



AWS ホワイトペーパー

AWS でのリアルタイム通信



AWS でのリアルタイム通信: AWS ホワイトペーパー

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

アマゾン の商標およびトレードドレスはアマゾン 以外の製品およびサービスに使用することはできません。また、お客様に誤解を与える可能性がある形式で、または アマゾン の信用を損なう形式で使用することもできません。Amazon が所有していない他のすべての商標は、それぞれの所有者の所有物であり、Amazon と提携、接続、または後援されている場合とされていない場合があります。

Table of Contents

要約	1
要約	1
Well-Architected の実現状況の確認	1
序章	2
RTC アーキテクチャの基本コンポーネント	3
ソフトスイッチ/PBX	4
セッションボーダーコントローラー (SBC)	4
PSTN 接続	4
PSTN ゲートウェイ	4
SIP トランク	4
メディアゲートウェイ (トランスコーダー)	5
WebRTC でのプッシュ通知	5
WebRTC および WebRTC ゲートウェイ	6
での高可用性とスケーラビリティ AWS	8
アクティブ/スタンバイステートフルサーバー間の HA のフローティング IP パターン	8
RTC ソリューションの適用性	9
RTC アーキテクチャの適用性	11
Application Load Balancer と Auto Scaling を使用した AWS for WebRTC のロードバランシング Application Load Balancer	11
Network Load Balancer または 製品を使用した SIP の AWS Marketplace 実装	12
クロスリージョン DNS ベースのロードバランシングとフェイルオーバー	13
永続的ストレージによるデータの耐久性と HA	15
Amazon Route 53 AWS Lambda、Amazon EC2 Auto Scaling による動的スケーリング	16
Amazon Kinesis Video Streams で高可用性 WebRTC	16
Amazon Chime Voice Connector での高可用性 SIP トランッキング	17
フィールドのベストプラクティス	18
SIP オーバーレイを作成する	18
詳細モニタリングを実行する	19
ロードバランシングに DNS を使用し、フェイルオーバーにフローティング IPs	20
複数のアベイラビリティーゾーンを使用する	22
トラフィックを 1 つのアベイラビリティーゾーンに保持し、EC2 プレイスメントグループを使用する	23
拡張ネットワーキング EC2 インスタンスタイプを使用する	24
セキュリティに関する考慮事項	25

結論	26
頭字語	27
寄稿者	29
ドキュメントの改訂	30
注意	31
AWS 用語集	32
.....	xxxiii

でのリアルタイム通信 AWS

で高可用性でスケーラブルなリアルタイム通信 (RTC) ワークロードを設計するためのベストプラクティス AWS

公開日: 2022 年 5 月 5 日 ([ドキュメントの改訂](#))

要約

現在、多くの組織は、リアルタイムの音声、メッセージング、マルチメディアワークロードのコストを削減し、スケーラビリティを実現しようとしています。このホワイトペーパーでは、Amazon Web Services () でリアルタイム通信 (RTC) ワークロードを管理するためのベストプラクティスの概要AWSと、これらの要件を満たすためのリファレンスアーキテクチャについて説明します。このホワイトペーパーは、これらのワークロードの高可用性とスケーラビリティを実現する方法について、リアルタイムのコミュニケーションに精通している個人向けのガイドです。

このホワイトペーパーには、で RTC ワークロードを設定する方法を示すリファレンスアーキテクチャと AWS、クラウド向けに最適化しながらエンドユーザーの要件を満たすようにソリューションを最適化するためのベストプラクティスが含まれています。進化パケットコア (EPC) はこのホワイトペーパーの対象外ですが、ここで説明するベストプラクティスは Virtual Network Functions (VNFs) に適用できます。

Well-Architected の実現状況の確認

[AWS Well-Architected フレームワーク](#)は、クラウド内でのシステム構築に伴う意思決定の長所と短所を理解するのに役立ちます。このフレームワークの 6 つの柱により、信頼性、安全性、効率、費用対効果、持続可能性の高いシステムを設計および運用するための、アーキテクチャのベストプラクティスを確認できます。で無料で利用できるを使用して [AWS マネジメントコンソール](#) (サインインが必要) [AWS Well-Architected Tool](#)、柱ごとに一連の質問に答えることで、これらのベストプラクティスに照らしてワークロードを確認できます。

クラウドアーキテクチャに関する専門的なガイダンスやベストプラクティス (リファレンスアーキテクチャのデプロイ、図、ホワイトペーパー) については、[AWS アーキテクチャセンター](#)を参照してください。

序章

音声、ビデオ、メッセージングをチャンネルとして使用する通信アプリケーションは、多くの組織とそのエンドユーザーにとって重要な要件です。これらのリアルタイム通信 (RTC) ワークロードには、特定のレイテンシーと可用性の要件があり、関連する設計のベストプラクティスに従うことで満たすことができます。これまで、RTC ワークロードは、専用のリソースを備えた従来のオンプレミスデータセンターにデプロイされてきました。

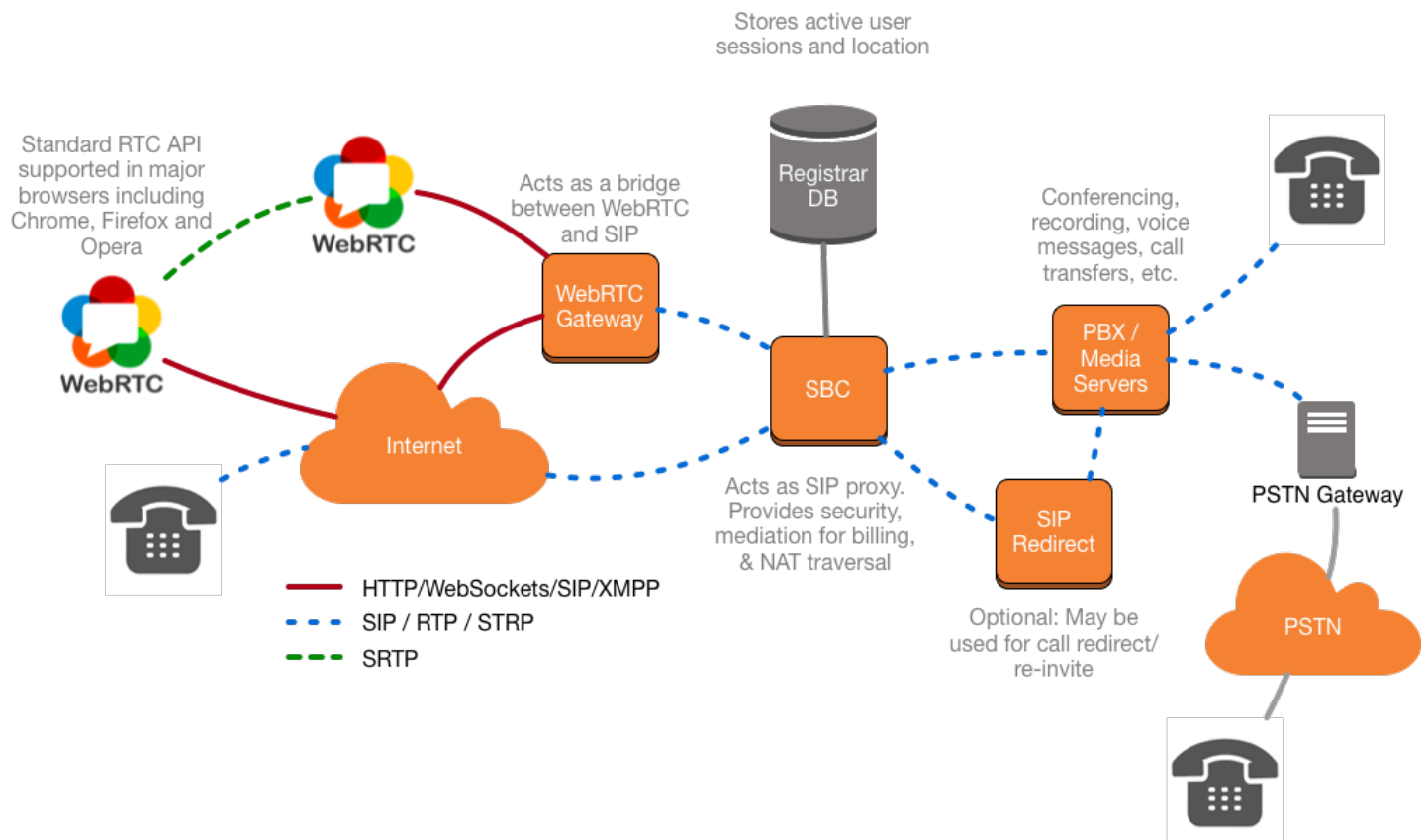
RTC ワークロードには、スケーラビリティ、耐障害性、可用性に優れた環境が必要です。現在、お客様は を使用して、コストを削減 AWS し、俊敏性、伸縮性、市場投入までの時間を短縮して RTC ワークロードを実行します。

RTC アーキテクチャの基本コンポーネント

通信業界では、RTC は通常、レイテンシーを最小限に抑えた 2 つのエンドポイント間のライブメディアセッションを指します。これらのセッションは、以下に関連している可能性があります。

- 2 者間の音声セッション (電話システム、モバイル、Voice over IP (VoIP) など)
- インスタントメッセージング (チャットやインスタントリレーチャット (IRC) など)
- ライブビデオセッション (ビデオ会議やテレプレゼンスなど)

上記の各ソリューションには、共通するコンポーネント (認証、認可、アクセスコントロール、トランスコード、バッファリング、リレーなどを提供するコンポーネントなど) と、送信されるメディアのタイプに固有のコンポーネント (ブロードキャストサービス、メッセージングサーバー、キューなど) があります。このセクションでは、次の図に示すように、音声およびビデオベースの RTC システムと関連するすべてのコンポーネントの定義に焦点を当てます。



RTC の基本的なアーキテクチャコンポーネント

ソフトスイッチ/PBX

ソフトスイッチまたは PBX は音声電話システムの頭であり、さまざまなコンポーネントを使用して、企業内外の音声通話を確認、維持、ルーティングするためのインテリジェンスを提供します。エンタープライズのすべてのサブスクリイバーは、コールを受信または発信するためにソフトスイッチに登録する必要があります。ソフトスイッチの重要な機能は、各サブスクリイバーと、音声ネットワーク内の他のコンポーネントを使用してサブスクリイバーに到達する方法を追跡することです。

