考量的系統工程方法。

私有託管區域

一種容器，它包含有關您希望 Amazon Route 53 如何回應一個或多個 VPC 內的域及其子域之 DNS 查詢的資訊。如需詳細資訊，請參閱 Route 53 文件中的[使用私有託管區域](#)。

主動控制

旨在防止部署不合規資源[的安全控制](#)。這些控制項會在佈建資源之前對其進行掃描。如果資源不符合控制項，則不會佈建。如需詳細資訊，請參閱 AWS Control Tower 文件中的[控制項參考指南](#)，並參閱實作安全控制項中的主動控制項。 AWS

產品生命週期管理 (PLM)

管理產品整個生命週期的資料和程序，從設計、開發和啟動，到成長和成熟，再到拒絕和移除。

生產環境

請參閱[環境](#)。

可程式設計邏輯控制器 (PLC)

在製造中，高度可靠、可調整的電腦，可監控機器並自動化製造程序。

提示鏈結

使用一個 [LLM](#) 提示的輸出做為下一個提示的輸入，以產生更好的回應。此技術用於將複雜任務分解為子任務，或反覆精簡或展開初步回應。它有助於提高模型回應的準確性和相關性，並允許更精細、個人化的結果。

擬匿名化

將資料集中的個人識別符取代為預留位置值的程序。假名化有助於保護個人隱私權。假名化資料仍被視為個人資料。

發佈/訂閱 (pub/sub)

一種模式，可啟用微服務之間的非同步通訊，以提高可擴展性和回應能力。例如，在微服務型 [MES](#) 中，微服務可以將事件訊息發佈到其他微服務可訂閱的頻道。系統可以新增新的微服務，而無需變更發佈服務。

Q

查詢計劃

一系列步驟，如指示，用於存取 SQL 關聯式資料庫系統中的資料。

查詢計劃迴歸

在資料庫服務優化工具選擇的計畫比對資料庫環境進行指定的變更之前的計畫不太理想時。這可能因為對統計資料、限制條件、環境設定、查詢參數繫結的變更以及資料庫引擎的更新所導致。

R

RACI 矩陣

請參閱 [負責、負責、諮詢、告知 \(RACI\)](#)。

RAG

請參閱 [擷取增強生成](#)。

勒索軟體

一種惡意軟體，旨在阻止對計算機系統或資料的存取，直到付款為止。

RASCI 矩陣

請參閱[負責、負責、諮詢、告知 \(RACI\)](#)。

RCAC

請參閱[資料列和資料欄存取控制](#)。

僅供讀取複本

用於唯讀用途的資料庫複本。您可以將查詢路由至僅供讀取複本以減少主資料庫的負載。

重新架構師

請參閱[7 個 R](#)。

復原點目標 (RPO)

自上次資料復原點以來可接受的時間上限。這會決定最後一個復原點與服務中斷之間可接受的資料遺失。

復原時間目標 (RTO)

服務中斷與服務還原之間的可接受延遲上限。

重構

請參閱[7 個 R](#)。

區域

地理區域中的 AWS 資源集合。每個 AWS 區域 都獨立於其他，以提供容錯能力、穩定性和彈性。
如需詳細資訊，請參閱[指定 AWS 區域 您的帳戶可以使用哪些](#)。

迴歸

預測數值的 ML 技術。例如，為了解決「這房子會賣什麼價格？」的問題 ML 模型可以使用線性迴歸模型，根據已知的房屋事實（例如，平方英尺）來預測房屋的銷售價格。

重新託管

請參閱[7 個 R](#)。

版本

在部署程序中，它是將變更提升至生產環境的動作。

重新定位

請參閱 [7 個 R。](#)

Replatform

請參閱 [7 個 R。](#)

回購

請參閱 [7 個 R。](#)

彈性

應用程式抵禦中斷或從中斷中復原的能力。[在 中 規劃彈性時，高可用性和災難復原](#)是常見的考量 AWS 雲端。如需詳細資訊，請參閱[AWS 雲端彈性](#)。

