

AWS でできる耐障害性の向上と カオスエンジニアリングのすすめ

菊池之裕

Sr. Solutions Architect Network Specialist
Amazon Web Services Japan G.K.



自己紹介



- 名前：菊池 之裕(きくち ゆきひろ)
- シニアソリューションアーキテクト ネットワークスペシャリスト
- ロール：Network系サービスについてのご支援
- 経歴：ISP,IXP,VPN運用、開発を経てネットワーク機器、仮想ルータ販売会社のプリセールス、プロダクトSEからAWSへ（8年目）
- 好きな AWS サービス:AWS Transit Gateway,AWS Direct Connect, AWS Marketplace



Agenda

- AWS概要
- グレー障害について
- グレー障害の検出
- グレー障害のAWSでの対応例
- カオスエンジニアリングと障害の訓練について
- まとめ

今日持って帰っていただきたいこと

- グレー障害とは何かを理解する
- 耐障害性をあげる方法について理解する
- カオスエンジニアリングと、訓練の方法について理解する

AWSとは何か？

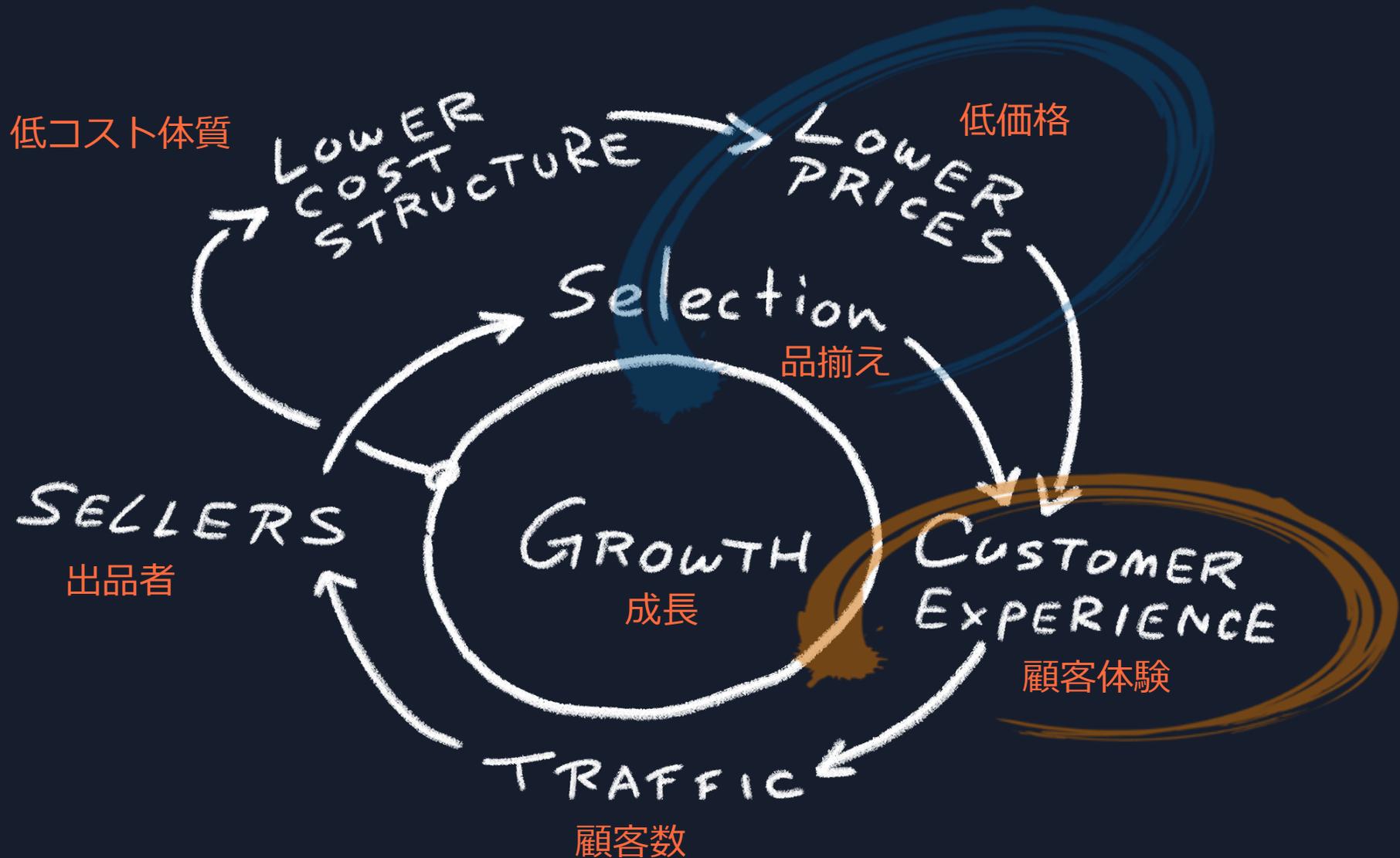
1. Amazonの課題から生まれたAWS
2. AWSの特徴



Our mission

地球上で、もっともお客様を大切にする企業であること

Amazon のビジネスモデル



Amazon の IT が抱えた課題の克服が AWS に

amazon



aws

課題

- ビジネスの**進化のスピード**に**基盤が追従しない**
- お客様からのフィードバック（市場のニーズ）をサービスの充実化に繋がられない

対応策

- 多目的・多機能の巨大なアプリケーションを**機能ごとにサービスとして細分化**（**マイクロサービス化**）
- **サービス間の依存関係を排除**し、個々に拡張・進化可能
- API ベースで運用できる**基盤に刷新**（自動化）

この基盤を元にして、
2006年に一般向けサービスとして
Amazon Web Services が事業を開始

クラウドの真価とは価値創造に集中できること



スモールスタートで
すぐに使い始められる



必要なときに必要なだけ
使うことが可能



アイディアから
実装までの時間を短縮

200 を超えるクラウドサービスで あらゆるワークロードをサポート

-  コンピューティング
-  モバイル
-  AR と VR
-  エンドユーザーコンピューティング
-  ストレージ
-  データベース
-  ネットワークとコンテンツ配信
-  AWS コスト管理
-  機械学習と AI
-  IoT
-  ロボット工学
-  ビジネスアプリケーション
-  メディアサービス
-  分析
-  マネジメントとガバナンス
-  開発者用ツール
-  サーバーレス
-  アプリケーション統合
-  Game Tech
-  量子テクノロジー
-  カスタマーイネーブルメント
-  移行と転送
-  ブロックチェーン
-  セキュリティ・ID・コンプライアンス
-  人工衛星
-  コンテナ

グレー障害の定義と発生シナリオ



気づきにくいユーザー体験の低下

- ある EC サイトで、一部のユーザーのみカート操作が失敗
- 商品は表示されているのに、カートに入れようとすると時々エラーになる
- テストユーザーは正常であり、運用チームによる動作確認では問題を再現できない

