

AWS

S U M M I T

EC2 Innovation at Scale

Raj Pai, Director of Product Management, EC2

June 2, 2017



Amazon Elastic Compute Cloud (EC2) - クラウドの伸縮自在な仮想サーバー

EC2インスタンス

ゲスト1

ゲスト2

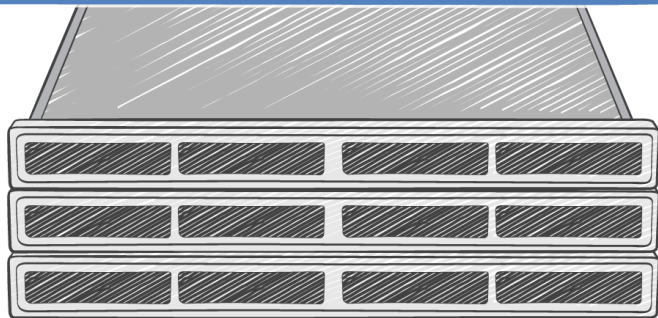
ゲスト n

ハイパーバイザー

ホストサーバー



AWSのグローバルリージョン
の物理サーバー



10年前のAmazon EC2...

シングルインスタンスファミリ／サイズ

- m1.small (1個のvCPU、1.7GiBのRAM、160GBのストレージ)

Linuxのみ

オンデマンド料金のみ



あれから10年

M4は64個のvCPU、256GBのRAM。当初の
m1.smallと比べて、vCPUは64倍、RAMは150倍！

それに加えて...

インスタンスの選択肢をさらに追加

リザーブドインスタンスとスポットインスタンス

OSとアプリケーションのサポート

Amazon Elastic Block Store (EBS)