セッションボーダーコントローラー (SBC)

セッションボーダーコントローラー (SBC) は、音声ネットワークのエッジに配置され、すべての送受信トラフィック (コントロールプレーンとデータプレーンの両方) を追跡します。SBC の主な責任の 1 つは、音声システムを悪意のある使用から保護することです。SBC は、外部接続のためにセッション開始プロトコル (SIP) トランクと相互接続するために使用できます。一部の SBCs には、[CODECs](#) から別の形式に変換するためのトランスコード機能も用意されています。ほとんどの SBCs には、ネットワークアドレス変換 (NAT) トラバーサル機能も用意されており、ファイアウォールされたネットワーク間でも通話が確立されるのに役立ちます。

PSTN 接続

Voice over IP (VoIP) ソリューションは、パブリック交換電話網 (PSTN) ゲートウェイと SIP トランクを使用してレガシー PSTN ネットワークに接続します。

PSTN ゲートウェイ

PSTN ゲートウェイは、CODEC トランスコーディングを使用して、SIP と SS7 間のシグナリングとリアルタイムトランスポートプロトコル (RTP) と時間分割多重化 (TDM) 間のメディアを変換します。PSTN ゲートウェイは常に PSTN ネットワークに近いエッジに配置されます。

SIP トランク

SIP トランクでは、エンタープライズは TDM (SS7 ベース) ネットワークへの呼び出しを終了せず、エンタープライズと通信事業者間のフローは IP 経路で維持されます。ほとんどの SIP トランクは SBCs。企業は、特定の範囲の IP アドレス、ポートを許可するなど、通信から事前定義されたセキュリティルールに同意する必要があります。

メディアゲートウェイ (トランスコーダー)

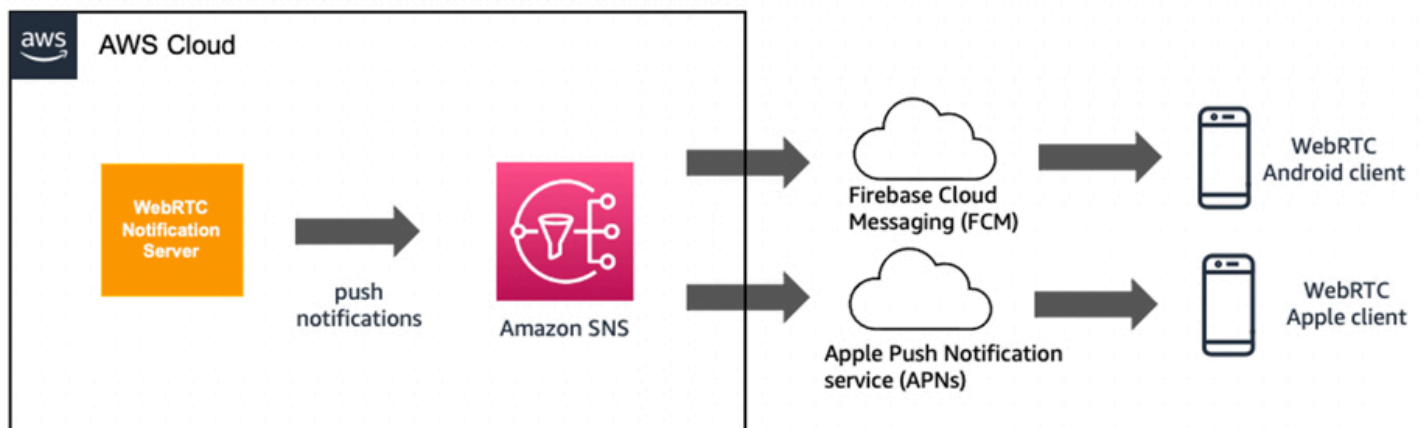
ユーザーは、オーディオやビデオ、オプションデータ、その他の情報を使用してリアルタイムで通信します。通信するには、2つのデバイスがメディアトラックごとに相互に理解されたコーデックについて合意でき、共有メディアを正常に通信して提示できる必要があります。すべての WebRTC 互換ブラウザは、ビデオのオーディオ、[VP8](#)、H.264 制約ベースラインプロファイルのオンライン測位ユーザーサポート (OPUS) と G711 をサポートしている必要があります。

WebRTC エコシステム外の一般的な音声ソリューションでは、さまざまなタイプの CODECs。一般的な CODECs には、北米向けの G.711 μ -law、G.711 A-law、G.729、G.722 などがあります。2つの異なる CODECs を使用している2つのデバイスが相互に通信する場合、メディアゲートウェイはデバイス間の CODEC フローを変換します。つまり、メディアゲートウェイはメディアを処理し、エンドデバイスが相互に通信できるようにします。

WebRTC でのプッシュ通知

WebRTC の実装は、モバイルデバイスで非常に一般的です。ウェブブラウザとは異なり、モバイルデバイスは WebSocket 接続を長時間開いたままにすることはできません。したがって、呼び出しやメッセージなど、すべての終了リクエストに対して WebRTC サーバーからのプッシュ通知に依存する必要があります。

[Amazon Simple Notification Service](#) (Amazon SNS) では、モバイルデバイスのアプリにプッシュ通知を送信できます。これらのアプリは、Apple iOS や Android などのさまざまなオペレーティングシステムで実行できます。次の図は、WebRTC 通知サーバーから WebRTC モバイルエンドポイントへのプッシュ通知フローの概要を示しています。

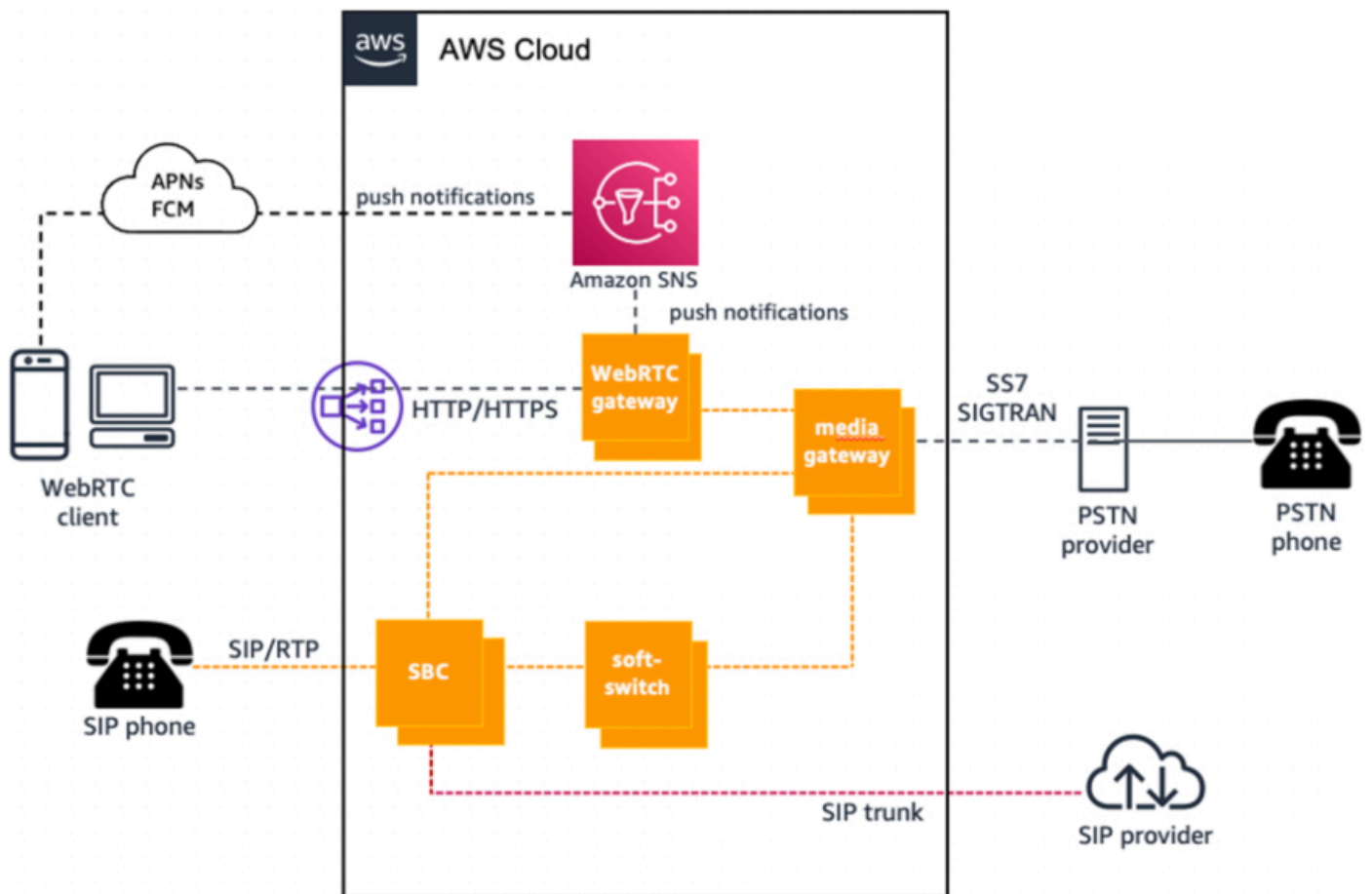


プッシュ通知用の Amazon SNS

WebRTC および WebRTC ゲートウェイ

ウェブリアルタイム通信 (WebRTC) を使用すると、ウェブブラウザからの呼び出しを確立したり、API を使用してバックエンドサーバーからリソースをリクエストしたりできます。このテクノロジーはクラウドテクノロジーを念頭に置いて設計されているため、呼び出しの確立に使用できるさまざまな APIs を提供します。すべての音声ソリューション (SIP を含む) がこれらの APIs をサポートしているわけではないため、API コールを SIP メッセージに変換するには WebRTC ゲートウェイが必要であり、その逆も同様です。