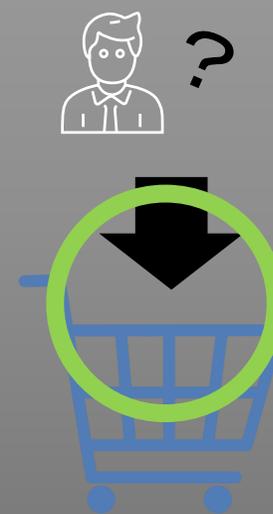
一部の一般ユーザー



他の一般ユーザー



テストユーザー

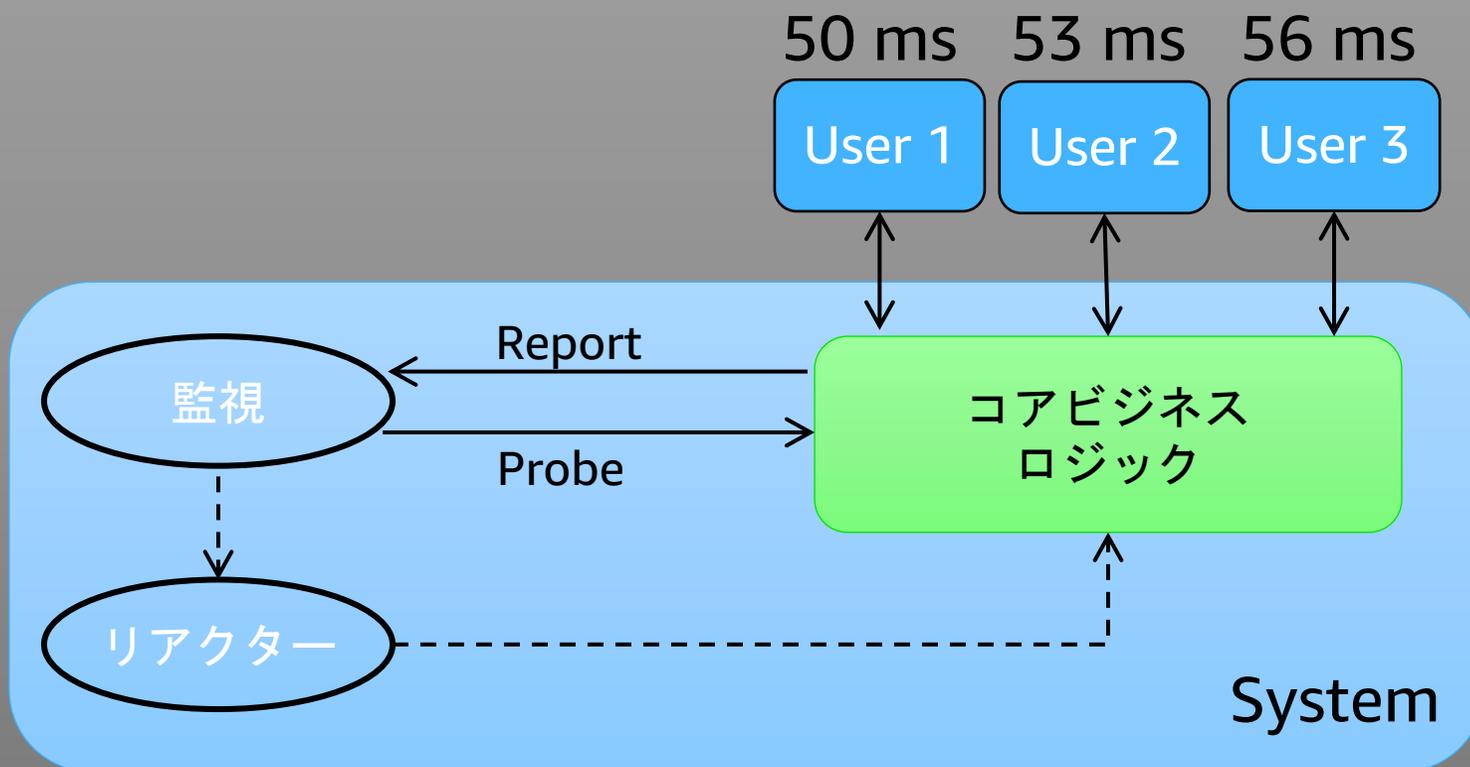


グレー障害

完全な停止ではないが、サービスの質が低下している状態

Average latency :
53 ms

Latency threshold :
60 ms

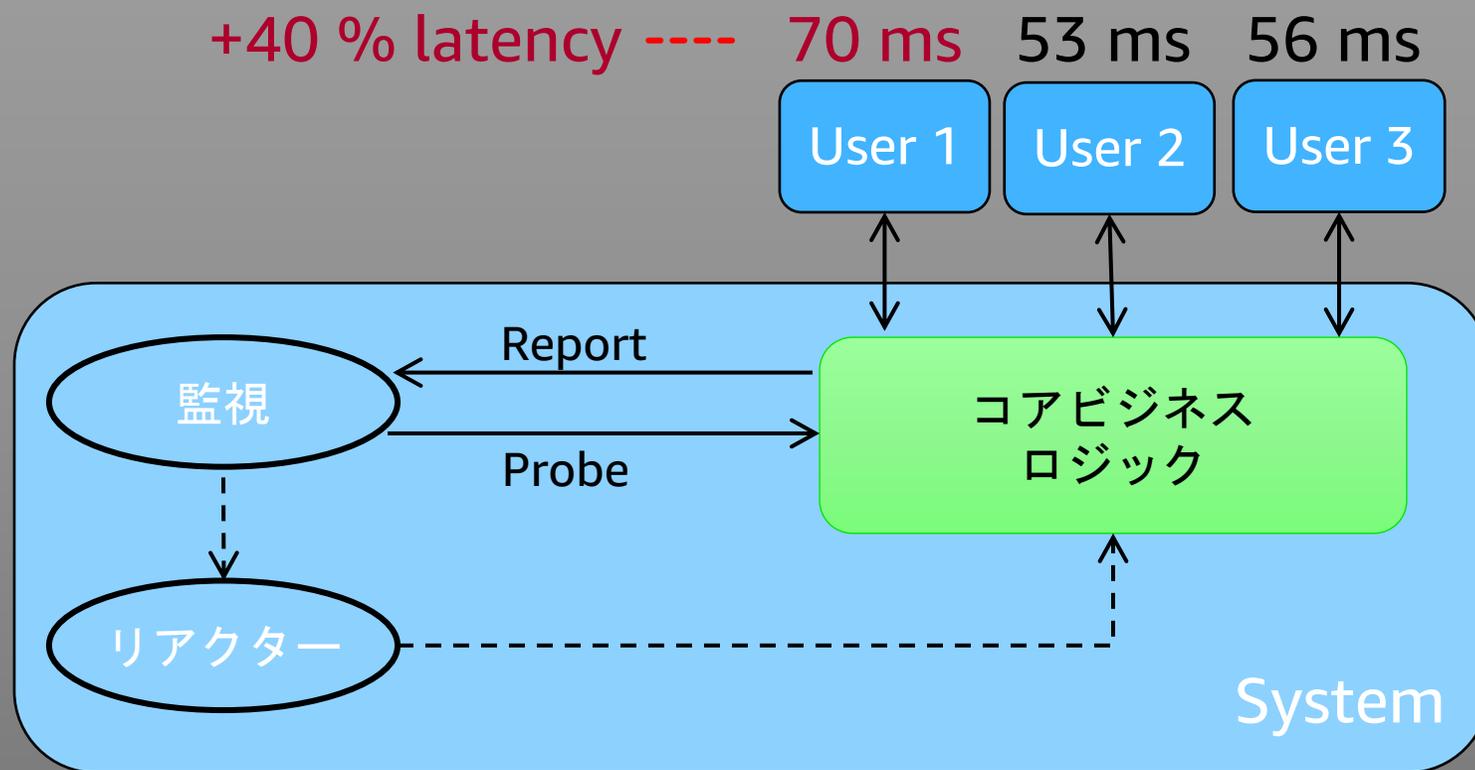


グレー障害

完全な停止ではないが、サービスの質が低下している状態

Average latency :
59.66 ms

Latency threshold :
60 ms



視点別の Observability

システムの状態は、誰がどこから見ているかによって、異なって見えることがある



ALL GOOD :

すべての視点で正常な状態

GRAY FAILURE (グレー障害) :

システム全体としては正常に見えるが、一部のユーザーやアプリケーションで問題が発生

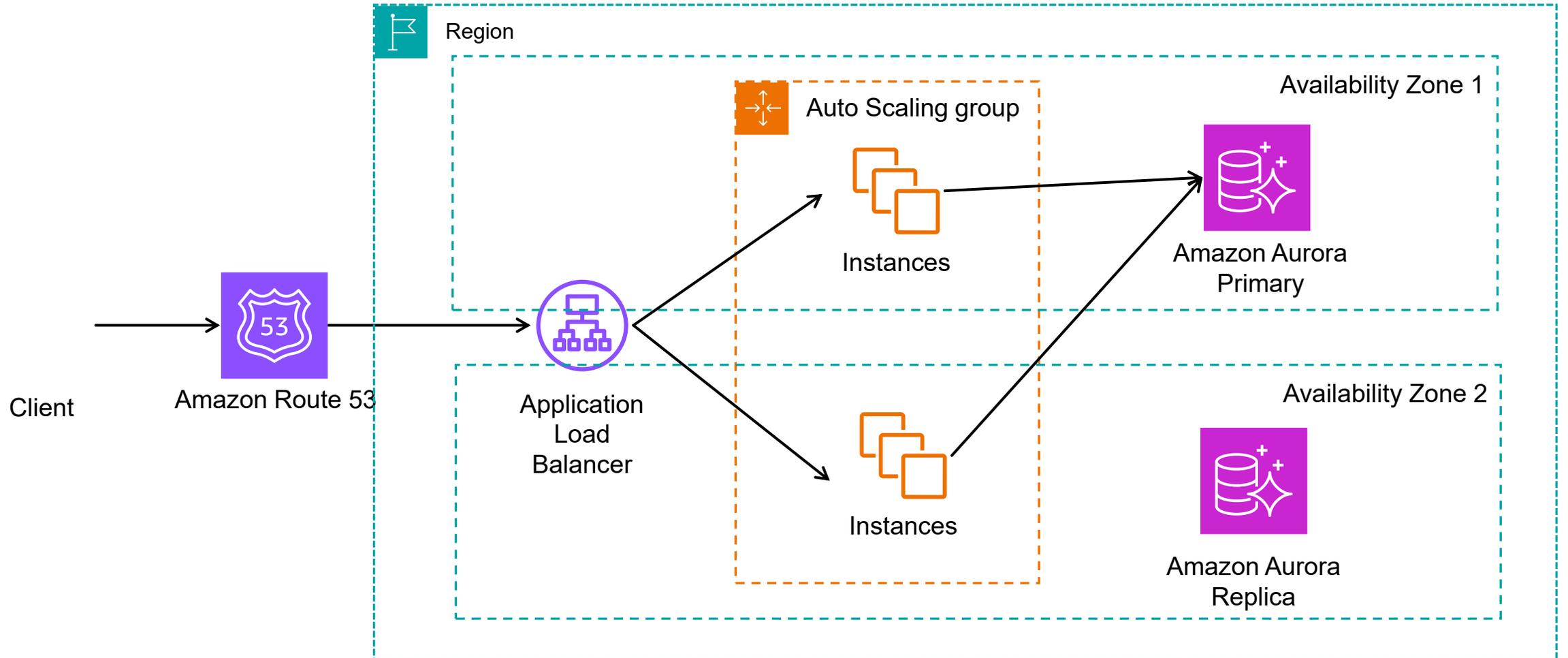
DETECTED FAILURE :

システムとして問題を検知できている状態

MASKED FAILURE :

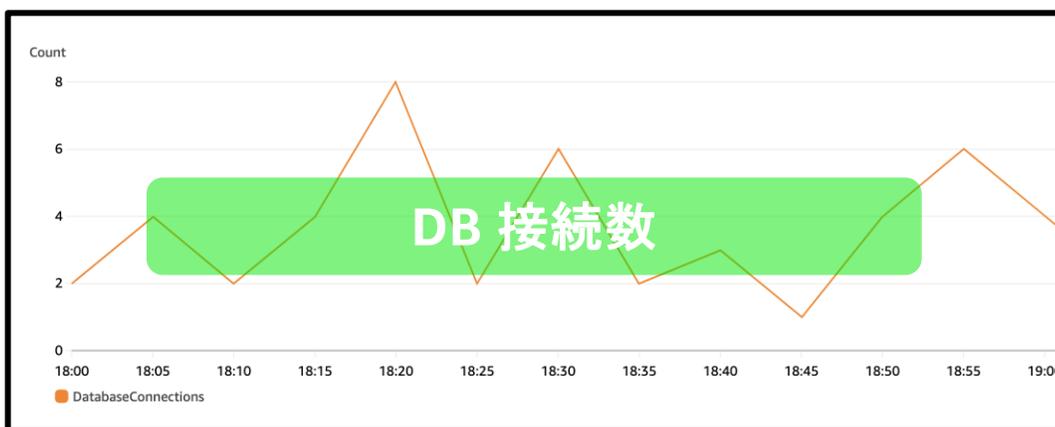
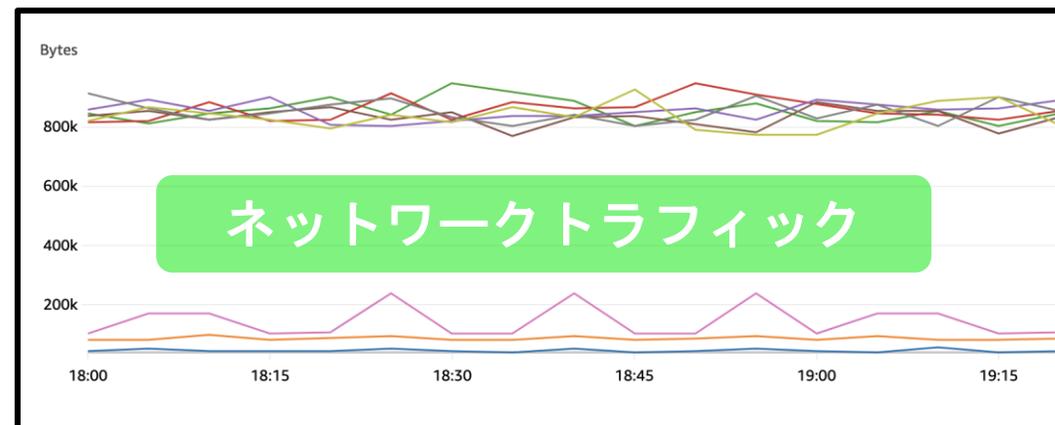
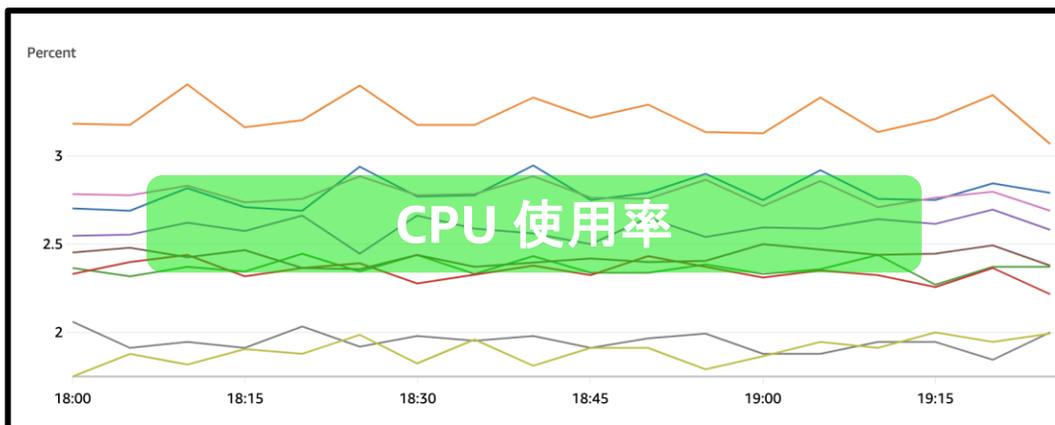
問題はあるがユーザーは検知できていない状態