Elastic IPアドレス

Amazon VPC

Auto Scaling

Elastic Load Balancing

パフォーマンス、セキュリティ、管理可能性、スケーラビリティの改善

Amazon ECS、Lambda

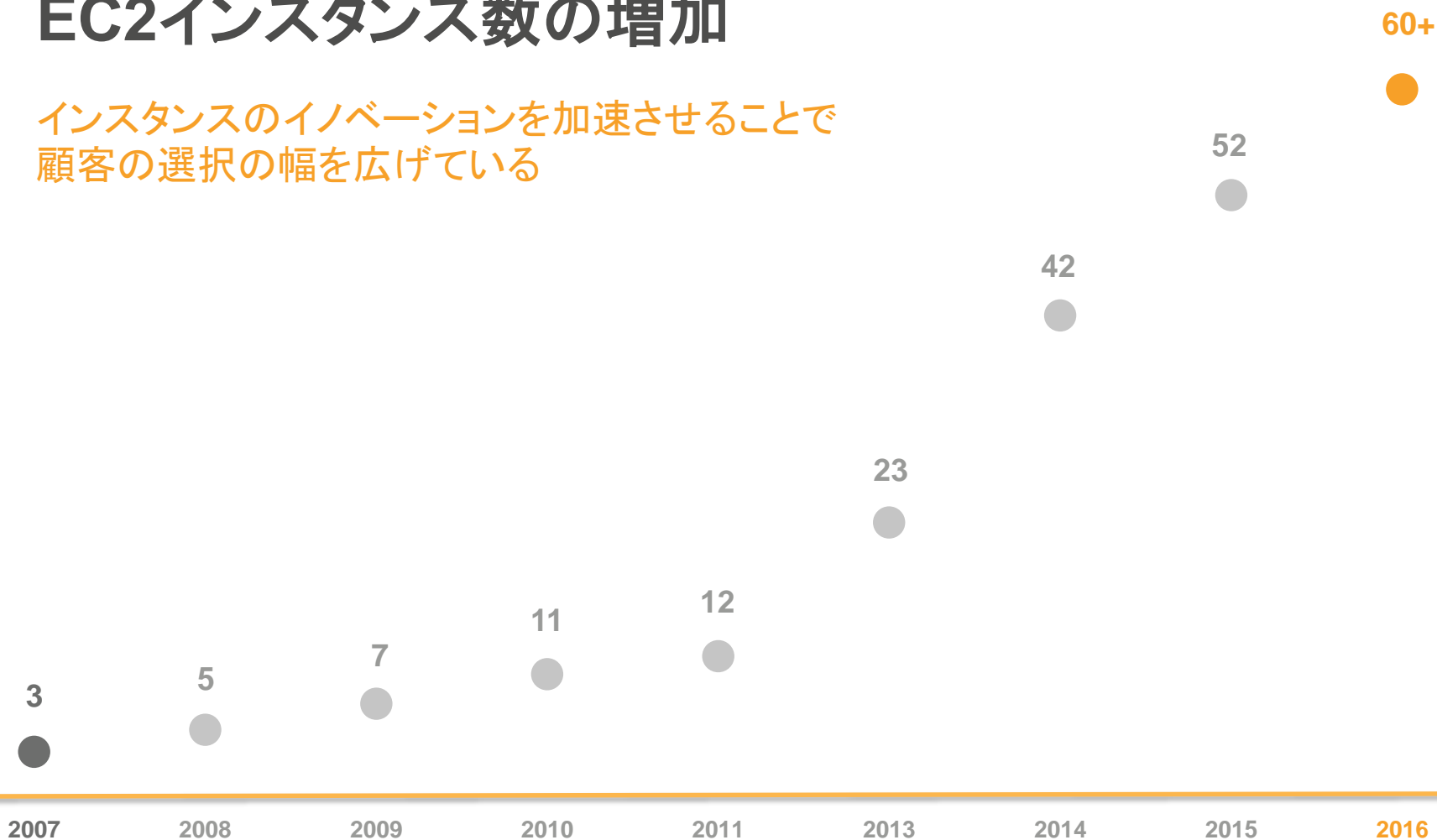
Amazon Machine Learning

他にもいろいろ



EC2インスタンス数の増加

インスタンスのイノベーションを加速させることで
顧客の選択の幅を広げている



現在のEC2インスタンスの特徴

CPU



メモリ



ストレージ



GPU

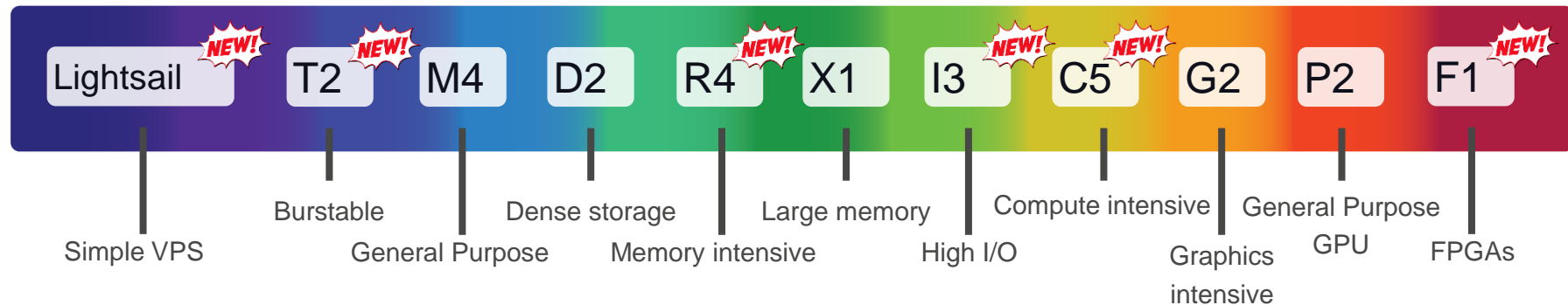
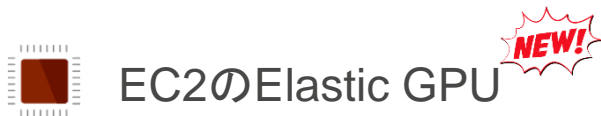