資源型政策

附接至資源的政策，例如 Amazon S3 儲存貯體、端點或加密金鑰。這種類型的政策會指定允許存取哪些主體、支援的動作以及必須滿足的任何其他條件。

負責者、當責者、事先諮詢者和事後告知者 (RACI) 矩陣

定義所有涉及遷移活動和雲端操作之各方的角色和責任的矩陣。矩陣名稱衍生自矩陣中定義的責任類型：負責人 (R)、責任 (A)、諮詢 (C) 和知情 (I)。支援 (S) 類型為選用。如果您包含支援，則矩陣稱為 RASCI 矩陣，如果您排除它，則稱為 RACI 矩陣。

回應性控制

一種安全控制，旨在驅動不良事件或偏離安全基準的補救措施。如需詳細資訊，請參閱在 AWS 上實作安全控制中的[回應性控制](#)。

保留

請參閱 [7 Rs。](#)

淘汰

請參閱 [7 個 R。](#)

檢索增強生成 (RAG)

一種生成式 AI 技術，其中 [LLM](#) 會在產生回應之前參考訓練資料來源以外的授權資料來源。例如，RAG 模型可能會對組織的知識庫或自訂資料執行語意搜尋。如需詳細資訊，請參閱[什麼是 RAG](#)。

輪換

定期更新[秘密](#)的程序，讓攻擊者更難存取登入資料。

資料列和資料欄存取控制 (RCAC)

使用已定義存取規則的基本、彈性 SQL 表達式。RCAC 包含資料列許可和資料欄遮罩。

RPO

請參閱[復原點目標](#)。

RTO

請參閱[復原時間目標](#)。

執行手冊

執行特定任務所需的一組手動或自動程序。這些通常是為了簡化重複性操作或錯誤率較高的程序而建置。

S

SAML 2.0

許多身分提供者 (IdP) 使用的開放標準。此功能可啟用聯合單一登入 (SSO)，讓使用者可以登入 AWS 管理主控台 或呼叫 AWS API 操作，而無需為您組織中的每個人在 IAM 中建立使用者。如需有關以 SAML 2.0 為基礎的聯合詳細資訊，請參閱 IAM 文件中的[關於以 SAML 2.0 為基礎的聯合](#)。

SCADA

請參閱[監督控制和資料擷取](#)。

SCP

請參閱[服務控制政策](#)。

秘密

您以加密形式存放的 AWS Secrets Manager 機密或限制資訊，例如密碼或使用者登入資料。它由秘密值及其中繼資料組成。秘密值可以是二進位、單一字串或多個字串。如需詳細資訊，請參閱 [Secrets Manager 文件中的 Secrets Manager 秘密中的什麼內容？](#)。

依設計的安全性

透過整個開發程序將安全性納入考量的系統工程方法。

安全控制

一種技術或管理防護機制，它可預防、偵測或降低威脅行為者利用安全漏洞的能力。安全控制有四種主要類型：[預防性](#)、[偵測性](#)、[回應性](#)和[主動性](#)。

安全強化

減少受攻擊面以使其更能抵抗攻擊的過程。這可能包括一些動作，例如移除不再需要的資源、實作授予最低權限的安全最佳實務、或停用組態檔案中不必要的功能。

安全資訊與事件管理 (SIEM) 系統

結合安全資訊管理 (SIM) 和安全事件管理 (SEM) 系統的工具與服務。SIEM 系統會收集、監控和分析來自伺服器、網路、裝置和其他來源的資料，以偵測威脅和安全漏洞，並產生提醒。

安全回應自動化

預先定義和程式設計的動作，旨在自動回應或修復安全事件。這些自動化可做為偵測或回應式安全控制，協助您實作 AWS 安全最佳實務。自動化回應動作的範例包括修改 VPC 安全群組、修補 Amazon EC2 執行個體或輪換登入資料。

伺服器端加密

由 AWS 服務 接收資料的 在其目的地加密資料。

服務控制政策 (SCP)

為 AWS Organizations 中的組織的所有帳戶提供集中控制許可的政策。SCP 會定義防護機制或設定管理員可委派給使用者或角色的動作限制。您可以使用 SCP 作為允許清單或拒絕清單，以指定允許或禁止哪些服務或動作。如需詳細資訊，請參閱 AWS Organizations 文件中的服務控制政策。