EC サイトのアーキテクチャ



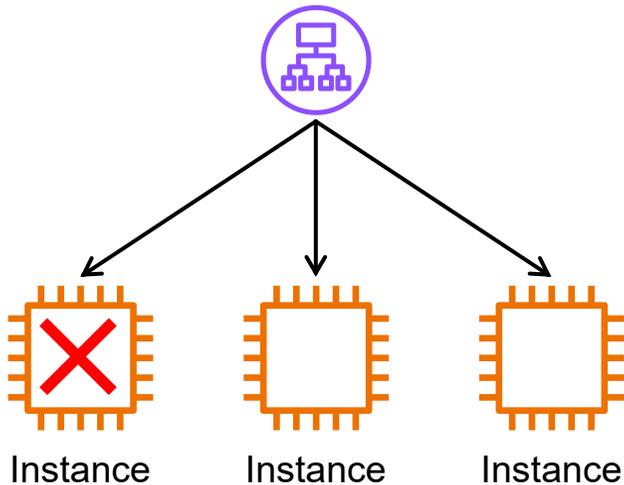
システムメトリクスは正常

監視メトリクスに問題はなく、システムは正常に動いているように見える

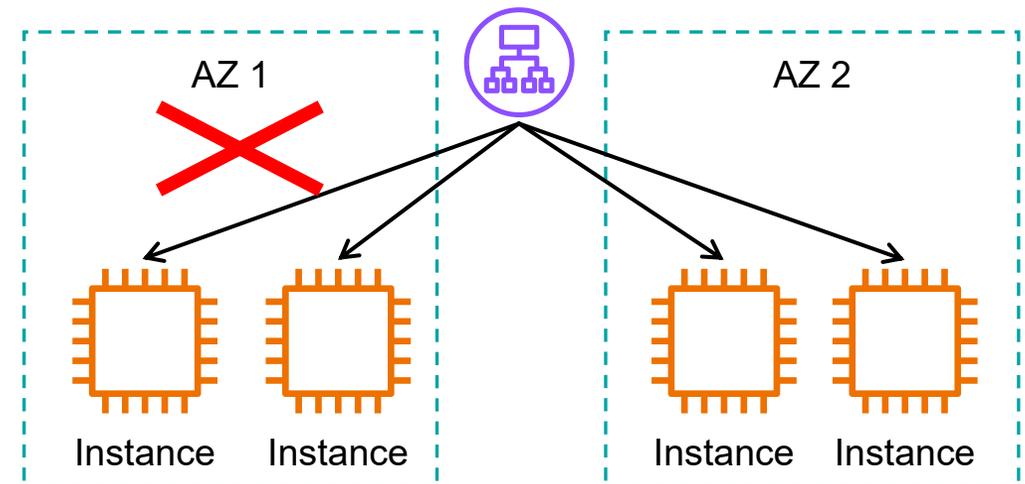


グレー障害の代表的なパターン

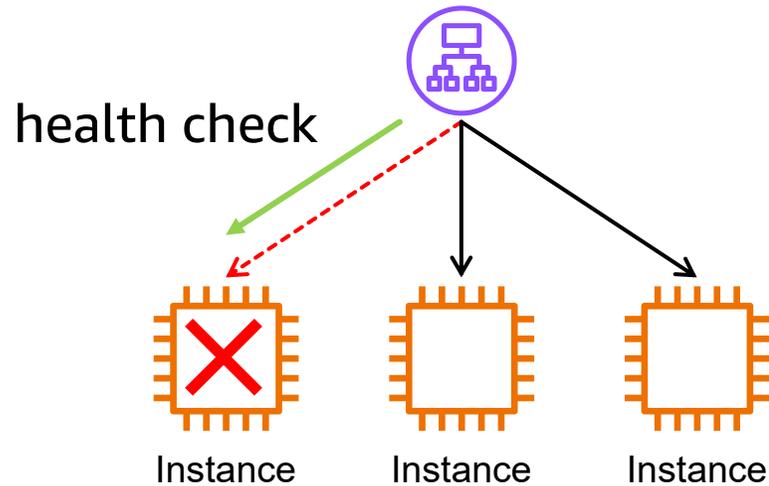
特定インスタンスでの部分的な障害



AZ レベルでの品質低下

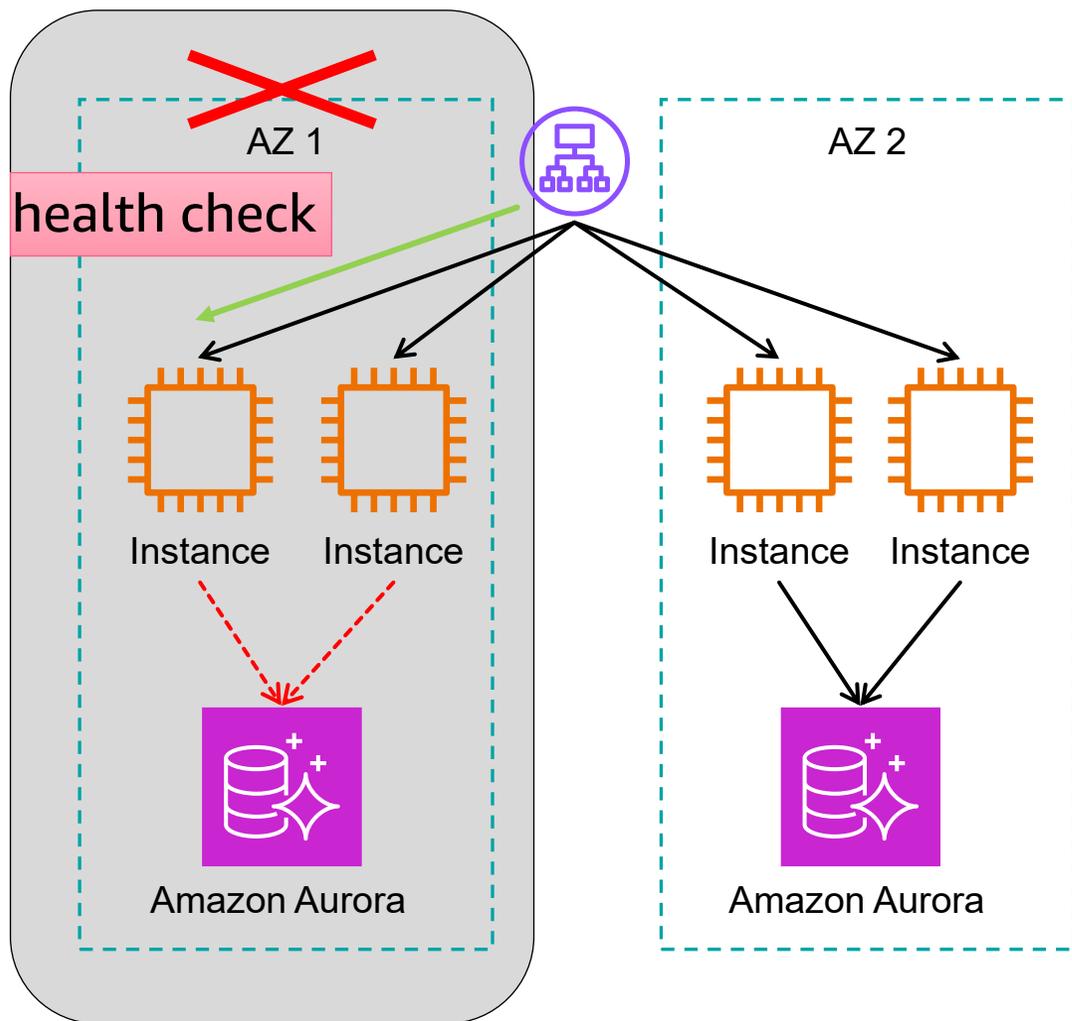


例 1：特定インスタンスで発生する部分的な障害



- Auto Scaling Group 内の複数インスタンスのうち、1 台で断続的な遅延が発生
- 他のインスタンスは正常に動作している
- 問題のあるインスタンスでも **ALB のヘルスチェックは通過**している
- 失敗しているのは一部のリクエストのみ
 - 商品一覧表示：○
 - カート追加：×
 - 注文処理：○

例 2: AZ レベルでの品質低下



- 特定 AZ 内のすべてのインスタンスで問題が発生している
- 他の AZ では正常に動作している
 - ユーザーはリクエストの半数が遅延やエラー
- 特定の AZ で発生する EC2 のネットワーク品質低下により、その AZ で稼働するインスタンスからデータベースへの接続が断続的に遅延
- **ALB のヘルスチェックは通過**