拡張



ネットワーキング

コンピュータイノベーション



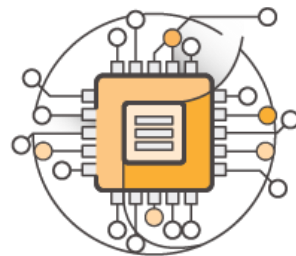
EC2インスタンスイノベーション: コンピュート最適化

大量の演算を行うワークロードに最適

- トラフィック量の多いフロントエンド群、MMOゲーミング、メディア処理、トランスコーディング、HPC (High Performance Computing) アプリケーション

C5 ... 次世代のコンピュート最適化インスタンス(まもなく登場！)

- 最新のIntel Skylakeプロセッサ
- 新しい高度なベクトル拡張命令セットであるAVX-512をサポート
ピークパフォーマンス(AVX2)では、クロックサイクルあたりのFLOPSが最大で2倍
- 最大サイズ。c5.18XLでは、72個のvCPU、144GiBのメモリ、20Gbpsの専用ネットワーク帯域幅



EC2 インスタンス・イノベーション: メモリーの最適化

R4 – 2017年12月出荷 – GiB と vCPUの比率は8:1

- より大きな新しいインスタンスのサイズ (r4.16XL)、64 vCPU 、 RAM 488 GiB
- Intel E5 v4 Broadwell プロセッサ(AVX2) 、TSX
- 改良された高性能メモリ、DDR4 メモリ
- ネットワーク帯域最大 20 Gbps

X1 – GiBあたりメモリ最大、価格は最安値 – GiB とvCPU の比率は16:1

- 2TB RAM/128 vCPU 、 1TB RAM/64 vCPU
- Intel E7 v3 Haswell 4-socket CPU (より高速な QPI 速度)
- SAP HANA インメモリデータベース / アナリティクス、シミュレーション、レンダリング

今年後半: X1E インスタンス 、4 TB の RAM

2018年までのロードマップ: 8TB と 16TB のメモリインスタンス!

EC2インスタンスイノベーション:ストレージ最適化インスタンス

I3・・・2017年2月に登場した次世代のHigh I/Oインスタンス

- Intel E5 v4 Broadwellプロセッサ (AVX2、TSX)
- VME (Non-Volatile Memory Express) ベースのSSD、ランダムリードは最大330万 IOPS、シーケンシャルリードのトータルスループットは16GB/s
- 1 × 475GBのNVME SSDを搭載したi3.largeから、8 × 1.9TBのNVME SSDを搭載したi3.16xlargeまで、サイズは6種類
- トランザクション対応のワークロード、ハイパフォーマンスデータベース、リアルタイム分析、NoSQLデータベースに最適

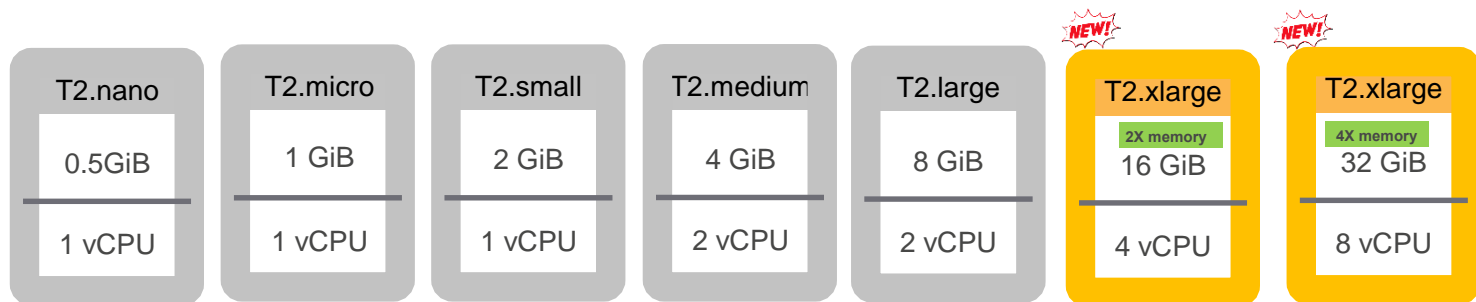
I2に比べて、IOPSが9倍、ストレージが2.3倍、メモリが2倍、vCPUが2倍！