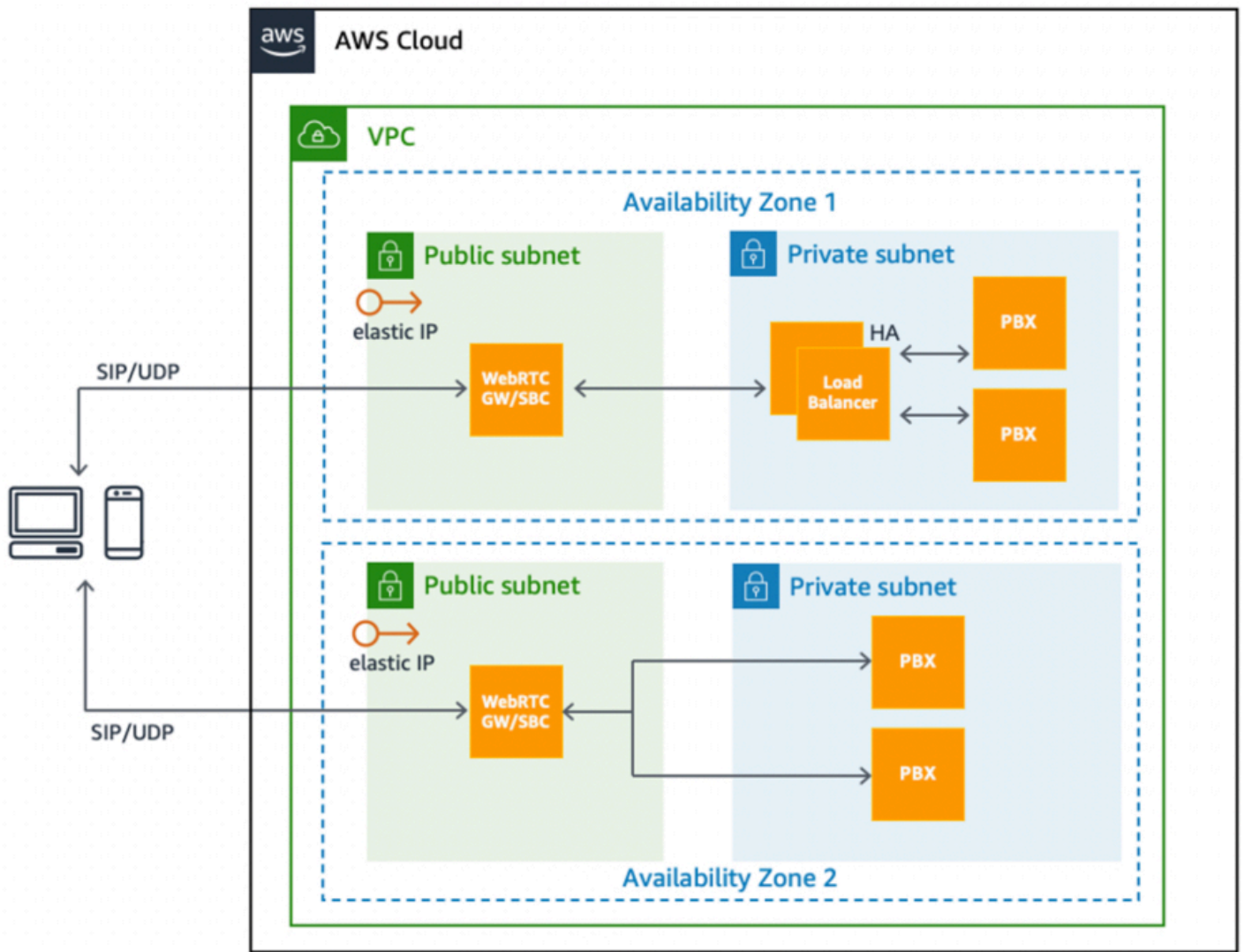
次の図は、高可用性 WebRTC アーキテクチャの設計パターンを示しています。WebRTC クライアントからの受信トラフィックは、[Amazon EC2 Auto Scaling](#) グループの一部である [Amazon Elastic Compute Cloud](#) (Amazon EC2) インスタンスで実行される WebRTC を備えた [Application Load Balancer](#) (ALB) によってバランスが取れています。



音声用の RTC システムの基本的なトポロジ

SIP トラフィックと RTP トラフィックのもう 1 つの設計パターンは、次の図に示すように、アベイラビリティゾーン全体でアクティブ/パッシブモードで Amazon EC2 の SBCs のペアを使用するこ

とです。ここでは、障害発生時に Elastic IP アドレスをインスタンス間で動的に移動できます。この場合、ドメインネームサービス (DNS) は使用できません。



Virtual Private Cloud (VPC) で Amazon EC2 を使用する RTC アーキテクチャ

での高可用性とスケーラビリティ AWS

リアルタイム通信のほとんどのプロバイダーは、99.9% から 99.999% までの可用性を提供するサービスレベルと一致しています。必要な高可用性 (HA) の度合いに応じて、アプリケーションのライフサイクル全体を通じて、ますます高度な対策を講じる必要があります。AWS では、堅牢な高可用性を実現するために、次のガイドラインに従うことを推奨しています。

- 単一障害点がないようにシステムを設計します。ステートレスコンポーネントとステートフルコンポーネントの両方で、自動モニタリング、障害検出、フェイルオーバーメカニズムを使用する
- 単一障害点 (SPOF) は、通常、N+1 または 2N 冗長設定で排除されます。N+1 はアクティブ/アクティブノード間の負荷分散によって実現され、2N はアクティブ/スタンバイ設定のノードのペアによって実現されます。
- AWS には、スケーラブルな負荷分散されたクラスターやアクティブスタンバイペアの引き受けなど、両方のアプローチで HA を達成するための方法がいくつかあります。
- システムの可用性を正しく計測してテストします。
- 障害に対応、軽減、復旧する手動メカニズムの運用手順を作成します。

このセクションでは、 の機能を使用して単一障害点を達成しない方法に焦点を当てます AWS。具体的には、このセクションでは、可用性の高いリアルタイム通信アプリケーションを構築できるコア AWS 機能と設計パターンのサブセットについて説明します。

アクティブ/スタンバイステートフルサーバー間の HA のフローティング IP パターン

フローティング IP 設計パターンは、ハードウェアノード (メディアサーバー) のアクティブペアとスタンバイペア間の自動フェイルオーバーを実現するよく知られたメカニズムです。静的セカンダリ仮想 IP アドレスがアクティブなノードに割り当てられます。アクティブノードとスタンバイノード間の継続的なモニタリングにより、障害が検出されます。アクティブなノードに障害が発生した場合、モニタリングスクリプトは仮想 IP を準備完了スタンバイノードに割り当て、スタンバイノードがプライマリアクティブ関数を引き継ぎます。このようにして、仮想 IP はアクティブノードとスタンバイノードの間で浮動します。

RTC ソリューションの適用性

N ノードのアクティブ/アクティブクラスターなど、同じコンポーネントの複数のアクティブインスタンスが稼働中であるとは限りません。アクティブスタンバイ設定は、HA に最適なメカニズムを提供します。例えば、メディアサーバーや会議サーバー、SBC やデータベースサーバーなど、RTC ソリューションのステートフルコンポーネントは、アクティブ/スタンバイ設定に適しています。SBC またはメディアサーバーには、一度に複数の長時間実行セッションまたはチャンネルがアクティブになっており、SBC アクティブインスタンスが失敗した場合、エンドポイントはフローティング IP によるクライアント側の設定なしでスタンバイノードに再接続できます。

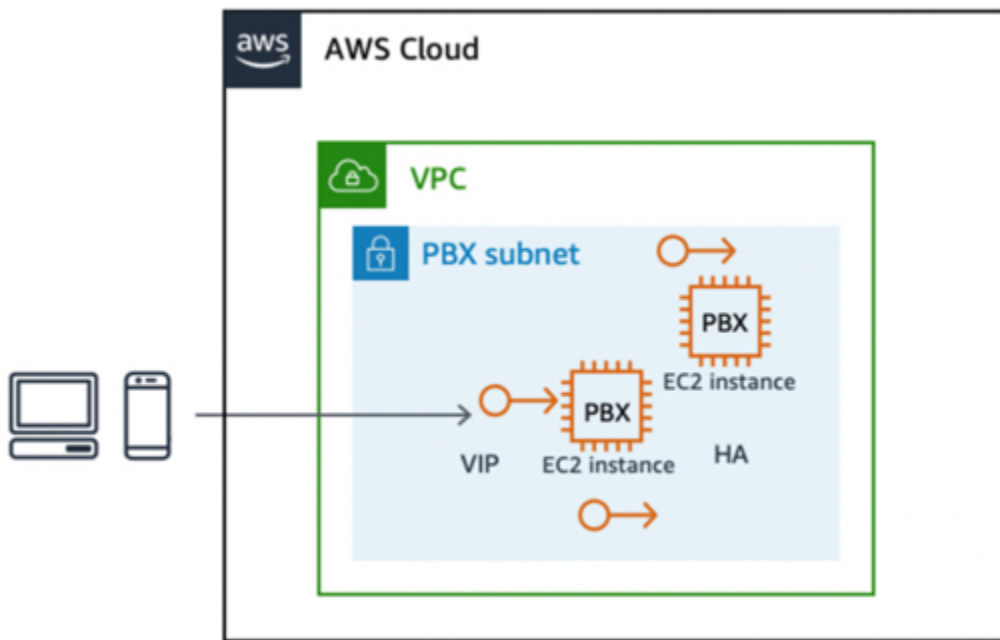
での実装 AWS

このパターンは、Amazon Elastic Compute Cloud (Amazon EC2)、Amazon EC2 API、Elastic IP アドレスのコア機能、および Amazon EC2 でのセカンダリプライベート IP アドレスのサポートを使用して AWS に実装できます。

フローティング IP パターンを実装するには AWS :