ヘルスチェックでは検出できないのか？

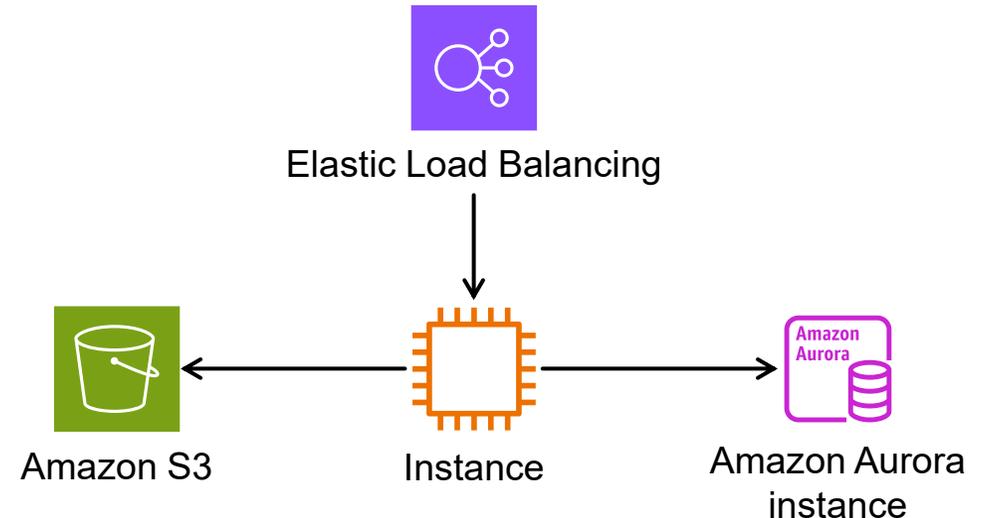
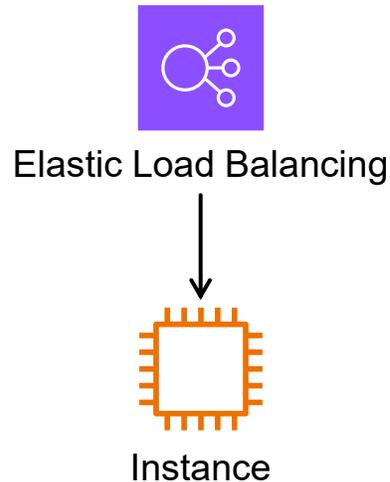
シャローヘルスチェックとディープヘルスチェック

シャロー（浅い）ヘルスチェック

- サーバーとアプリケーションの状態を確認
- TCP 接続確認、HTTP 200 応答確認、EC2 ステータスチェック
- 外部依存関係はチェックしない

ディープ（深い）ヘルスチェック

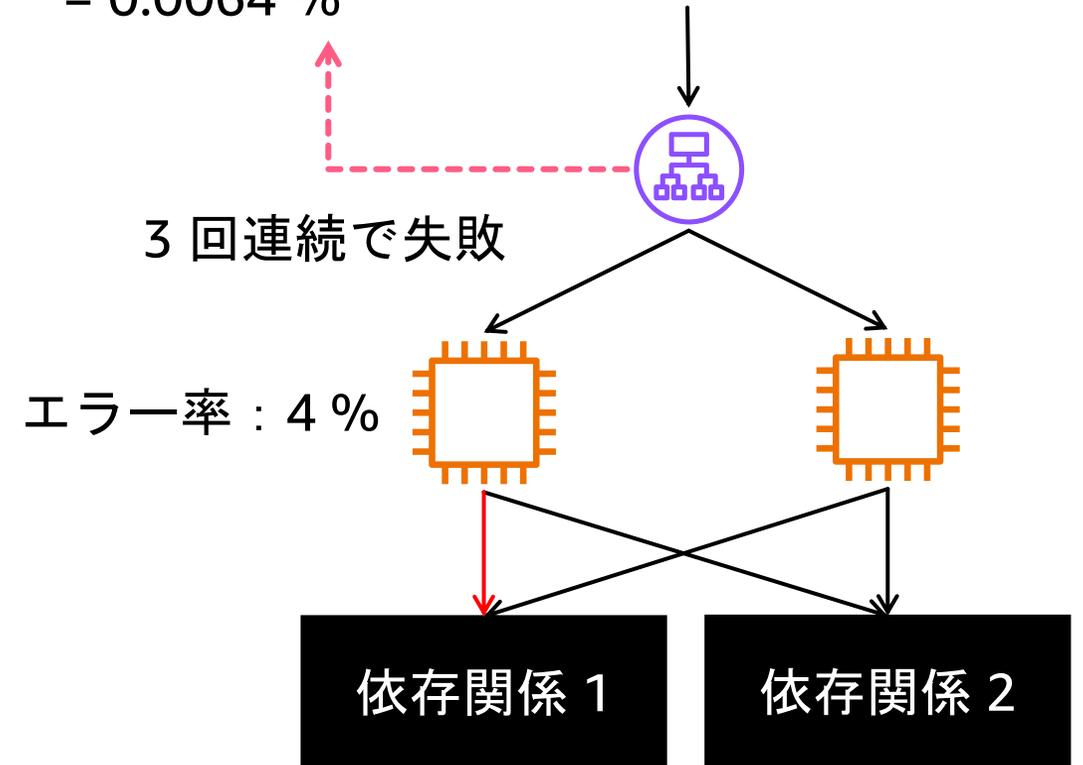
- サーバーの状態に加え、外部依存関係も確認
- DB へのクエリ実行や S3 アクセスなども含む
- エンドツーエンドでの機能確認が可能
- 依存関係の問題について過剰に検知



ヘルスチェックで検出できない小さなエラー

- あるサーバーでのエラーレート 4 %
- 3 回連続でヘルスチェックが失敗したらアラームを発行
- 3 回連続で失敗する確率は 0.0064 %

$$4 \% \times 4 \% \times 4 \% = 0.0064 \%$$



リクエストの 4 % に問題が起きているが
検出できず、サービス品質低下

グレー障害の検出



見えづらい問題をどう検出するか

- システム全体では正常に見える中、部分的に発生している異常を見つけ出す
- ヘルスチェックのような基本的な監視では、問題を捉えきれない



Amazon CloudWatch を活用したグレー障害の検出



Amazon CloudWatch

- 1 Contributor Insights
- 2 複合アラーム

ディメンションを活用した柔軟な監視

- ディメンションは、メトリクスを分類するためのラベル（InstanceID 等）
- ディメンションには、障害分離境界（AZ-ID 等）やビジネスコンテキスト（CustomerID 等）を含めることができる
- EMF（Embedded Metric Format）を使用することでログデータに直接メトリクスを埋め込むことが可能

インスタンス名 1...	InstanceID	メトリクス名	アラーム
webservers	i-09371c03af0776...	StatusCheckFailed_Instance ⓘ	アラームなし
webservers	i-09371c03af0776...	CPUCreditBalance ⓘ	アラームなし
webservers	i-09371c03af0776...	EBSReadOps ⓘ	アラームなし
webservers	i-09371c03af0776...	NetworkOut ⓘ	アラームなし
webservers	i-09371c03af0776...	NetworkIn ⓘ	アラームなし
webservers	i-09371c03af0776...	NetworkPacketsOut ⓘ	アラームなし

EMF のディメンション設定例

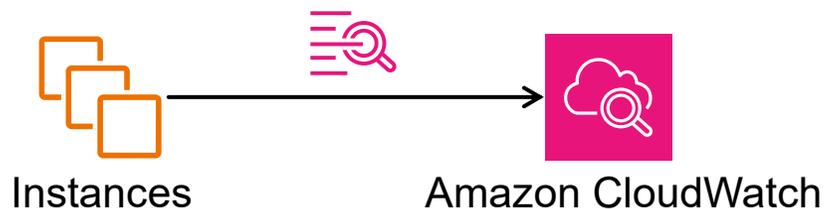
```
"Dimensions": [  
  ["AZ-ID", "CustomerId", "Region"],  
  ["AZ-ID", "Region"]  
]
```

高カーディナリティな（値の種類が非常に多い）ディメンションをどう分析するか？

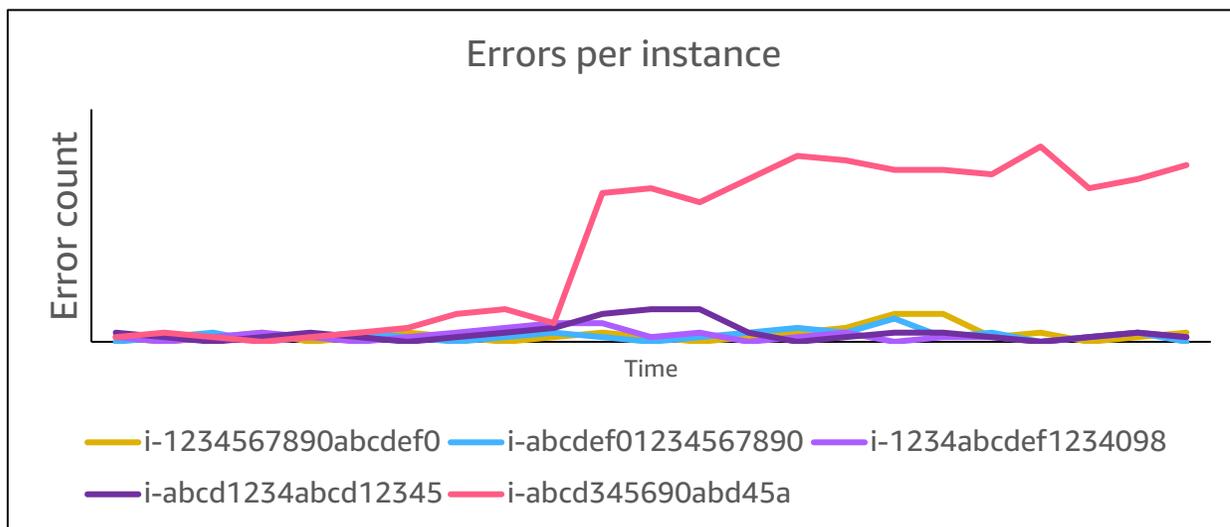
Contributor Insights の活用

Contributor Insights を使用することで、膨大な数の要素から重要な情報を抽出

高カーディナリティデータを分析する



- CloudWatch Logs の時系列データに対する上位 N のコントリビューターの分析を簡素化
- システムとアプリケーションのパフォーマンスに影響を与えている人またはものをリアルタイムに確認
- グラフ化やアラームの設定も可能