EC2インスタンスイノベーション: バースタブルインスタンス

T2バースタブルパフォーマンスインスタンスは、ベースラインレベルのCPUパフォーマンスを提供する一方で、ベースラインを超える需要にも対応できる機能を提供

t2.xlargeとt2.2xlargeは2016年12月から提供開始

- 最大8個のvCPUと32GiBのメモリ
- 開発環境、データベース、アプリケーションサーバー、Webサーバーに最適



大量の演算を実行するワークロード

- CPUでのスケーリング
 - バッチジョブ: スポットインスタンス
- もっとうまく行うことは可能？
 - ワークロードによっては、事実上、CPUだけで実行することは不可能 - 数週間かかる
 - 実行時の遅延の削減
 - パフォーマンスとコスト最適化

ハードウェアアクセラレーションとは何か？

- 一部の機能をCPUで実行されるソフトウェアよりも効率よく実行するための特別なハードウェア(ハードウェアアクセラレータ)を使用



ハードウェアアクセラレーションとは何か？

CPUがスイスアーミーナイフのようなものだとならば...



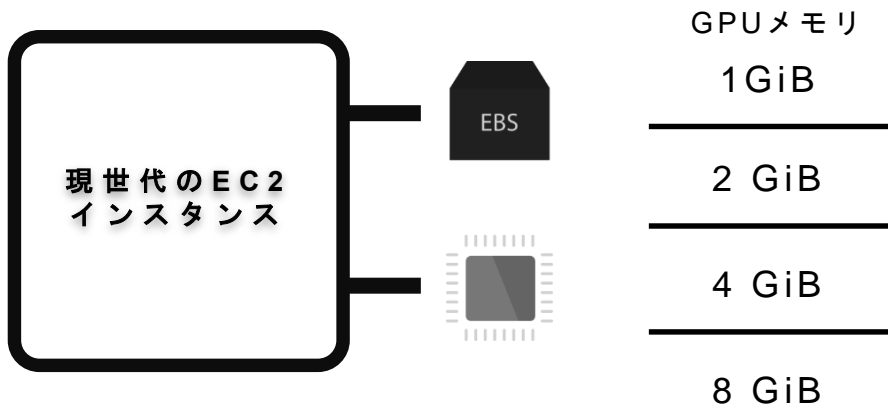
ハードウェアアクセラレータはエッグスライサーのようなもの



グラフィックアクセラレーション

Elastic GPU - プレビュー

- Amazon EC2インスタンスに安価なグラフィックアクセラレーションをネットワーク経由で追加することが可能
- 幅広いサイズ展開。GPUをさまざまなEC2インスタンスにアタッチすることで、最適なパフォーマンスを実現
- どのようなグラフィックスアプリケーションでも実行できる自信につながるOpenGLへの準拠



GPUによる高速コンピューティング

- ユビキタス
- 高度なデータ並列処理
- 浮動小数点演算の割合が高い
- 一貫性の高い優れたAPIドキュメント(CUDA、OpenGL)
- 幅広いISVとオープンソースフレームワークによるサポート