1. 2 つの EC2 インスタンスを起動して、プライマリノードとセカンダリノードのロールを引き受けます。プライマリノードはデフォルトでアクティブ状態であると想定されます。
2. プライマリ EC2 インスタンスに追加のセカンダリプライベート IP アドレスを割り当てます。
3. 仮想 IP (VIP) に似た Elastic IP アドレスは、セカンダリプライベートアドレスに関連付けられます。このセカンダリプライベートアドレスは、外部エンドポイントがアプリケーションにアクセスするために使用するアドレスです。
4. 一部のオペレーティングシステム (OS) 設定は、セカンダリ IP アドレスをプライマリネットワークインターフェイスにエイリアスとして追加するために必要です。
5. アプリケーションはこの Elastic IP アドレスにバインドする必要があります。アスタリスクソフトウェアの場合、高度なアスタリスク SIP 設定を使用してバインディングを設定できます。
6. 各ノードでカスタム、Linux の KeepAlive、Corosync などのモニタリングスクリプトを実行して、ピアノードの状態をモニタリングします。現在のアクティブなノードに障害が発生した場合、ピアはこの障害を検出し、Amazon EC2 API を呼び出してセカンダリプライベート IP アドレスを自身に再割り当てします。

したがって、セカンダリプライベート IP アドレスに関連付けられた VIP でリッスンしていたアプリケーションは、スタンバイノードを介してエンドポイントで使用できるようになります。



Elastic IP アドレスを使用したステートフル EC2 インスタンス間のフェイルオーバー

利点

このアプローチは、EC2 インスタンス、インフラストラクチャ、またはアプリケーションレベルでの障害から保護する信頼性の高い予算の少ないソリューションです。

制限と拡張性

この設計パターンは通常、単一のアベイラビリティゾーン内に限定されます。2つのアベイラビリティゾーンに実装できますが、バリエーションがあります。この場合、利用可能な Elastic IP アドレスの再関連付け API を介して、異なるアベイラビリティゾーンのアクティブノードとスタンバイノード間でフローティング Elastic IP アドレスが再関連付けされます。前の図に示すフェイルオーバー実装では、進行中の呼び出しは削除され、エンドポイントは再接続する必要があります。基盤となるセッションデータのレプリケーションを使用してこの実装を拡張し、セッションのシームレスなフェイルオーバーやメディア継続性を実現することもできます。

WebRTC と SIP によるスケーラビリティと HA の負荷分散

ラウンドロビン、アフィニティ、レイテンシーなどの事前定義されたルールに基づくアクティブなインスタンスのクラスターの負荷分散は、HTTP リクエストのステートレスな性質によって広く普及している設計パターンです。実際、ロードバランシングは、多くの RTC アプリケーションコンポーネントの場合に実行可能なオプションです。

ロードバランサーは、目的のアプリケーションへのリクエストのリバースプロキシまたはエン트리ポイントとして機能し、それ自体が複数のアクティブなノードで同時に実行されるように設定されています。任意の時点で、ロードバランサーは定義されたクラスター内のアクティブなノードのいずれかにユーザーリクエストを送信します。ロードバランサーは、ターゲットクラスター内のノードに対してヘルスチェックを実行し、ヘルスチェックに失敗したノードに受信リクエストを送信しません。したがって、ロードバランシングによって、基本的なレベルの高可用性が達成されます。また、ロードバランサーはすべてのクラスターノードに対してアクティブおよびパッシブのヘルスチェックを1秒未満の間隔で実行するため、フェイルオーバーの時間はほぼ瞬時に行われます。

どのノードを指示するかは決定は、ロードバランサーで定義されたシステムルールに基づきます。

- ラウンドロビン
- セッションまたは IP アフィニティ。セッション内または同じ IP からの複数のリクエストがクラスター内の同じノードに送信されます。
- レイテンシーベース
- 負荷ベース

RTC アーキテクチャの適用性

WebRTC プロトコルを使用すると、[Elastic Load Balancing](#) (ELB)、[Application Load Balancer](#) (ALB)、[Network Load Balancer](#) (NLB) などの HTTP ベースのロードバランサーを介して WebRTC Gateway を簡単にロードバランシングできます。[Load Balancer](#) ほとんどの SIP 実装では、Transmission Control Protocol (TCP) と User Datagram Protocol (UDP) の両方を介したトランスポートに依存しているため、TCP および UDP ベースのトラフィックの両方をサポートするネットワークレベルまたは接続レベルの負荷分散が必要です。

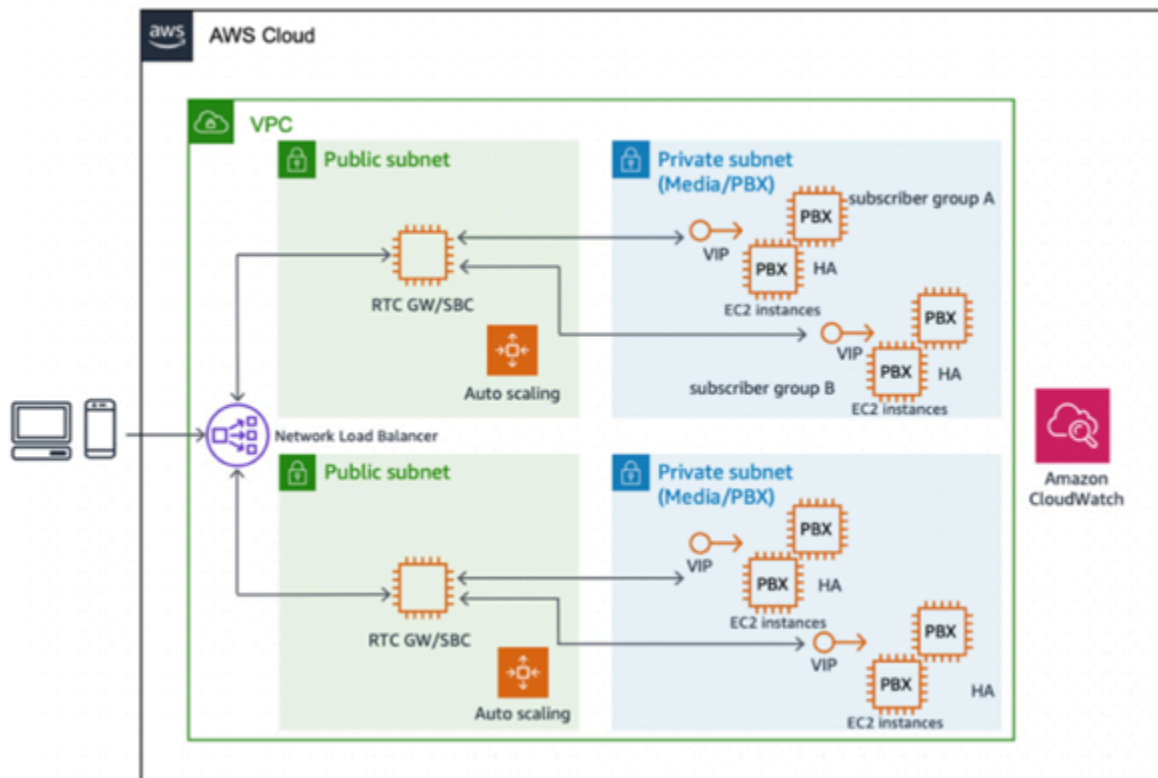
Application Load Balancer と Auto Scaling を使用した AWS for WebRTC のロードバランシング

WebRTC ベースの通信の場合、Elastic Load Balancing は、フルマネージドで可用性が高くスケラブルなロードバランサーを提供し、リクエストのエン트리ポイントとして機能します。ロードバランサーは、Elastic Load Balancing に関連付けられた EC2 インスタンスのターゲットクラスターに送信されます。WebRTC リクエストはステートレスであるため、Amazon EC2 Auto Scaling を使用して、完全に自動化され制御可能なスケラビリティ、伸縮性、高可用性を提供できます。

Application Load Balancer は、複数のアベイラビリティーゾーンを使用して可用性が高く、スケラブルなフルマネージド型のロードバランシングサービスを提供します。これにより、WebRTC アプ

リケーションのシグナリングを処理する WebSocket WebSocket リクエストのロードバランシングと、長時間実行される TCP 接続を使用したクライアントとサーバー間の双方向通信がサポートされます。Application Load Balancer は、コンテンツベースのルーティングと [スティッキーセッション](#) もサポートし、ロードバランサーが生成した Cookie を使用して、同じクライアントから同じターゲットにリクエストをルーティングします。スティッキーセッションを有効にすると、同じターゲットがリクエストを受け取り、Cookie を使用してセッションコンテキストを復元できます。

次の図は、ターゲットトポロジを示しています。



WebRTC スケーラビリティと高可用性アーキテクチャ

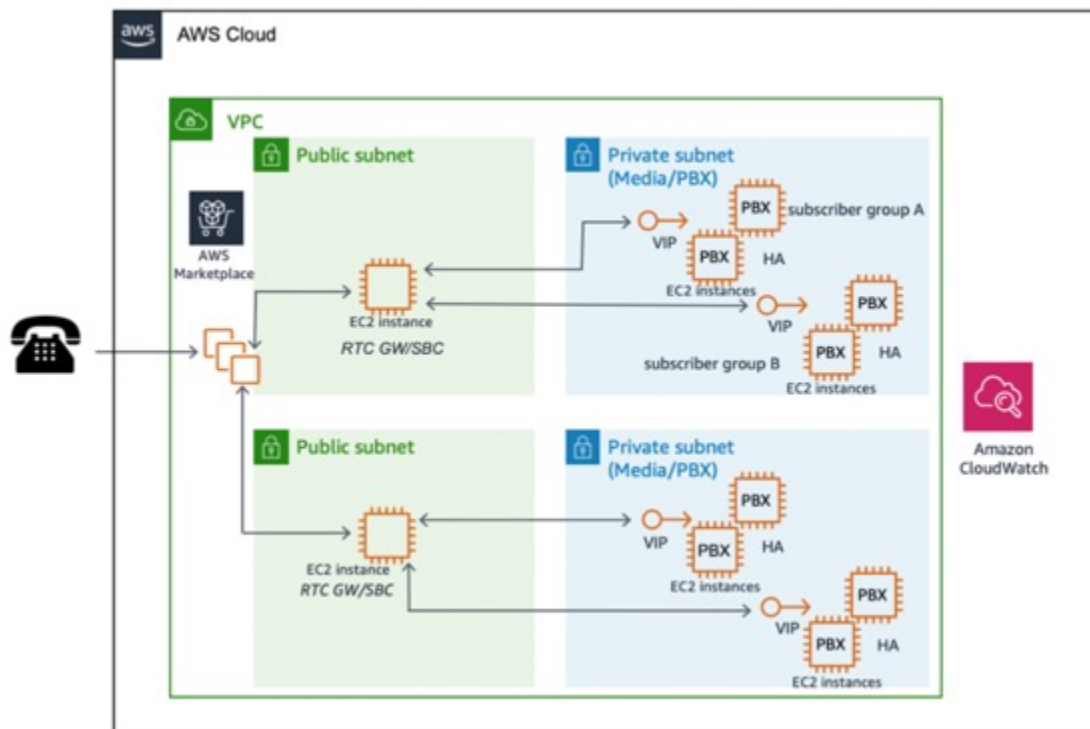
Network Load Balancer または 製品を使用した SIP の AWS Marketplace 実装