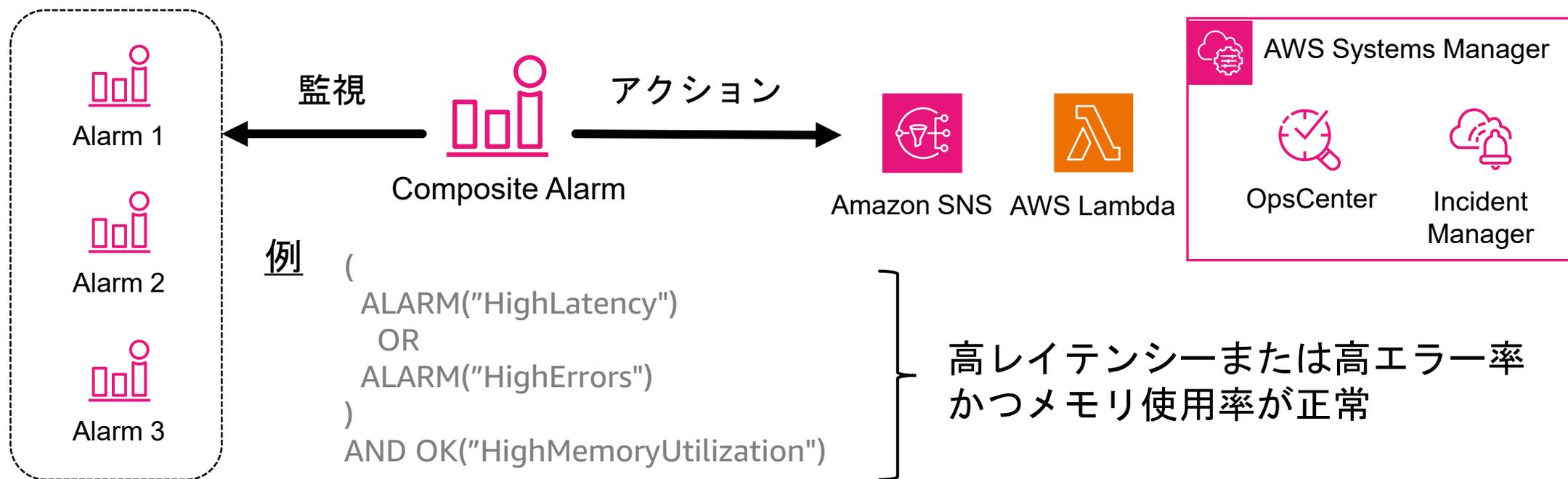


例 :

- エラーを最も多く返す API / インスタンス
- 特定のステータスコードが頻出する URL

複合アラーム

- 複数のアラームの状態を組み合わせて監視
- 特定の条件の組み合わせが満たされた場合にのみアクションを実行
- AND / OR / NOT というシンプルなアラームロジックで設定



複合アラームを応用して特定 AZ の問題を検出

**AZ 1 内のサービス障害を
検知するアラーム**

az1-availability
ALARM("az1-home-availability") **OR**
ALARM("az1-listproducts-availability")

**AZ 1 で高レイテンシーを
検知するアラーム**

az1-latency
ALARM("az1-home-latency") **OR**
ALARM("az1-listproducts-latency")

**AZ 1 だけの問題か
チェックするアラーム**

use1-az1-impact
(ALARM("az1-availability") **OR** **ALARM**("az1-latency")) **AND NOT**
(ALARM("az2-availability") **OR** **ALARM**("az2-latency") **OR**
ALARM("az3-availability") **OR** **ALARM**("az3-latency"))

* : Contributor Insights ルール

**AZ 1 の複数インスタンスでの
エラー発生*を検知するアラーム**

not-single-instance-use1-az1
INSIGHT_RULE_METRIC("5xx-errors-use1-
az1", "UniqueContributors") >= 2

**AZ 1 の障害を検知する
アラーム**

use1-az1-isolated-impact
ALARM("use1-az1-impact") **AND**
ALARM("not-single-instance-use1-az1")

AZ 1 の障害アラーム

グレー障害のAWSでの対応例



グレー障害対応の選択肢

状況に応じて適切な手段を選択する

	インスタンス単位	アプリケーション	ロケーション切り替え
特徴	<ul style="list-style-type: none">• 比較的シンプルな実施手順• 効果は該当リソースに限定	<ul style="list-style-type: none">• 事前の実装と設計が必要• より柔軟な対応が可能	<ul style="list-style-type: none">• 物理的な問題から完全な退避が可能• 広範な問題に有効
例	<ul style="list-style-type: none">• インスタンスの再起動• 新規インスタンスへの置き換え	<ul style="list-style-type: none">• グレースフルデグラデーション• サーキットブレーカー• エクスポネンシャルバックオフによるリトライ 等	<ul style="list-style-type: none">• AZ のシフト• リージョンフェイルオーバー

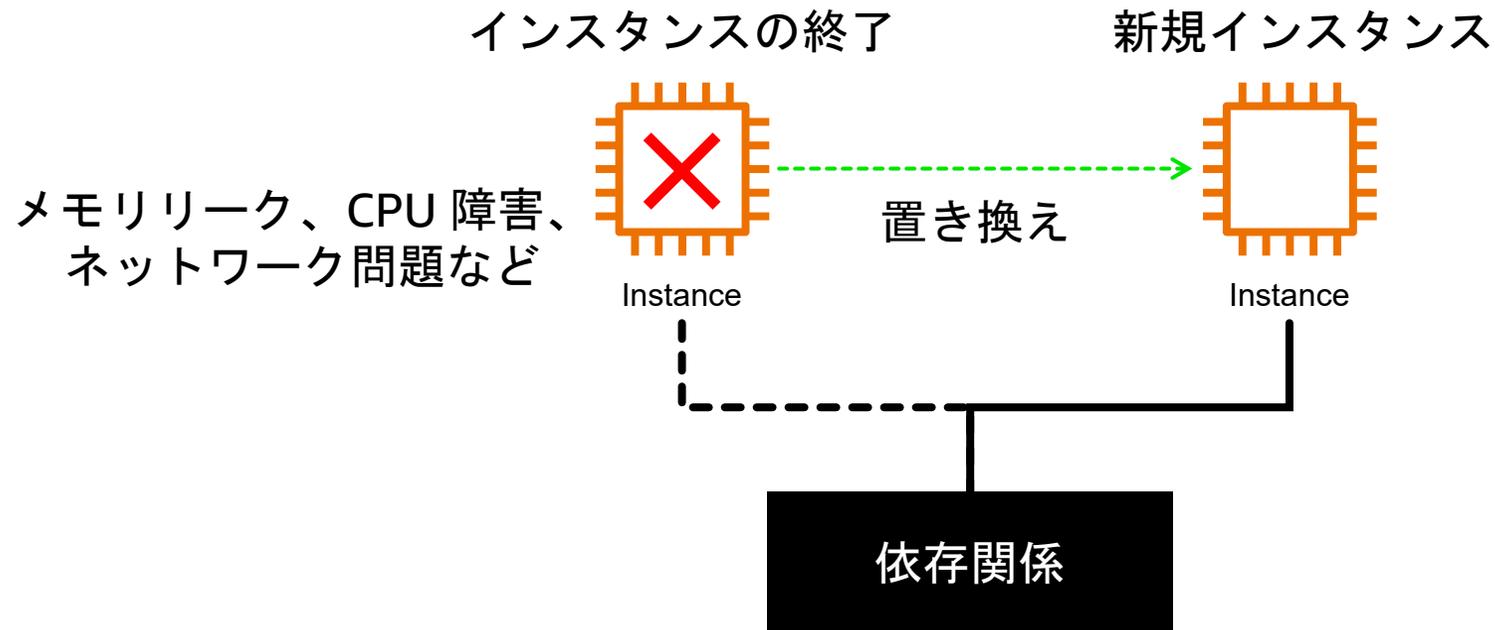
グレー障害対応の選択肢

状況に応じて適切な手段を選択する

	インスタンス単位	アプリケーション	ロケーション切り替え
特徴	<ul style="list-style-type: none">比較的シンプルな実施手順効果は該当リソースに限定	<ul style="list-style-type: none">事前の実装と設計が必要より柔軟な対応が可能	<ul style="list-style-type: none">物理的な問題から完全な退避が可能広範な問題に有効
例	<ul style="list-style-type: none">インスタンスの再起動新規インスタンスへの置き換え	<ul style="list-style-type: none">グレースフルデグラデーションサーキットブレーカーエクスポネンシャルバックオフによるリトライ 等	<ul style="list-style-type: none">AZ のシフトリージョンフェイルオーバー

インスタンス単位の対応

- 新しいインスタンスに置き換えることで、解消する可能性もある
- AZ全体に影響が出ているような場合には適していない
- 新しいインスタンス起動と初期化に時間がかかる



グレー障害対応の選択肢

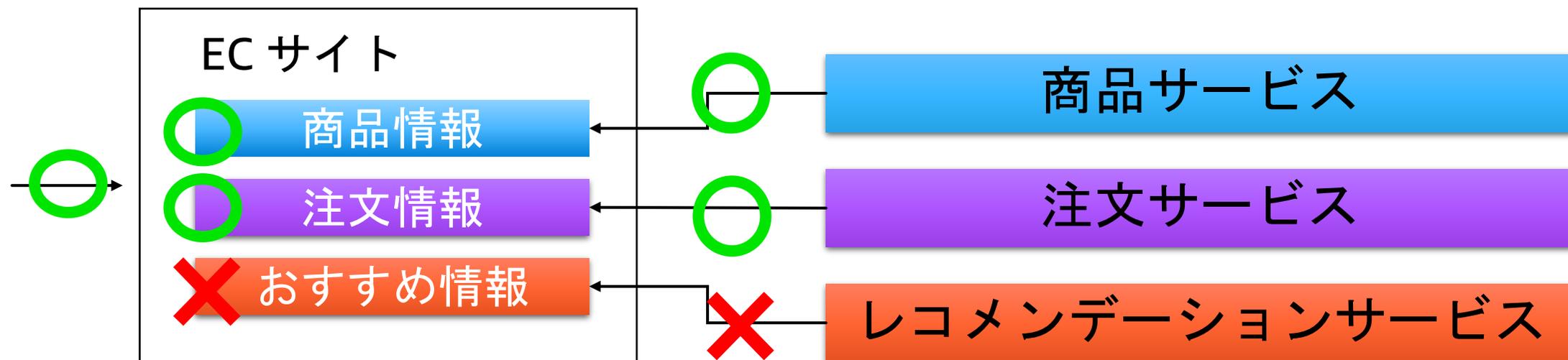
状況に応じて適切な手段を選択する

	インスタンス単位	アプリケーション	ロケーション切り替え
特徴	<ul style="list-style-type: none">• 比較的シンプルな実施手順• 効果は該当リソースに限定	<ul style="list-style-type: none">• 事前の実装と設計が必要• より柔軟な対応が可能	<ul style="list-style-type: none">• 物理的な問題から完全な退避が可能• 広範な問題に有効
例	<ul style="list-style-type: none">• インスタンスの再起動• 新規インスタンスへの置き換え	<ul style="list-style-type: none">• グレイスフル デグラデーション• サーキットブレーカー• エクスポンネンシャル バックオフによるリトライ 等	<ul style="list-style-type: none">• AZ のシフト• リージョンフェイルオーバー