服務端點

的進入點 URL AWS 服務。您可以使用端點，透過程式設計方式連接至目標服務。如需詳細資訊，請參閱 AWS 一般參考 中的 [AWS 服務 端點](#)。

服務水準協議 (SLA)

一份協議，闡明 IT 團隊承諾向客戶提供的服務，例如服務正常執行時間和效能。

服務層級指標 (SLI)

服務效能方面的測量，例如其錯誤率、可用性或輸送量。

服務層級目標 (SLO)

代表服務運作狀態的目標指標，由服務層級指標測量。

共同責任模式

描述您與 共同 AWS 承擔雲端安全與合規責任的模型。 AWS 負責雲端的安全，而 負責雲端的安全。如需詳細資訊，請參閱[共同責任模式](#)。

SIEM

請參閱[安全資訊和事件管理系統](#)。

單一故障點 (SPOF)

應用程式的單一關鍵元件故障，可能會中斷系統。

SLA

請參閱[服務層級協議](#)。

SLI

請參閱[服務層級指標](#)。

SLO

請參閱[服務層級目標](#)。

先拆分後播種模型

擴展和加速現代化專案的模式。定義新功能和產品版本時，核心團隊會進行拆分以建立新的產品團隊。這有助於擴展組織的能力和服務，提高開發人員生產力，並支援快速創新。如需詳細資訊，請參閱[中的階段式應用程式現代化方法 AWS 雲端](#)。

SPOF

請參閱[單一故障點](#)。

星狀結構描述

使用一個大型事實資料表來存放交易或測量資料的資料庫組織結構，並使用一或多個較小的維度資料表來存放資料屬性。此結構旨在用於[資料倉儲](#)或商業智慧用途。

Strangler Fig 模式

一種現代化單一系統的方法，它會逐步重寫和取代系統功能，直到舊式系統停止使用為止。此模式源自無花果藤，它長成一棵馴化樹並最終戰勝且取代了其宿主。該模式由[Martin Fowler 引入](#)，作為重寫單一系統時管理風險的方式。如需有關如何套用此模式的範例，請參閱[使用容器和 Amazon API Gateway 逐步現代化舊版 Microsoft ASP.NET \(ASMX\) Web 服務](#)。

子網

您 VPC 中的 IP 地址範圍。子網必須位於單一可用區域。

監控控制和資料擷取 (SCADA)

在製造中，使用硬體和軟體來監控實體資產和生產操作的系統。

對稱加密

使用相同金鑰來加密及解密資料的加密演算法。

合成測試

以模擬使用者互動的方式測試系統，以偵測潛在問題或監控效能。您可以使用 [Amazon CloudWatch Synthetics](#) 來建立這些測試。

系統提示

一種向 [LLM](#) 提供內容、指示或指導方針以指示其行為的技術。系統提示有助於設定內容，並建立與使用者互動的規則。

T

標籤

做為中繼資料以組織 AWS 資源的鍵值對。標籤可協助您管理、識別、組織、搜尋及篩選資源。如需詳細資訊，請參閱 [標記您的 AWS 資源](#)。

目標變數

您嘗試在受監督的 ML 中預測的值。這也被稱為結果變數。例如，在製造設定中，目標變數可能是產品瑕疵。

任務清單

用於透過執行手冊追蹤進度的工具。任務清單包含執行手冊的概觀以及要完成的一般任務清單。對於每個一般任務，它包括所需的預估時間量、擁有者和進度。

測試環境

請參閱 [環境](#)。

訓練

為 ML 模型提供資料以供學習。訓練資料必須包含正確答案。學習演算法會在訓練資料中尋找將輸入資料屬性映射至目標的模式（您想要預測的答案）。它會輸出擷取這些模式的 ML 模型。可以使用 ML 模型，來預測您不知道的目標新資料。