SIP ベースの通信の場合、接続は TCP または UDP 経由で行われ、RTC アプリケーションの大部分は UDP を使用します。SIP/TCP が最適なシグナルプロトコルである場合は、Network Load Balancer を使用して、フルマネージド、高可用性、スケーラブル、パフォーマンスの負荷分散を行うことができます。

Network Load Balancer は接続レベル (レイヤー 4) で動作し、IP プロトコルデータに基づいて Amazon EC2 インスタンス、コンテナ、IP アドレスなどのターゲットへの接続をルーティングし

ます。TCP または UDP トラフィック負荷分散に最適です。ネットワーク負荷分散は、非常に低いレイテンシーを維持しながら、1 秒あたり数百万のリクエストを処理できます。Amazon EC2 Auto Scaling、Amazon [Elastic Container Service \(Amazon ECS\)](#)、[Amazon Elastic Kubernetes Service \(Amazon EKS\)](#)、など、他の一般的な AWS サービスと統合されています [AWS CloudFormation](#)。

SIP 接続が開始された場合のもう 1 つのオプションは、[AWS Marketplace](#) 市販 off-the-shelf ソフトウェア (COTS) を使用することです。は、UDP やその他のタイプのレイヤー 4 接続負荷分散を処理できる多くの製品 AWS Marketplace を提供します。COTS には通常、高可用性のサポートが含まれており、Amazon EC2 Auto Scaling などの機能と統合して、可用性とスケーラビリティをさらに強化します。次の図は、ターゲットトポロジを示しています。



AWS Marketplace 製品による SIP ベースの RTC スケーラビリティ

クロスリージョン DNS ベースのロードバランシングとフェイルオーバー

[Amazon Route 53](#) は、RTC クライアントがメディアアプリケーションに登録して接続するためのパブリックエンドポイントまたはプライベートエンドポイントとして使用できるグローバル DNS サービスを提供します。Amazon Route 53 では、DNS ヘルスチェックを設定して、トラフィックを正常なエンドポイントにルーティングしたり、アプリケーションのヘルスを個別にモニタリングしたりできます。

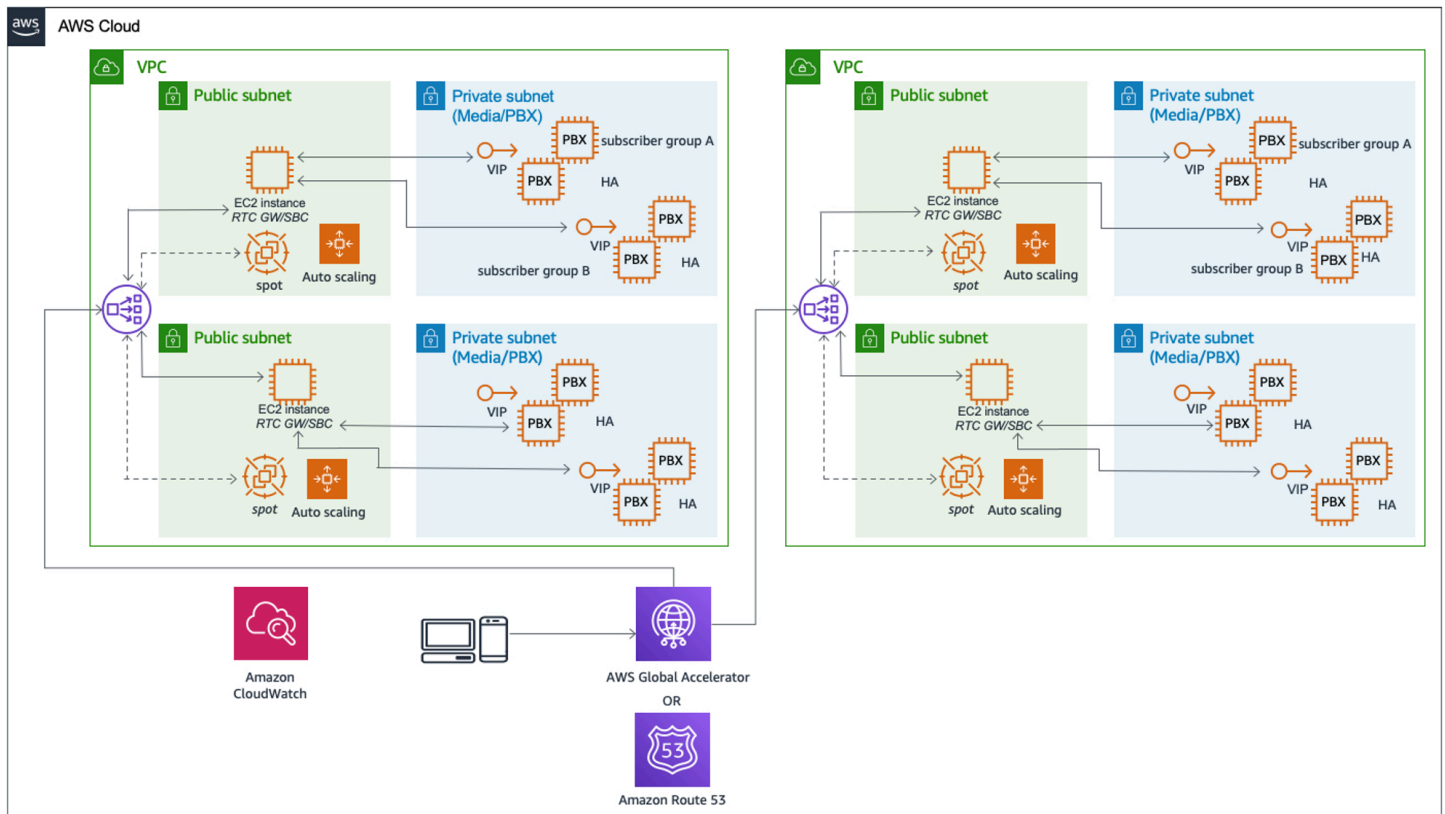
Amazon Route 53 トラフィックフロー機能を使用すると、レイテンシーベースのルーティング、地理的 DNS、地理的近接性、加重ラウンドロビンなど、さまざまなルーティングタイプを通じてトラフィックをグローバルに管理することが容易になります。これらはすべて DNS フェイルオーバーと組み合わせることで、低レイテンシーで耐障害性のあるさまざまなアーキテクチャを実現できます。Amazon Route 53 トラフィックフローのシンプルなビジュアルエディタを使用すると、エンドユーザーがアプリケーションのエンドポイントにルーティングされる方法を、単一の AWS リージョン内か世界中に分散しているかにかかわらず管理できます。

グローバルデプロイの場合、Route 53 のレイテンシーベースのルーティングポリシーは、リアルタイムメディア交換に関連するサービスの品質を向上させるために、メディアサーバーの最も近い場所に顧客を誘導するために特に役立ちます。

新しい DNS アドレスへのフェイルオーバーを適用するには、クライアントキャッシュをフラッシュする必要があります。また、DNS 変更がグローバル DNS サーバーに伝播されるため、遅延が発生する可能性があります。Time to Live 属性を使用して、DNS ルックアップの更新間隔を管理できます。この属性は、DNS ポリシーの設定時に設定できます。

グローバルユーザーにすばやく連絡するため、または単一のパブリック IP を使用する要件を満たすために、をクロスリージョンフェイルオーバーに使用する AWS Global Accelerator こともできます。[AWS Global Accelerator](#)は、ローカルとグローバルの両方のリーチを持つアプリケーションの可用性とパフォーマンスを向上させるネットワークサービスです。は、Application Load Balancer、Network Load Balancer、Amazon EC2 インスタンスなどのアプリケーションエンドポイントへの固定エン트리ポイントとして機能する静的 IP アドレスを単一または複数の AWS リージョンで AWS Global Accelerator 提供します。AWS グローバルネットワークを使用して、ユーザーからアプリケーションへのパスを最適化し、TCP および UDP トラフィックのレイテンシーなどのパフォーマンスを向上させます。

AWS Global Accelerator はアプリケーションエンドポイントの正常性を継続的にモニタリングし、現在のエンドポイントが異常になった場合にトラフィックを最も近い正常なエンドポイントに自動的にリダイレクトします。追加のセキュリティ要件として、高速 Site-to-Site VPN は AWS Global Accelerator を使用して、AWS グローバルネットワークと AWS エッジロケーションを介してトラフィックをインテリジェントにルーティングすることで、VPN 接続のパフォーマンスを向上させます。



AWS Global Accelerator または Amazon Route 53 を使用したリージョン間の高可用性設計

永続的ストレージによるデータの耐久性と HA

ほとんどの RTC アプリケーションは、認証、認可、アカウントリング (セッションデータ、通話詳細レコードなど)、運用モニタリング、ログ記録のためにデータを保存およびアクセスするために永続的ストレージに依存しています。従来のデータセンターでは、永続的ストレージコンポーネント (データベース、ファイルシステムなど) の高可用性と耐久性を確保するには、通常、ストレージエリアネットワーク (SAN)、独立ディスクの冗長配列 (RAID) 設計、バックアップ、復元、フェイルオーバー処理のプロセスのセットアップによる手間のかかる作業が必要です。は、データの耐久性と可用性に関する従来のデータセンターのプラクティス AWS クラウド を大幅に簡素化し、強化します。