グレイスフルデグラデーション

Graceful Degradation 優美（上品）な劣化

一部のコンポーネントが停止しても、システム全体は停止せず
機能制限（デグレート）しながらサービスを提供し続ける能力



一部の情報が表示されなくても顧客は
引き続き ECサイトの主要機能を利用可能

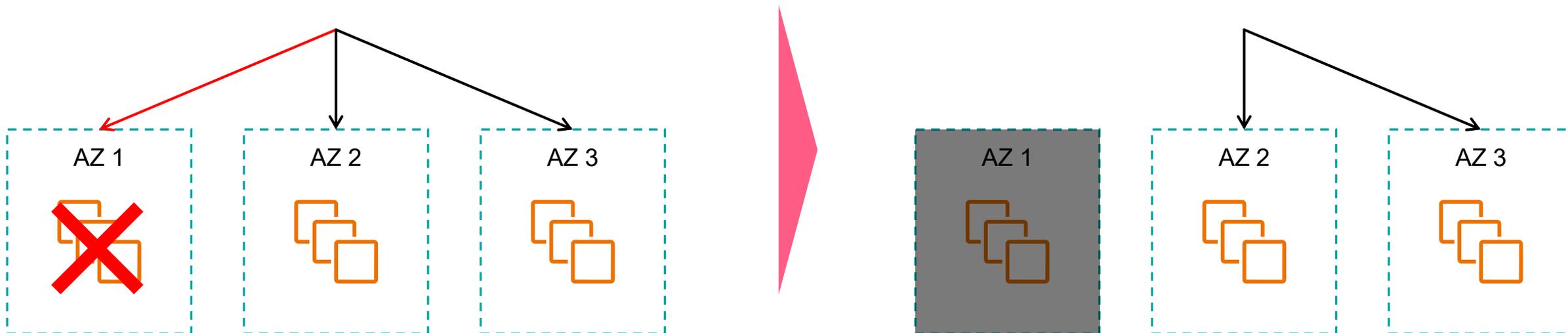
グレー障害対応の選択肢

状況に応じて適切な手段を選択する

	インスタンス単位	アプリケーション	ロケーション切り替え
特徴	<ul style="list-style-type: none">• 比較的シンプルな実施手順• 効果は該当リソースに限定	<ul style="list-style-type: none">• 事前の実装と設計が必要• より柔軟な対応が可能	<ul style="list-style-type: none">• 物理的な問題から完全な退避が可能• 広範な問題に有効
例	<ul style="list-style-type: none">• インスタンスの再起動• 新規インスタンスへの置き換え	<ul style="list-style-type: none">• グレースフルデグラデーション• サーキットブレーカー• エクスポネンシャルバックオフによるリトライ 等	<ul style="list-style-type: none">• AZ のシフト• リージョンフェイルオーバー

アベイラビリティゾーンの切り離し

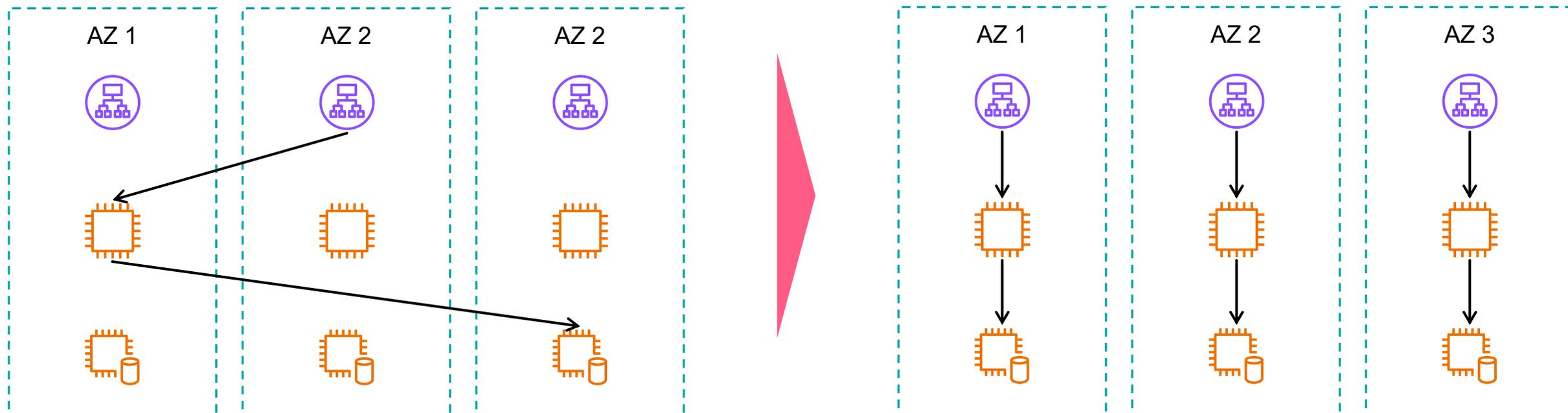
- 問題のある AZ を一時的に切り離すことで、システム全体の健全性を保つ
- API レスポンス低下やコントロールプレーン障害など**広範な問題への対応**が可能
- 障害の原因特定に時間をかけずに、**まず影響を回避**できる



アベイラビリティゾーンの独立性

Availability Zone independence (AZI)

- 異なる AZ 間の不要な通信を防ぎ、AZ 内でリソースの利用を完結させる
- AZ 間のデータ転送による遅延とコストを最小限に抑える
- AZ 障害等の際、影響を受けている AZ からの退避を可能にする

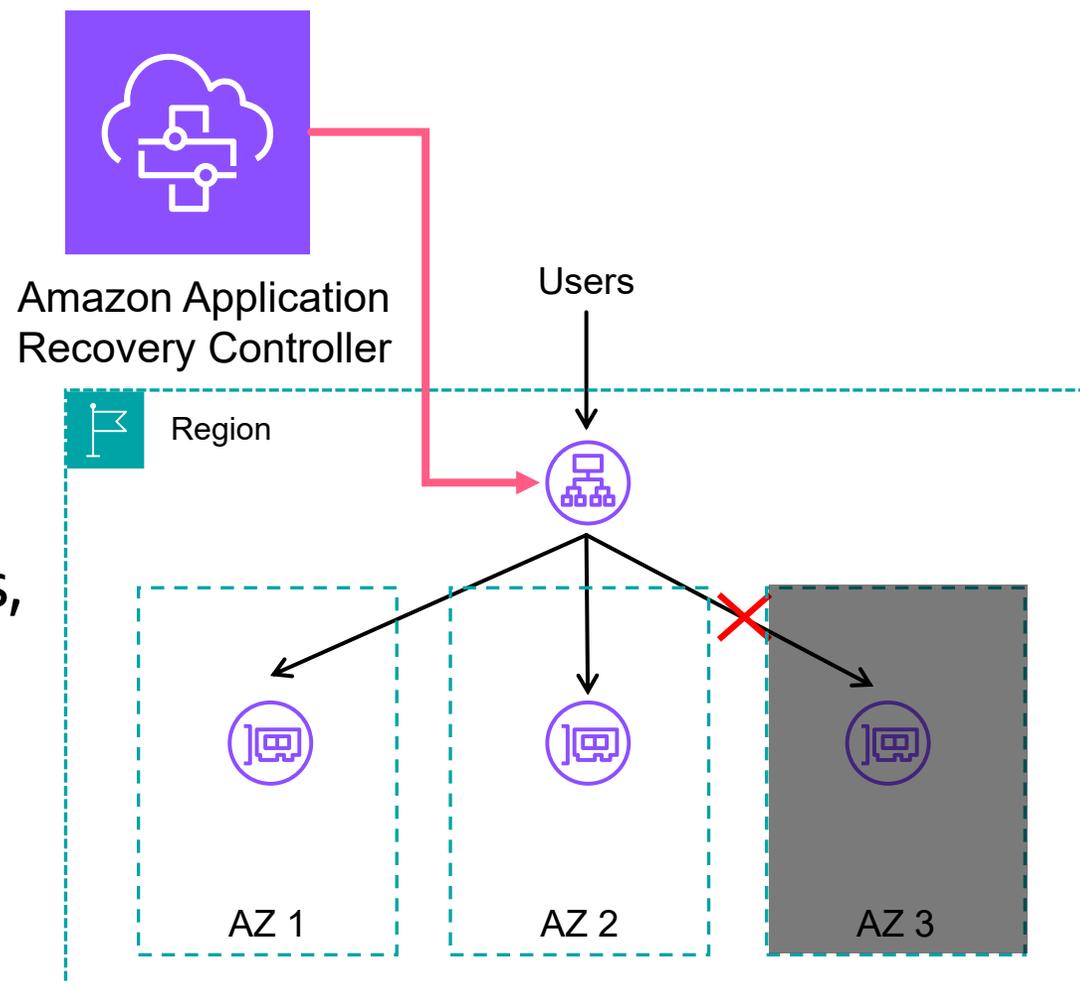


どうやって問題のある AZ
を切り離すか？

Amazon Application Recovery Controller (ARC)