傳輸閘道

可以用於互連 VPC 和內部部署網路的網路傳輸中樞。如需詳細資訊，請參閱 AWS Transit Gateway 文件中的 [什麼是傳輸閘道](#)。

主幹型工作流程

這是一種方法，開發人員可在功能分支中本地建置和測試功能，然後將這些變更合併到主要分支中。然後，主要分支會依序建置到開發環境、生產前環境和生產環境中。

受信任的存取權

將許可授予您指定的服務，以代表您在組織中 AWS Organizations 及其帳戶中執行任務。受信任的服務會在需要該角色時，在每個帳戶中建立服務連結角色，以便為您執行管理工作。如需詳細資訊，請參閱文件中的 AWS Organizations [搭配使用 AWS Organizations 與其他 AWS 服務](#)。

調校

變更訓練程序的各個層面，以提高 ML 模型的準確性。例如，可以透過產生標籤集、新增標籤、然後在不同的設定下多次重複這些步驟來訓練 ML 模型，以優化模型。

雙比薩團隊

兩個比薩就能吃飽的小型 DevOps 團隊。雙披薩團隊規模可確保軟體開發中的最佳協作。

U

不確定性

這是一個概念，指的是不精確、不完整或未知的資訊，其可能會破壞預測性 ML 模型的可靠性。有兩種類型的不確定性：認知不確定性是由有限的、不完整的資料引起的，而隨機不確定性是由資料中固有的噪聲和隨機性引起的。如需詳細資訊，請參閱[量化深度學習系統的不確定性](#)指南。

未區分的任務

也稱為繁重工作，這是建立和操作應用程式的必要工作，但不為最終使用者提供直接價值或提供競爭優勢。未區分任務的範例包括採購、維護和容量規劃。

較高的環境

請參閱 [環境](#)。

V

清空

一種資料庫維護操作，涉及增量更新後的清理工作，以回收儲存並提升效能。

版本控制

追蹤變更的程序和工具，例如儲存庫中原始程式碼的變更。

VPC 對等互連

兩個 VPC 之間的連線，可讓您使用私有 IP 地址路由流量。如需詳細資訊，請參閱 Amazon VPC 文件中的[什麼是 VPC 對等互連](#)。

漏洞

危害系統安全性的軟體或硬體瑕疵。

W

暖快取

包含經常存取的目前相關資料的緩衝快取。資料庫執行個體可以從緩衝快取讀取，這比從主記憶體或磁碟讀取更快。

暖資料

不常存取的資料。查詢這類資料時，通常可接受中等速度的查詢。

視窗函數

SQL 函數，對與目前記錄在某種程度上相關的資料列群組執行計算。視窗函數適用於處理任務，例如根據目前資料列的相對位置計算移動平均值或存取資料列的值。

工作負載

提供商業價值的資源和程式碼集合，例如面向客戶的應用程式或後端流程。

工作串流

遷移專案中負責一組特定任務的功能群組。每個工作串流都是獨立的，但支援專案中的其他工作串流。例如，組合工作串流負責排定應用程式、波次規劃和收集遷移中繼資料的優先順序。組合工作串流將這些資產交付至遷移工作串流，然後再遷移伺服器和應用程式。

WORM

請參閱[寫入一次，多次讀取](#)。

WQF

請參閱[AWS 工作負載資格架構](#)。

寫入一次，讀取許多 (WORM)

儲存模型，可一次性寫入資料，並防止刪除或修改資料。授權使用者可以視需要多次讀取資料，但無法變更資料。此資料儲存基礎設施被視為不可變。

Z

零時差入侵

利用零時差漏洞的攻擊，通常是惡意軟體。

零時差漏洞

生產系統中未緩解的缺陷或漏洞。威脅行為者可以使用這種類型的漏洞來攻擊系統。開發人員經常因為攻擊而意識到漏洞。

零鏡頭提示

提供 LLM 執行任務的指示，但沒有可協助引導任務的範例 (快照)。LLM 必須使用其預先訓練的知識來處理任務。零鏡頭提示的有效性取決於任務的複雜性和提示的品質。另請參閱少量擷取提示。

殞屍應用程式

CPU 和記憶體平均使用率低於 5% 的應用程式。在遷移專案中，通常會淘汰這些應用程式。

本文為英文版的機器翻譯版本，如內容有任何歧義或不一致之處，概以英文版為準。