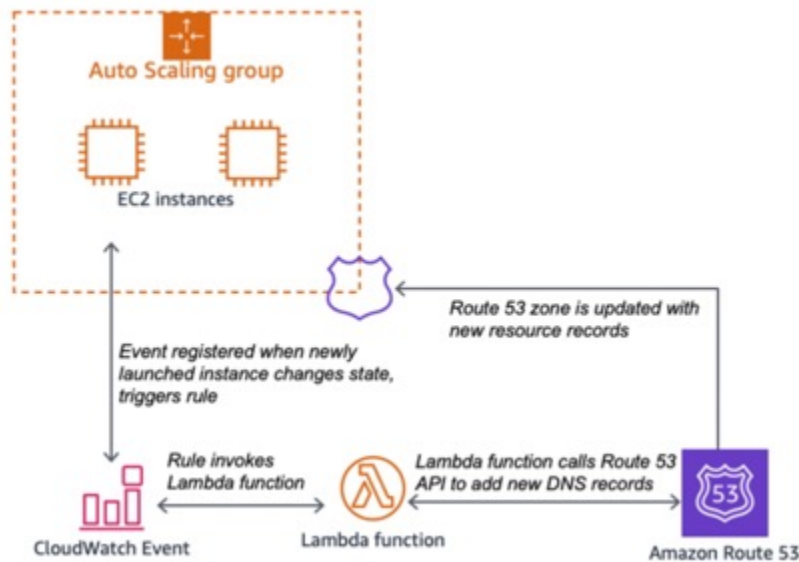
オブジェクトストレージとファイルストレージの場合、[Amazon Simple Storage Service](#) (Amazon S3) や [Amazon Elastic File System](#) (Amazon EFS) などの AWS サービスは、マネージド型の高可用性とスケーラビリティを提供します。Amazon S3 のデータ耐久性は 99.999999999% (11 nines) です。

トランザクションデータストレージでは、高可用性デプロイで Amazon Aurora、PostgreSQL、MySQL、MariaDB、Oracle、Microsoft SQL Server をサポートするフル

マネージド Amazon Relational Database Service (Amazon RDS) を利用できます。PostgreSQL、MySQL、MariaDB レジストラ関数、サブスクリバードプロファイル、またはアカウントレコードストレージ (CDRs) の場合、Amazon RDS は耐障害性があり、可用性が高く、スケーラブルなオプションを提供します。

Amazon Route 53、AWS Lambda、Amazon EC2 Auto Scaling による動的スケーリング

AWS では、機能の連鎖と、インフラストラクチャイベントに基づくサービスとしてカスタムサーバーレス関数を組み込む機能を使用できます。RTC アプリケーションで多くの汎用性がある設計パターンの 1 つは、Auto Scaling ライフサイクルフックと [Amazon CloudWatch Events](#)、Amazon Route 53、および [AWS Lambda functions](#) の組み合わせです。functions は、任意のアクションまたはロジックを埋め込むことができます。次の図は、これらの機能が連携して、自動化によってシステムの信頼性とスケーラビリティを向上させる方法を示しています。



Amazon Route 53 の動的更新による自動スケーリング

Amazon Kinesis Video Streams で高可用性 WebRTC

[Amazon Kinesis Video Streams](#) は、WebRTC を介したリアルタイムのメディアストリーミングを提供するため、ユーザーは再生、分析、機械学習のためにメディアストリームをキャプチャ、処理、保存できます。これらのストリームは可用性が高く、スケーラブルで、WebRTC 標準に準拠しています。Amazon Kinesis Video Streams には、高速ピア検出と安全な接続確立のための WebRTC シグナリングエンドポイントが含まれています。これには、NAT (STUN) 用のマネージドセッション

トラバーサルユーティリティと、ピア間でメディアをリアルタイムで交換するための NAT (TURN) エンドポイントに関するリレーを使用したトラバーサルが含まれます。また、無料のオープンソース SDK が含まれており、カメラファームウェアと直接統合して Amazon Kinesis Video Streams エンドポイントとの安全な通信を可能にし、ピア検出とメディアストリーミングを可能にします。最後に、Android、iOS、JavaScript 用のクライアントライブラリを提供します。これにより、WebRTC 準拠のモバイルおよびウェブプレーヤーは、メディアストリーミングと双方向通信のためにカメラデバイスを安全に検出して接続できます。

Amazon Chime Voice Connector での高可用性 SIP トランキング

[Amazon Chime Voice Connector](#) は、pay-as-you-go SIP トランキングサービスを提供します。これにより、企業は電話システムを使用して安全で安価な電話を発信または受信できます。Amazon Chime Voice Connector は、サービスプロバイダーの SIP トランクまたは統合サービスデジタルネットワーク (ISDN) プライマリレートインターフェイス (PRIs)。お客様は、インバウンド通話、アウトバウンド通話、またはその両方を有効にすることができます。

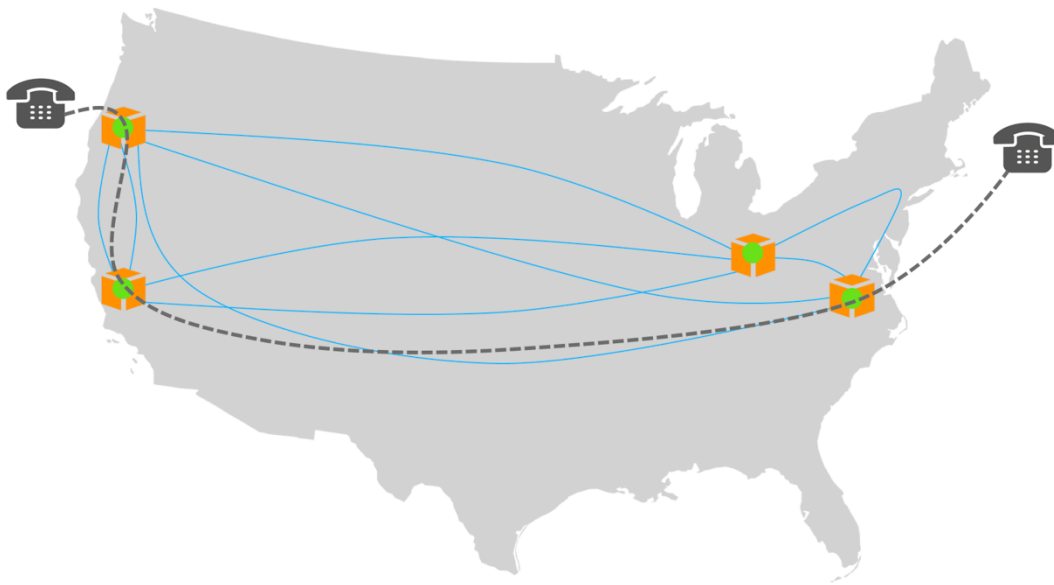
このサービスは、AWS ネットワークを使用して、複数の にわたって高可用性の通話エクスペリエンスを提供します AWS リージョン。SIP トランキング通話、または転送された SIP ベースのメディア録画 (SIPREC) フィードから Amazon Kinesis Video Streams に音声を実時間ストリーミングして、ビジネス通話からリアルタイムでインサイトを得ることができます。 [Amazon Transcribe](#) やその他の一般的な機械学習ライブラリとの統合により、オーディオ分析用のアプリケーションをすばやく構築できます。

フィールドのベストプラクティス

このセクションでは、大規模なリアルタイムセッション開始プロトコル (SIP) ワークロードを実行する、大規模で最も成功した AWS 顧客の一部が実装したベストプラクティスを要約します。AWS のお客様は、さまざまな種類の障害が発生した場合にシステムの信頼性と回復性を向上させることができるため、独自の SIP インフラストラクチャをパブリッククラウドで実行したいと考えているため、これらのベストプラクティスが役立ちます。これらのベストプラクティスの一部は SIP 固有ですが、そのほとんどは で実行されているリアルタイム通信アプリケーションに適用されます AWS。

SIP オーバーレイを作成する

AWS には、異なる 間の接続を提供する堅牢でスケーラブルな冗長ネットワークバックボーンがあります AWS リージョン。ファイバーカットなどのネットワークイベントによって AWS バックボーンリンクが劣化すると、ボーダーゲートウェイプロトコル (BGP) などのネットワークレベルのルーティングプロトコルを使用して、トラフィックが冗長パスにすばやくフェイルオーバーされます。このネットワークレベルのトラフィックエンジニアリングは AWS お客様にとってブラックボックスであり、ほとんどののはこれらのフェイルオーバーイベントに気付くことさえありません。ただし、音声、高品質の動画、低レイテンシーのメッセージングなどのリアルタイムワークロードを実行するお客様は、これらのイベントに気付くことがあります。では、ネットワーク AWS レベルで が提供するものに加えて、 の AWS お客様が独自のトラフィックエンジニアリングを実装するにはどうすればよいですか？ このソリューションは、SIP インフラストラクチャをさまざまな方法でデプロイします AWS リージョン。SIP は、通話制御機能の一部として、特定の SIP プロキシを介して通話をルーティングする機能も提供します。

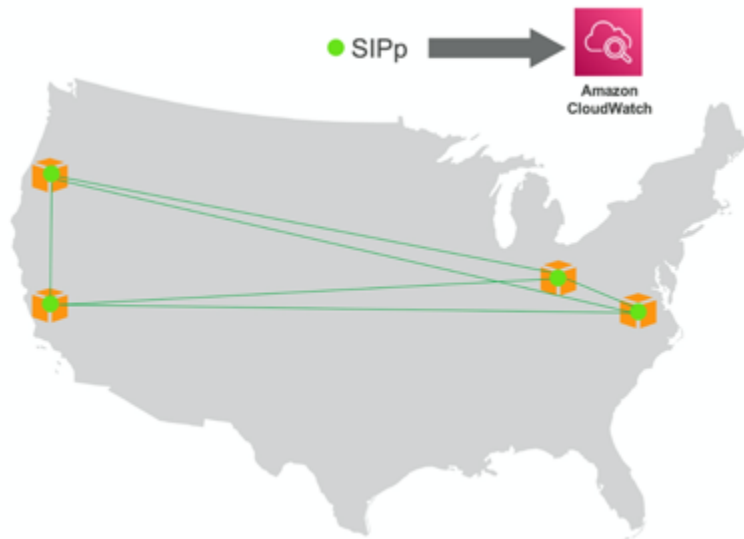


SIP ルーティングを使用したネットワークルーティングの上書き

前の図では、SIP インフラストラクチャ (キューブ内の緑色のドットで表されます) は 4 つの米国リージョンすべてで実行されています。青い実線は、AWS バックボーンの架空の表現を表します。SIP ルーティングが実装されていない場合、米国西部沿岸を起点とし、米国東部沿岸を起点とする通話は、オレゴンリージョンとバージニアリージョンを直接接続しているバックボーンリンクを経由します。この図は、顧客がネットワークレベルのルーティングをオーバーライドし、SIP ルーティングを使用してカリフォルニア経由でルーティングされたオレゴン州とバージニア州の間で同じ通話を行う方法を示しています。このタイプの SIP トラフィックエンジニアリングは、SIP 再送信や顧客固有のビジネス設定などのネットワークメトリクスに基づいて、SIP プロキシとメディアゲートウェイを使用して実装できます。