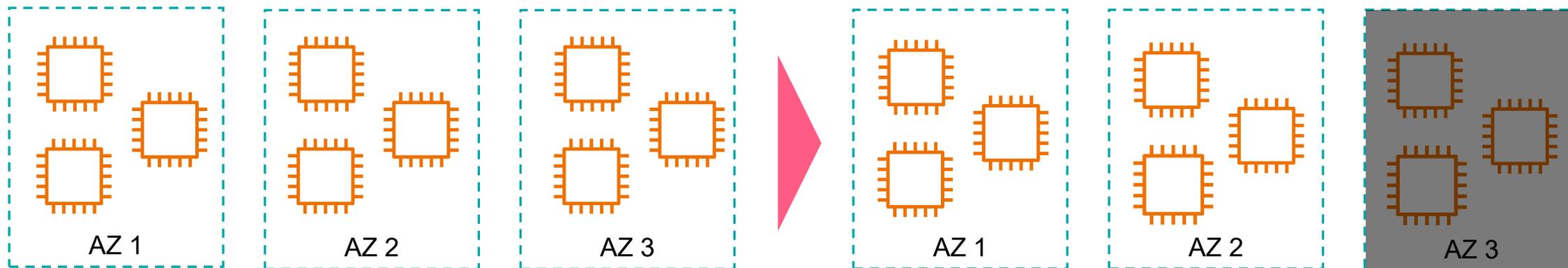
ゾーンシフト

- リージョン内のリソースに対して、**特定 AZ へのトラフィックルーティングを回避**
- トラフィックを遠ざける AZ とサポートされているリソースを選択する
- 指定有効期間の後、自動的に復帰
- ALB/NLB, Amazon EC2 Auto Scaling groups, Amazon EKS をサポート



ゾーンシフトにおけるリソース容量の考慮

- 1つのAZで障害が発生しても残りのAZで負荷を処理できるよう、**事前に十分な容量をプロビジョニング**する
- 3つのAZを使用する場合、各AZは通常負荷の66%で動作するようにオーバープロビジョニング
 - アプリケーションに6インスタンスが必要な場合、各AZに3インスタンス（計9インスタンス）を配置する
 - 1つのAZがダウンしても、残りのAZで合計6インスタンスを維持可能



ゾーンシフト実行の判断基準

ゾーンシフトの実行タイミングを判断するには、明確な SLO/SLI の設定が不可欠

サービスレベル目標 (SLO)

- サービスの信頼性に関する目標値
- 顧客に提供するサービス品質の目標を定量的に示したもの

例：

- 注文処理の成功率を 99.9 % 以上に保つ
- ページ読み込み時間を 2 秒以内に抑える
- 買い物かご機能の可用性を 99.99 % 維持する
- 決済処理のエラー率を 0.1 % 未満に抑える

サービスレベル指標 (SLI)

- SLO の達成度を測定するための具体的な指標
- 実際に測定可能な定量的なメトリクス

例：

- 注文 API の HTTP 200 レスポンス率
- トップページでの Time to First Byte (TTFB)
- 買い物かご API のエラー率
- 決済トランザクションの成功率

例：AZ 障害が疑われ、SLO 違反が 15 分以上継続する場合、ゾーンシフトを開始

ゾーンシフトは手動で実行できる

The screenshot shows the 'ロードバランサー属性を編集' (Edit Load Balancer Attributes) page for a load balancer named 'test'. The 'アベイラビリティゾーンのルーティング設定' (Availability Zones Routing Settings) section is expanded. It includes options for 'クライアントルーティングポリシー (DNS レコード)' (Client Routing Policy (DNS Record)), 'DNS 解決' (DNS Resolution) with 'Route 53 Resolver' selected, and 'ロードバランサーのターゲット選択ポリシー' (Load Balancer Target Selection Policy) with 'クロスゾーン負荷分散を無効にする - デフォルト' (Disable cross-zone load balancing - default) selected. A diagram illustrates traffic distribution across three Availability Zones (AZ1, AZ2, AZ3) via a Network Load Balancer (NLB). A legend indicates that orange lines represent traffic and grey boxes represent Availability Zones.

The screenshot shows the 'ゾーンシフトを開始' (Start Zone Shift) page in the Amazon Application Recovery Controller console. It includes a 'ゾーンシフトのオプションを選択' (Select Zone Shift Options) section with a dropdown for 'アベイラビリティゾーンを選択' (Select Availability Zone) set to 'apne1-az1 (ap-northeast-1c)'. Below, the 'apne1-az1 (ap-northeast-1c) のリソース (1)' (Resources (1)) section shows a table with one resource: 'test' (ARN: arn:aws:arc:ap-northeast-1:1b774b404aae:nlb). The 'ゾーンシフトの有効期限を設定' (Set Zone Shift Validity Period) section has '12 時間' (12 hours) selected. A 'コメント' (Comment) field contains the text 'テスト'. At the bottom, there is a '同意' (I agree) checkbox and 'キャンセル' (Cancel) / '開始' (Start) buttons.

ゾーンシフトを有効にしておけば、手動でゾーンシフトを実行可能
APIで自動実行も可能



ゾーンオートシフト

AZ に潜在的な影響がある場合、**AWS が自動的にトラフィックを別の AZ に移動**

設定するリソース 情報

ゾーンオートシフト用に設定するリソースを選択します。AWS は練習実行中、リソースのためにトラフィックを移動させます。リソースのためにゾーンオートシフトを有効にすることもできます。

ゾーンオートシフトのために設定するリソースを選択

test

ゾーンオートシフトのステータス

ゾーンオートシフトを有効にすると、AWS は、アベイラビリティゾーンで潜在的な問題が発生した場合、または練習実行中に、リソースについてのトラフィックを AZ から自動的にシフトします。練習実行は大体 1 週間に 1 回の頻度で行われます。ゾーンオートシフトを無効にすると、AWS は、練習実行のためにのみトラフィックをアベイラビリティゾーンからシフトします。この設定は後で変更できません。

有効化

AWS は、アベイラビリティゾーンの潜在的な問題が検出され、練習実行用に自動的にトラフィックを移動します。

無効化

AWS は練習実行用のトラフィックのみを移動します。

i アプリケーションが各アベイラビリティゾーンで十分なキャパシティを備えており、オートシフトおよび練習実行中に 1 つのアベイラビリティゾーンのキャパシティが削除されても、正常に動作し続けるようにしてください。

詳細はこちら [🔗](#)

- AWS の内部モニタリングツールを使用して早期に問題を検知
- 電源障害やネットワーク障害などの**潜在的な問題に自動対応**
- 明示的な有効化が必要
- 毎週 30 分間の練習実行が走る
(練習実行をブロックする時間帯を定義することも可能)
- Amazon EventBridge を使用してモニタリングも可能

今日持って帰っていただきたいこと（再掲）

- グレー障害とは何かを理解する
- 耐障害性をあげる方法について理解する
- カオスエンジニアリングと、訓練の方法について理解する

カオスエンジニアリングと 障害の訓練について

カオスエンジニアリングとは

Netflix 社がクラウド上で分散したアプリケーションを構築していく中でスタートした、本番環境に意図的に障害を注入し、その影響を観察する実験を通して、システムの可用性、信頼性を高めていく取り組み

<https://principlesofchaos.org/ja/>

「カオスエンジニアリングの原則」では、システムが本番環境における不安定な状態に耐える能力へ自信を持つためにシステム上で実験を行う訓練方法 と定義

カオスエンジニアリングの詳細原則

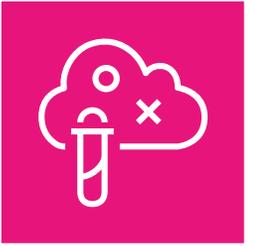
- 定常状態における振る舞いの仮説を立てる
- 実世界の事象は多様である
- 本番環境で検証を実行する
- 継続的に実行する検証の自動化
- 影響範囲を局所化する

<https://principlesofchaos.org/ja/>

カオスエンジニアリングの詳細原則

- 定常状態における振る舞いの仮説を立てる
- 実世界の事象は多様である
- 本番環境で検証を実行する
- 継続的に実行する検証の自動化
- 影響範囲を局所化する

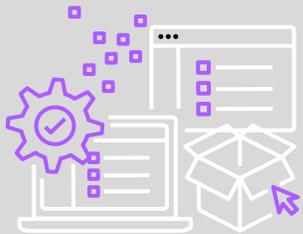
これらを簡単に実現できないか？



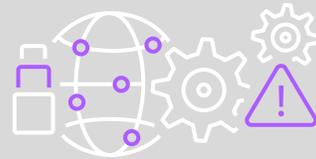
AWS Fault Injection Service (FIS)

障害注入実験を実行するためのフルマネージド サービス

簡単に
始められる



実環境の
条件で

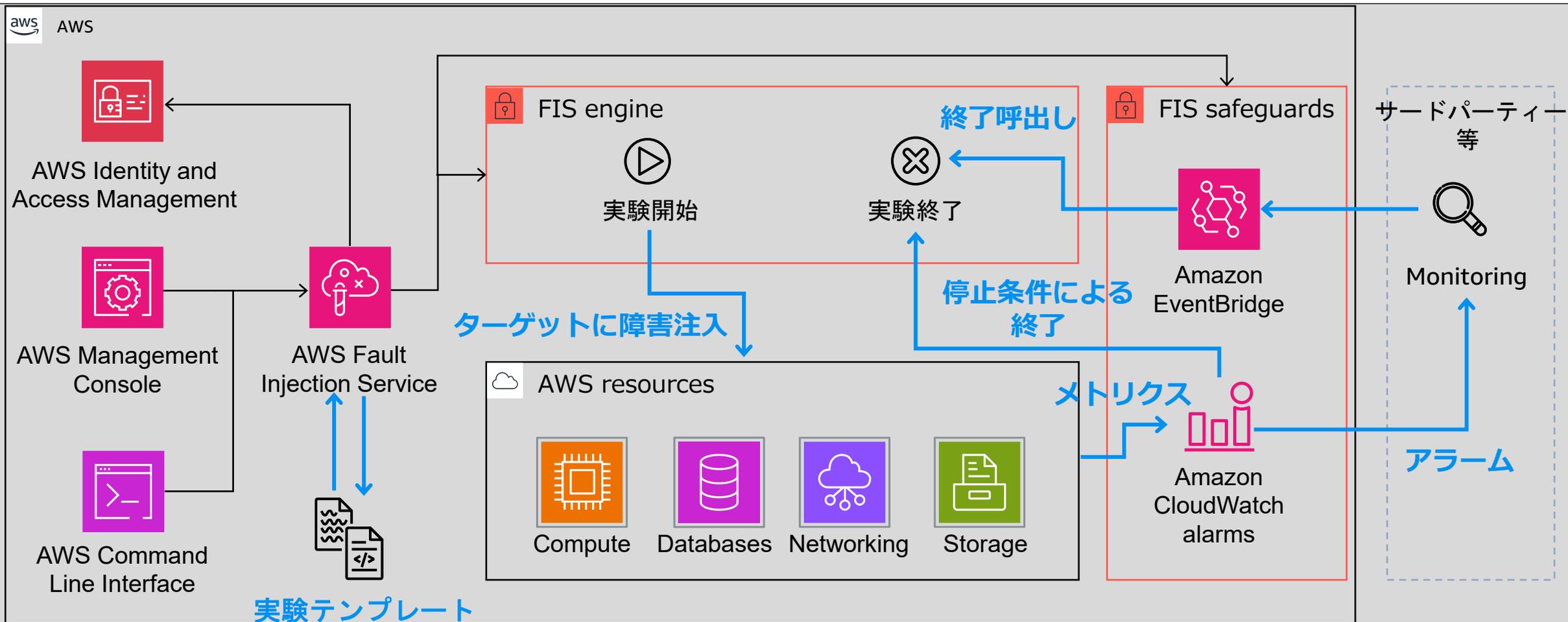


セーフガード



※ AWS Resilience Hub の1つの機能

AWS FIS を使った実験の流れ



AWS FIS で AZ 障害をシミュレートして訓練する

- お客様がアプリケーションの耐障害性をテストするために適用できるイベントを定義したシナリオライブラリ
- **AZの可用性: 電源の中断** のシナリオを使用して、AZ の電力が中断される状況をシミュレート