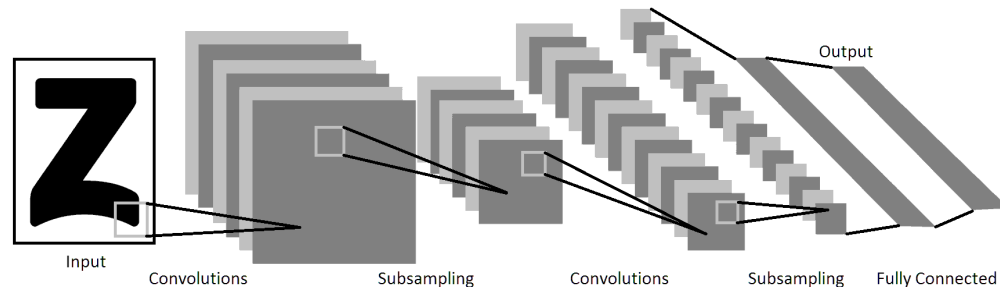
GPUコンピュータインスタンス

P2 ... GPGPUインスタンスとDeep Learning AMI - 9月から提供

- 16個のNVIDIA Tesla K80 GPUと192GBのGPUメモリ
- 完全なGPUDirect P2P機能
- 40,000個のCUDAコア、70TFLOPSの単精度浮動小数点演算性能、23TFLOPSを超える倍精度浮動小数点演算性能
- 機械学習、数値流体力学、金融工学、地震解析、分子モデル構築、ゲノミクス、レンダリングのための並列処理を実現

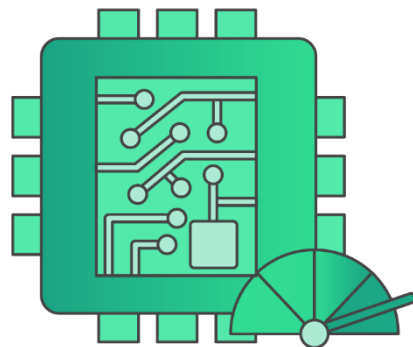
次世代のGPGPUインスタンス

- NVIDIAとAWSはクラウドでのGPUアクセラレーションを最適化するために数年前から緊密な共同作業を行っており、NVIDIA Volta GPUのローンチパートナーになることをとても楽しみにしている
- NVIDIAとAWSの次のGPGPUインスタンスファミリーは、今年後半にVoltaの提供が開始された時点で、Voltaベースとなる



FPGAによる高速コンピューティング

- 特別なアルゴリズムのためのカスタムハードウェア
- 標準以外のデータ構造のサポート
- フィールドリプログラマビリティに基づくより容易なメンテナンス
- データフロープログラミング
- スレッド間の依存性が高いアプリケーション向き
- 大きなローカルメモリと高いメモリ帯域幅を提供
- コスト効率

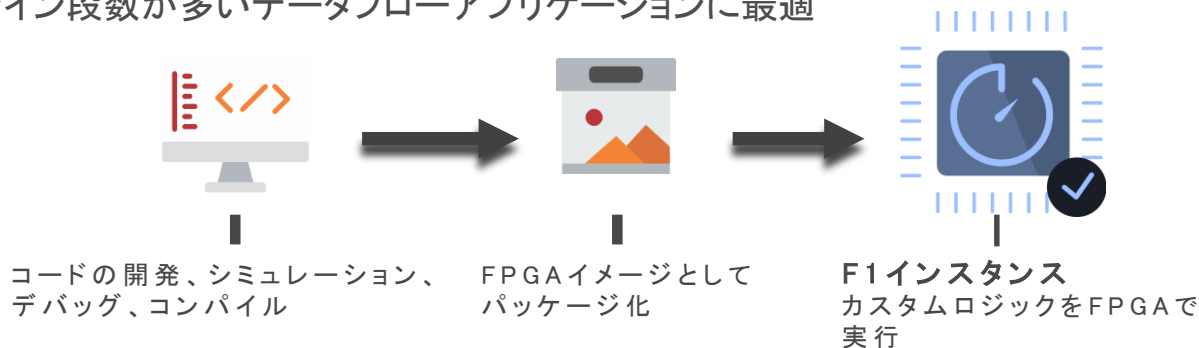


Amazon FPGA ImageとF1インスタンス

F1は顧客によるプログラミングが可能なアプリケーションアクセラレーションのためのFPGAハードウェアを初めて搭載したコンピュータインスタンス

FPGAイノベーションの開発とデプロイを容易にするAmazon FPGA Image

- F1インスタンスから高性能なFPGAへの専用アクセスにより、最大30倍までアプリケーションを高速化
- HDK（Hardware Developer Kit）と開発者AMIにより、開発時間を大幅に短縮
- AWS Marketplaceとの統合により、100万人以上の顧客にFPGAイノベーションを提供
- トランスコーディング、金融リスクモデリング、ゲノム解析、ビッグデータ処理、大規模なシミュレーションを含め、パイプライン段数が多いデータフローアプリケーションに最適