詳細モニタリングを実行する

リアルタイム音声およびビデオアプリケーションのエンドユーザーは、従来のテレフォニーサービスで達成したのと同じレベルのパフォーマンスを期待します。したがって、アプリケーションに問題が発生すると、プロバイダーの評価が低下します。事後対応型ではなく事前対応型であるためには、エンドユーザーに提供するシステムのあらゆる部分に詳細モニタリングをデプロイすることが不可欠です。



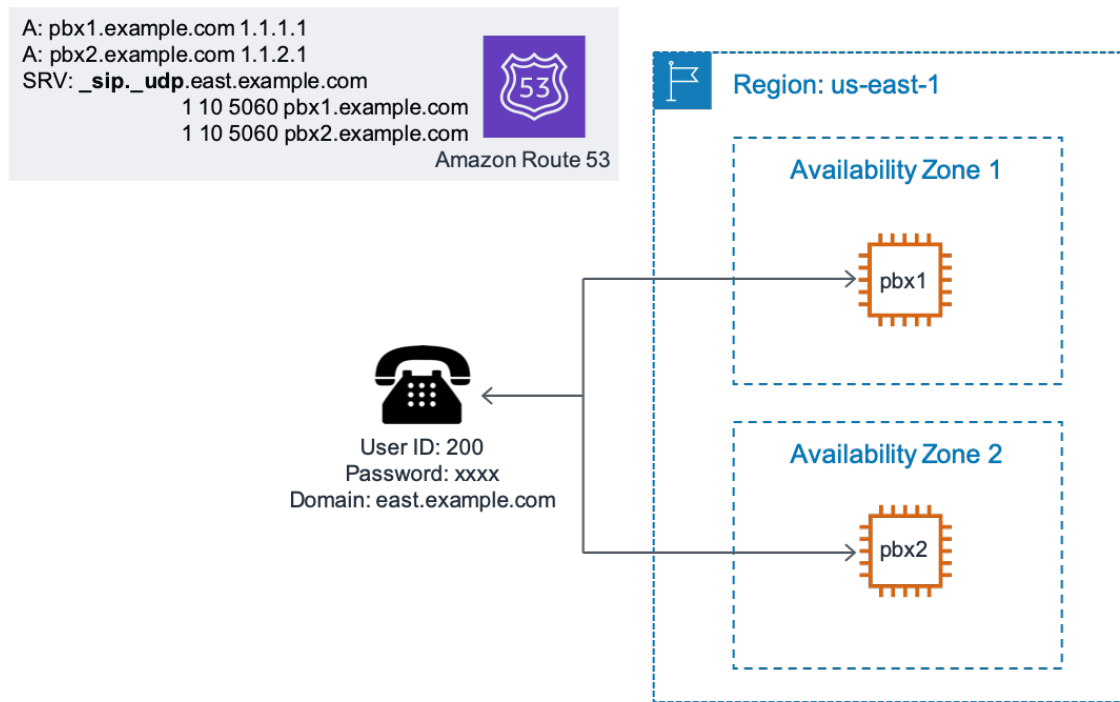
SIPp を使用した VoIP インフラストラクチャのモニタリング

[iPerf](#) や [SIPp](#)、[VOIPMonitor](#) などの多くのオープンソースツールは、SIP/RTP トラフィックのモニタリングに使用できます。前の例では、クライアントモードとサーバーモードで SIP を実行しているノードは、4 つのすべての米国間で成功した通話や SIP 再送信などの SIP メトリクスを測定しています AWS リージョン。これらのメトリクスは、カスタムスクリプトを使用して Amazon CloudWatch にエクスポートできます。CloudWatch を使用すると、特定のしきい値に基づいてこれらのカスタムメトリクスにアラームを作成できます。その後、これらの CloudWatch アラームの状態に基づいて、自動または手動の修復アクションを実行できます。

カスタムモニタリングシステムの開発と保守に必要なエンジニアリングリソースを割り当てたくないお客様にとって、[ThousandEyes](#) など、多くの優れた VoIP モニタリングソリューションが市場に用意されています。修復アクションの例は、SIP 再送信の増加に基づいて SIP ルーティングを変更することです。

ロードバランシングに DNS を使用し、フェイルオーバーにフローティング IPs

DNS SRV 機能をサポートする IP テレフォニークライアントは、異なる SBCs/PBXs にクライアントをロードバランシングすることで、インフラストラクチャに組み込まれている冗長性を効率的に使用できます。



DNS SRV レコードを使用した SIP クライアントの負荷分散

前の図は、SRV レコードを使用して SIP トラフィックを負荷分散する方法を示しています。SRV 標準をサポートする IP テレフォニークライアントは、SRV タイプの DNS レコードで `sip.<transport protocol>` プレフィックスを探します。この例では、DNS の回答セクションに、異なる AWS アベイラビリティゾーンで実行されている両方の PBXs が含まれています。ただし、エンドポイント URIs に加えて、SRV レコードには 3 つの追加情報が含まれています。

- 最初の数値は Priority (上記の例では 1) です。優先度は、優先度よりも低いほうが推奨されます。
- 2 番目の数値は重み (上の例では 10) です。
- 3 番目の番号は、使用するポート (5060) です。

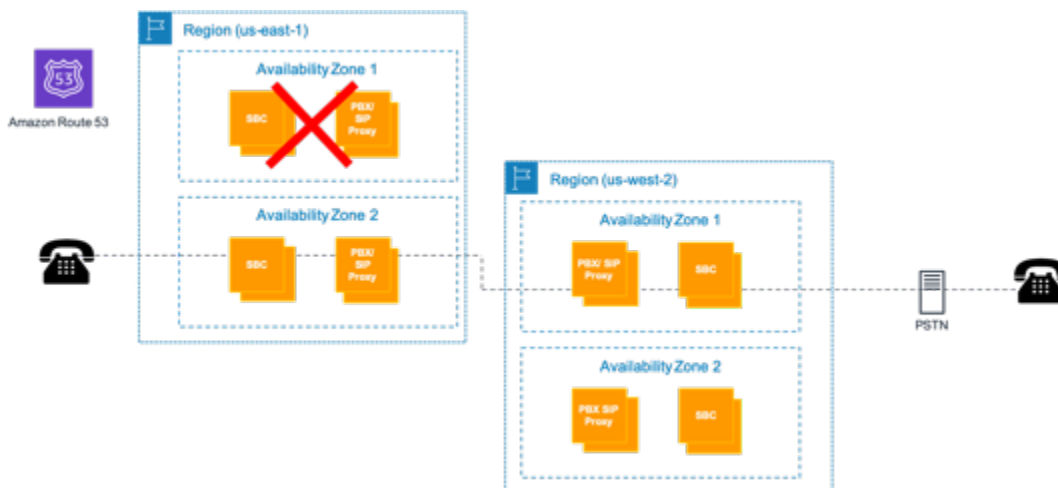
優先度は両方の PBXs サーバーで同じ (1) であるため、クライアントは重みを使用して 2 つの PBXs。この場合、重みは同じであるため、SIP トラフィックは 2 つの PBXs 間で均等に負荷分散する必要があります。

DNS はクライアントロードバランシングに適したソリューションかもしれませんが、DNS の「A」レコードを変更/更新してフェイルオーバーを実装するにはどうすればよいですか？ この方法は、クライアントノードと中間ノード内の DNS キャッシュ動作で不整合が見つかったため、推奨されません。SIP ノードのクラスター間の AZ 内フェイルオーバーのより良い方法は、EC2 API を使用して、障害のあるホストの IP アドレスが正常なホストに即座に再割り当てされる EC2 IP 再割り当てを使

用することです。詳細なモニタリングおよびヘルスチェックソリューションと組み合わせることで、障害が発生したノードの IP 再割り当てにより、トラフィックが正常なホストにタイムリーに移行され、エンドユーザーの中断を最小限に抑えることができます。

複数のアベイラビリティーゾーンを使用する

各 AWS リージョンは別々のアベイラビリティーゾーンに分割されます。各アベイラビリティーゾーンには独自の電源、冷却、ネットワーク接続があるため、独立した障害ドメインを形成します。この構造内では AWS、複数のアベイラビリティーゾーンでワークロードを実行することをお勧めします。これにより、お客様のアプリケーションは完全なアベイラビリティーゾーンの障害にも耐えることができます。これは非常にまれなイベントです。このレコメンデーションは、リアルタイム SIP インフラストラクチャも表しています。



アベイラビリティーゾーンの障害の処理

壊滅的なイベント (カテゴリ 5 の台風など) により、us-east-1 リージョンでアベイラビリティーゾーンが完全に停止したとします。図に示されているようにインフラストラクチャが実行されている場合、障害が発生したアベイラビリティーゾーンのノードに最初に登録されたすべての SIP クライアントは、アベイラビリティーゾーン #2 で実行されている SIP ノードに再登録する必要があります。(SIP クライアント/電話を使用してこの動作をテストし、サポートされていることを確認します。) アベイラビリティーゾーンの停止時のアクティブな SIP コールは失われますが、新しいコールはアベイラビリティーゾーン 2 を介してルーティングされます。