シナリオライブラリ (12) 情報 シナリオ使用してテンプレートを作成

シナリオを選択すると、その詳細が表示されます。

Q リソースタイプ、サービス、アクションで検索します

< 1 > ⚙

AZ の可用性: 電源の中断

1つのAZの複数のリソースタイプに影響を与え、タグと明示的なARNによってターゲットリングし、1つのAZの停電を近似します。

EC2 ストレス: CPU

増加するCPU負荷を1つ以上のインスタンスに注入します(インスタンスタグに基づいてターゲットを設定します)

EC2 ストレス: インスタンス障害

1つ以上のインスタンスを5分間停止します(インスタンスタグに基づいてターゲットを設定します)。

EC2 ストレス: ディスク

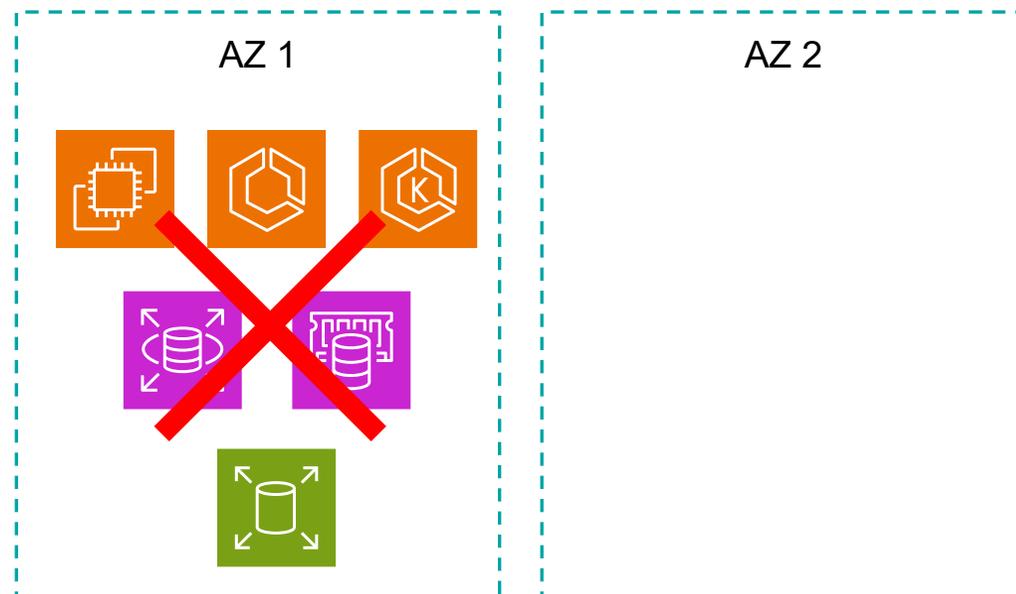
増加するディスク使用率を1つ以上のインスタンスに注入します(インスタンスタグに基づいてターゲットを設定します)

EC2 ストレス: ネットワークレイテンシー

増加するネットワークレイテンシーを1つ以上のインスタンスに注入します(インスタンスタグに基づいてターゲットを設定します)

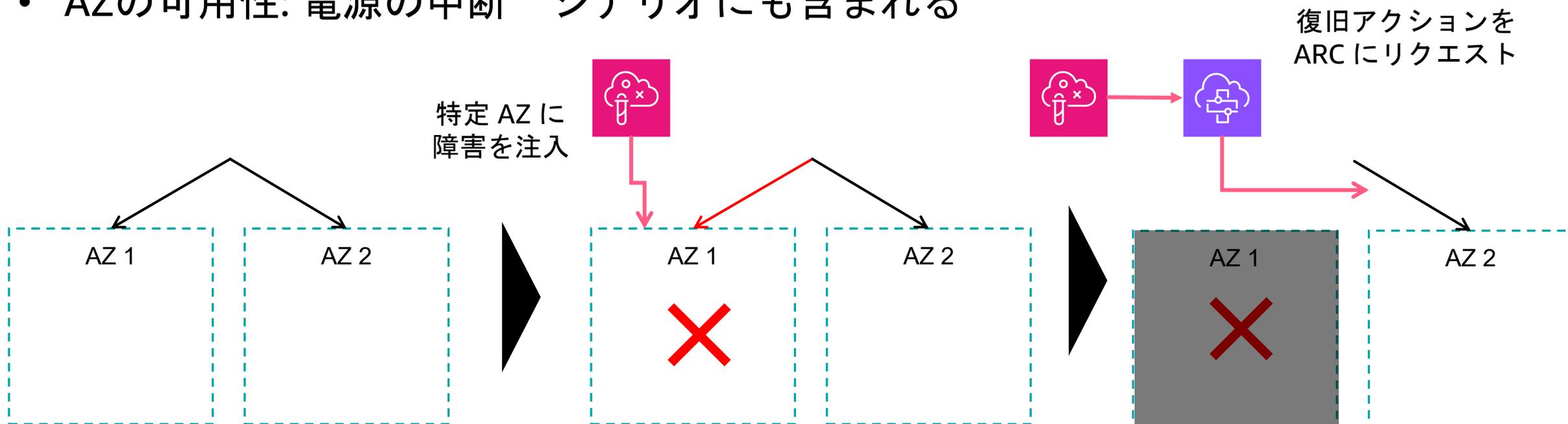
EC2 ストレス: メモリ

増加するメモリ負荷を1つ以上のインスタンスに注入します(インスタンスタグに基づいてターゲットを設定します)



AWS FIS の復旧アクション

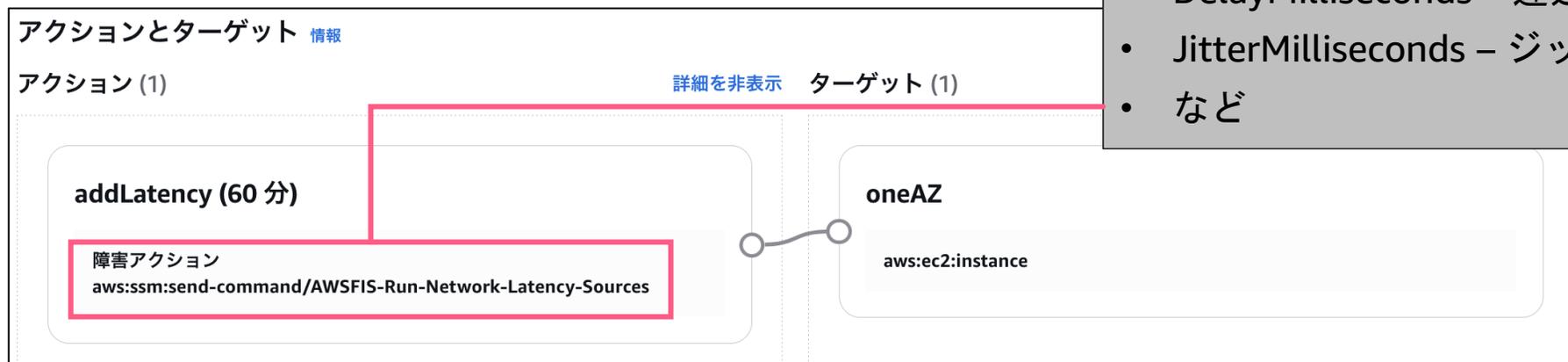
- `aws:arc:start-zonal-autoshift` という復旧アクションが提供されている
- 特定の AZ からのトラフィックシフトをシミュレートすることが可能
- FIS が ARC の API を使用してリクエストを送りゾーンシフトが開始
- AZの可用性: 電源の中断 シナリオにも含まれる



SSM ドキュメントでグレー障害シミュレート

AWS FIS の SSM ドキュメントアクションを使用してグレー障害をシミュレート

- SSM Agent がインストールされたインスタンスで実行可能
- グレー障害をシミュレートできる SSM ドキュメントアクションの例：
 - **AWSFIS-Run-Network-Latency-Sources :**
特定ソースとの通信にネットワークレイテンシーとジッターを追加
 - **AWSFIS-Run-Network-Packet-Loss-Sources :**
特定ソースとの通信にパケットロスを発生



ドキュメントパラメータ

- DelayMilliseconds – 遅延 (ミリ秒)
- JitterMilliseconds – ジッター (ミリ秒)
- など

組織全体で取り組むレジリエンス文化の醸成

AWS Well-Architected フレームワーク：信頼性の柱



- 障害訓練を特別なイベントではなく **日常的な活動として組織に定着させる**
- 技術チームだけでなく、**実際の障害時に判断を求められる関係者**も含める
- 訓練で見つかった課題をプロセス改善につなげ、次の訓練で検証

まとめ

- グレー障害とは何かを理解する
- 耐障害性をあげる方法について理解する
- カオスエンジニアリングと、訓練の方法について理解する

参考資料

- AWSにおけるグレー障害の検出と対策
<https://www.youtube.com/watch?v=Nn07ZgWqI8I>
- ミッションクリティカルシステムを AWS に載せるには?
https://pages.awscloud.com/rs/112-TZM-766/images/AWS-43_AWS-Summit-2023_Resilience.pdf
- **第三十七回 ちょっぴり DD - AWS Fault Injection Service で AZ 障害を体験しよう - データベースの可用性向上対策の検証**
<https://www.youtube.com/watch?v=9KYuib6NNFM>
- **カオスエンジニアリング—回復力のあるシステムの実践**
<https://www.oreilly.co.jp//books/9784873119885/>

Thank you!

Yukihiro Kikuchi

kyukihi@amazon.co.jp