インスタンスの拡張機能: ネットワーク機能

Elastic Network Adapter(ENA)・・・X1と新世代のインスタンスで提供

- スケーラビリティの改善、高いスループットとpps(packet per second)パフォーマンス、一貫した低遅延を目的として、Amazonによって構築されたカスタムネットワークドライバ

ネットワークパフォーマンスの改善

- より大きなインスタンスでのスループット・・・X1、P2、M4、R4、I3、C5では20Gbps
- より小さなインスタンスでのスループット・・・より小さなインスタンス(R4、I3、C5、およびその他のインスタンス)でのピーク帯域幅は10Gbps

IPv6

- PCにIPv6 CIDRブロックを関連付けることで、VPCのEC2インスタンスでIPv6アドレスを使用することが可能



インスタンスの拡張機能：EBSストレージ

新しいスループット最適化HDDボリューム

- ST1 … 最大スループットは500MB/s、ベースラインは40MB/s (\$0.045/GB)
- SC1 … 最大スループットは250MB/s、ベースラインは12MB/s (\$0.025/GB)

パフォーマンスの改善

- PIOPS … IOPSとGBの比率を30:1から50:1に改善

EBSの暗号化とカスタムキー

- AWSのリージョンおよびアカウントの間で暗号化されたスナップショットをコピー
- 暗号化されたブートボリューム

EC2のコストを最適化

Amazon EC2の購入オプション

オンデマンド

コンピュータキャパシティを
時間単位で支払い、
長期契約なし

スパイキーなワークロードに、
またはニーズを定義するために



リザーブド

契約期間は1年または3年、
オンデマンドと比べて
大幅な割引

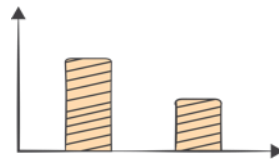
確定済みのワークロードに、
またはベースラインでの使用に



スポット

未使用のコンピュータキャパシティ
に対して市場価格で支払い、
オンデマンドと比べてかなり割引

フォールトトレラントなワークロード、
時間的に余裕のあるワークロード、
または一過性のワークロードに



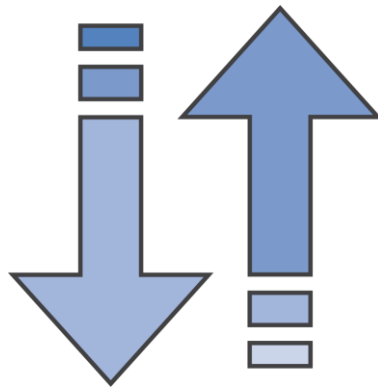
EC2オンデマンドインスタンスの料金



低コストで柔軟



開発とテスト

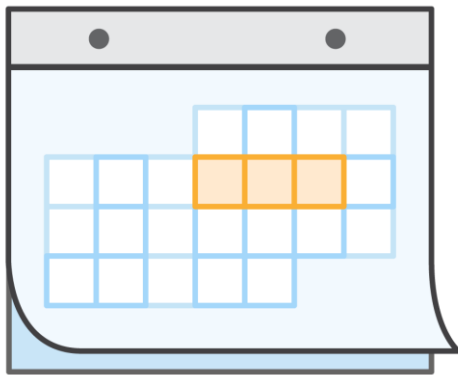


短期間、スパイキー、
予測不能

EC2リザーブドインスタンスの料金



定常的



必要に応じてキャパシティを
予約



前払いによる
コスト削減

リザーブドインスタンスの概要

ニーズに最適なRIオプションを決定

節約の可能性

インスタンスファミリー、
OS、
テナンシーの変更

リース期間

AZ、
インスタンスサイズ (Linux)、
ネットワークタイプの変更

支払い

リージョンに関する特典

スタンダード

最大75%

なし

1年または3年

あり

前払いなし
一部前払い
全額前払い

あり

コンバーティブル

最大66%

あり

3年のみ

あり

前払いなし
一部前払い
全額前払い

あり

コンバーティブルリザーブドインスタンス(RI)