要約すると、DNS SRV レコードは、各アベイラビリティーゾーンに 1 つずつ、複数の「A」レコードをクライアントにポイントする必要があります。これらの「ASBCs/PBXs の複数の IP アドレスを指し、アベイラビリティーゾーン内とアベイラビリティーゾーン間の回復性の両方を提供する必要があります。アベイラビリティーゾーン内およびアベイラビリティーゾーン間のフェイルオーバー

は、IP がパブリックである場合に IPs の再割り当てを使用して実装できます。ただし、アベイラビリティゾーン間でプライベート IPs を再割り当てすることはできません。お客様がプライベート IP アドレス指定を使用している場合は、SIP クライアントがバックアップ SBC/PBX に再登録して、アベイラビリティゾーン間のフェイルオーバーを行う必要があります。

トラフィックを 1 つのアベイラビリティゾーンに保持し、EC2 プレイスマントグループを使用する

アベイラビリティゾーンのアフィニティとも呼ばれるこのベストプラクティスは、アベイラビリティゾーンに障害が発生した場合のまれなイベントにも適用されます。1 つのアベイラビリティゾーンに入る SIP または RTP トラフィックがリージョンを終了するまでそのアベイラビリティゾーンに留まるように、クロス AZ トラフィックを削除することをお勧めします。



アベイラビリティゾーンのアフィニティ (アクティブな呼び出しの最大 50% が失われます)

上の図は、アベイラビリティゾーンのアフィニティを使用する簡略化されたアーキテクチャを示しています。このアプローチの比較上の利点は、アベイラビリティゾーンが完全に停止した場合の影響を考慮すると明らかになります。図に示すように、アベイラビリティゾーン 2 が失われた場合、アクティブな呼び出しの 50% が最大でも影響を受けます (アベイラビリティゾーン間の負荷分散が等しいと仮定)。アベイラビリティゾーンのアフィニティが実装されていない場合、一部の呼び出しは 1 つのリージョンのアベイラビリティゾーン間で流れ、障害はアクティブな呼び出しの 50% 以上に影響を与える可能性が高くなります。

トラフィックのレイテンシーを最小限に抑えるために、AWS では各アベイラビリティゾーン内で [EC2 プレイスマントグループ](#) を使用することもお勧めします。同じ EC2 プレイスマントグループ内で起動されるインスタンスは、EC2 が相互にこれらのインスタンスのネットワーク近接性を確保するため、帯域幅が大きくなり、レイテンシーが短縮されます。

拡張ネットワーキング EC2 インスタンスタイプを使用する

Amazon EC2 で適切なインスタンスタイプを選択すると、システムの信頼性とインフラストラクチャの効率的な使用が保証されます。EC2 は、さまざまなユースケースに合わせて最適化された幅広いインスタンスタイプを提供します。インスタンスタイプは CPU、メモリ、ストレージ、およびネットワーク容量のさまざまな組み合わせで構成され、アプリケーションに適したリソースの組み合わせを柔軟に選択できます。これらの拡張ネットワーキングインスタンスタイプにより、それらで実行されている SIP ワークロードは、一貫した帯域幅と比較的低い集約レイテンシーにアクセスできます。Amazon EC2 への最近の追加は、最大 100 Gbps の帯域幅を提供する Elastic Network Adapter (ENA) の可用性です。EC2 インスタンスタイプの最新のカタログと関連する機能は、[EC2 インスタンスタイプページ](#)にあります。

ほとんどのお客様にとって、最新世代の [Compute Optimized インスタンス](#) はコストに最適な価値を提供する必要があります。例えば、C5N は、1 秒あたり数百万パケット (PPS) の帯域幅が最大 100 Gbps の新しい Elastic Network Adapter をサポートします。ほとんどのリアルタイムアプリケーションでは、ネットワークパケット処理を大幅に向上できる [Intel Data Plane Developer Kit \(DPDK\)](#) を使用することもメリットになります。

ただし、要件に応じてさまざまな EC2 インスタンスタイプをベンチマークして、どのインスタンスタイプが最適かを確認することが常にベストプラクティスです。ベンチマークを使用すると、特定のインスタンスタイプが一度に処理できる呼び出しの最大数など、他の設定パラメータも検索できます。

セキュリティに関する考慮事項

RTC アプリケーションコンポーネントは、通常、インターネットに接続する Amazon EC2 インスタンスで直接実行されます。TCP に加えて、フローは UDP や SIP などのプロトコルを使用します。このような場合、は、UDP リフレクション DDoS 攻撃、DNS リフレクション、NTP リフレクション、SSDP リフレクションなどの一般的なインフラストラクチャレイヤー (レイヤー 3 および 4) DDoS 攻撃から Amazon EC2 インスタンス AWS Shield Standard を保護します。は、明確に定義された DDoS 攻撃シグネチャが検出されたときに自動的にエンゲージされる優先度ベースのトラフィックシェーピングなどのさまざまな手法 AWS Shield Standard を使用します。

AWS また、は、Elastic IP アドレス AWS Shield Advanced で を有効にすることで、これらのアプリケーションの大規模で高度な DDoS 攻撃に対する高度な保護を提供します。は、EC2 インスタンスの AWS リソースのタイプとサイズを自動的に検出し、SYN または UDP フラッドに対する保護を備えた適切な事前定義された緩和策を適用する拡張 DDoS 検出 AWS Shield Advanced を提供します。を使用すると AWS Shield Advanced、お客様は 24 時間 365 日の AWS DDoS レスポンスチーム (DRT) を関与させることで、独自のカスタム緩和プロファイルを作成することもできます。AWS Shield Advanced また、DDoS 攻撃中に、すべての Amazon VPC ネットワークアクセスコントロールリスト (ACLs) が AWS ネットワークの境界に自動的に適用され、追加の帯域幅とスクラブ容量にアクセスして大規模な DDoS 攻撃を軽減できます。

結論

リアルタイム通信 (RTC) ワークロードを にデプロイ AWS して、主要な要件を満たしながらスケーラビリティ、伸縮性、高可用性を実現できます。現在、AWS、そのパートナー、およびオープンソースソリューションを使用して、コストを削減し、俊敏性を高め、グローバルフットプリントを削減した RTC ワークロードを実行しているお客様もいます。

このホワイトペーパーに記載されているリファレンスアーキテクチャとベストプラクティスは、お客様が で RTC ワークロードを正常にセットアップ AWS し、エンドユーザーの要件を満たすソリューションを最適化しながら、クラウドに合わせて最適化するのに役立ちます。

頭字語

このドキュメントで使用されている頭字語は次のとおりです。

ACL — アクセスコントロールリスト

ALB — Application Load Balancer

APNs Apple Push Notification サービス

BGP — ボーダーゲートウェイプロトコル

CDR — 通話詳細レコード

COTS — off-the-shelfソフトウェア

DDoS — 分散型denial-of-service

DNS — ドメインネームシステム

DPDK — インテルデータプレーン開発者キット

DRT — DDoS 対応チーム

ENA — Elastic Network Adapter

EPC – 進化したパケットコア

FCM — Firebase クラウドメッセージング

HA — 高可用性

IRC — インターネットリレーチャット

ISDN — 統合サービスデジタルネットワーク

NAT — ネットワークアドレス変換

OPUS — オンライン測位ユーザーサポート

PBX — プライベートブランチ交換

PRI — プライマリレートインターフェイス

PSTN — パブリックスイッチドテレフォンネットワーク

RAID — 独立ディスクの冗長配列

RTC — リアルタイム通信

RTP — リアルタイムトランスポートプロトコル

SAN — ストレージエリアネットワーク

SBC — セッションボーダーコントローラー

SIP — セッション開始プロトコル

SPOF — 単一障害点

SRV — サービス

SS7 — シグナリングシステム n.7

STUN — NAT のセッショントラバーサルユーティリティ

SYN — 同期

TCP — Transmission Control Protocol

TDM — 時間除算多重化

TURN — NAT のリレーを使用したトラバーサル

UDP — ユーザーデータグラムプロトコル

URI — ユニフォームリソース識別子

VIP — 仮想 IP

VNF — 仮想ネットワーク関数

VoIP — Voice Over IP

VPC — 仮想プライベートクラウド

WebRTC — ウェブリアルタイム通信

寄稿者

このドキュメントには、次の個人および組織が貢献しました。

- Amazon Web Services、Senior Solutions Architect、Mounir Chinnana
- Amazon Web Services、Senior Solutions Architect、Mohammed Al-Mehdar
- Amazon Web Services、Senior Solutions Architect、Ejaz Sial
- Amazon Web Services、Senior Solutions Architect、Ahmad Khan
- Amazon Web Services、プリンシパルエンジニア AWS サポート、Tipu Qureshi
- Hasan Khan、Amazon Web Services、シニアテクニカルアカウントマネージャー
- Amazon Web Services、Telecom、WW Technical Leader、Shooma Chakravarty

ドキュメントの改訂

このホワイトペーパーの更新に関する通知を受け取るには、RSS フィードにサブスクライブしてください。

変更	説明	日付
ホワイトペーパーの更新	最新の サービスと機能を更新しました。	2022 年 5 月 5 日
ホワイトペーパーの更新	最新の サービスと機能を更新しました。	2020 年 2 月 13 日
初版発行	ホワイトペーパーの初回発行。	2018 年 10 月 1 日

注意

お客様は、本書に記載されている情報を独自に評価する責任を負うものとし、本書は、(a) 情報提供のみを目的とし、(b) AWS の現行製品と慣行について説明しており、これらは予告なしに変更されることがあり、(c) AWS およびその関連会社、サプライヤー、またはライセンサーからの契約上の義務や保証をもたらすものではありません。AWS の製品やサービスは、明示または黙示を問わず、一切の保証、表明、条件なしに「現状のまま」提供されます。お客様に対する AWS の責任は AWS 契約によって規定されています。また、本文書は、AWS とお客様との間の契約に属するものではなく、また、当該契約が本文書によって修正されることもありません。

© 2022 Amazon Web Services, Inc. or its affiliates. All rights reserved.

AWS 用語集

最新の AWS 用語については、「AWS の用語集 リファレンス」の [AWS 「用語集」](#) を参照してください。

翻訳は機械翻訳により提供されています。提供された翻訳内容と英語版の間で齟齬、不一致または矛盾がある場合、英語版が優先します。