コンバーティブルリザーブドインスタンスでは以下が可能:

新しいインスタンスファミリに交換: R3→C3→T2→M4など

新しいインスタンス料金に交換: AWSがインスタンスの定価を値下げした場合

新しいOSに交換: WindowsからLinuxへの交換など

新しいインスタンスサイズに交換

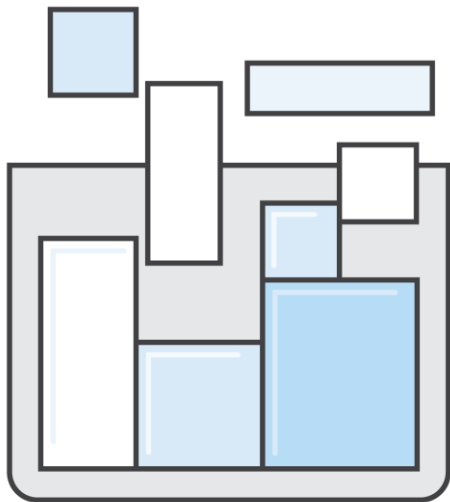
テナンシーを交換: ハードウェア占有(dedicated)インスタンスからデフォルトへの交換など

別の支払いオプションに交換: 「前払いなし」から「一部前払い」への交換など

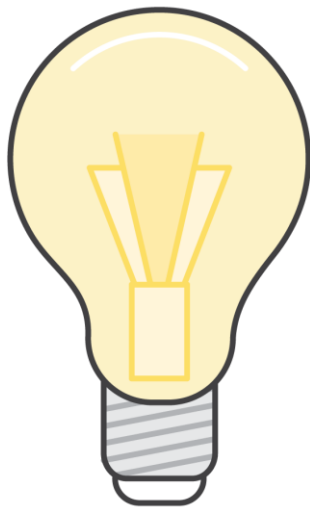
リージョンに関する特典の概要

	キャパシティ予約の特典？	AZ間の自動的な割引？	インスタンスサイズ間の自動的な割引？	RI Marketplaceでの販売？
ゾーナル	あり	なし	なし	あり
リージョナル	なし	あり	あり	なし

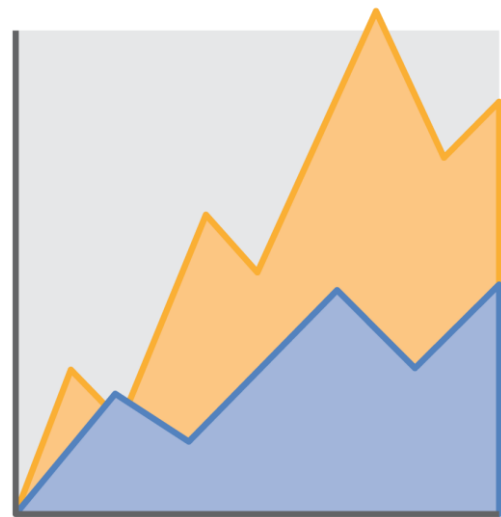
EC2スポットインスタンスの料金



時間または
インスタンスが柔軟



コスト重視のビジネスを
実験または構築



緊急性の高いコンピューティング
ニーズを持つ、または追加の
キャパシティを大量に必要として
いるユーザー

スポットインスタンスの詳細

90%節約!*

オプション

- インスタンスの可用性を維持するためのスポットフリート
- 継続的に実行しなければならないワークロードに対するスポットブロックの継続期間(1～6時間)

コミットレベル

- なし

* 特定のEC2インスタンスタイプ、リージョン、AZに基づくオンデマンド価格との比較

スポットインスタンスのルール

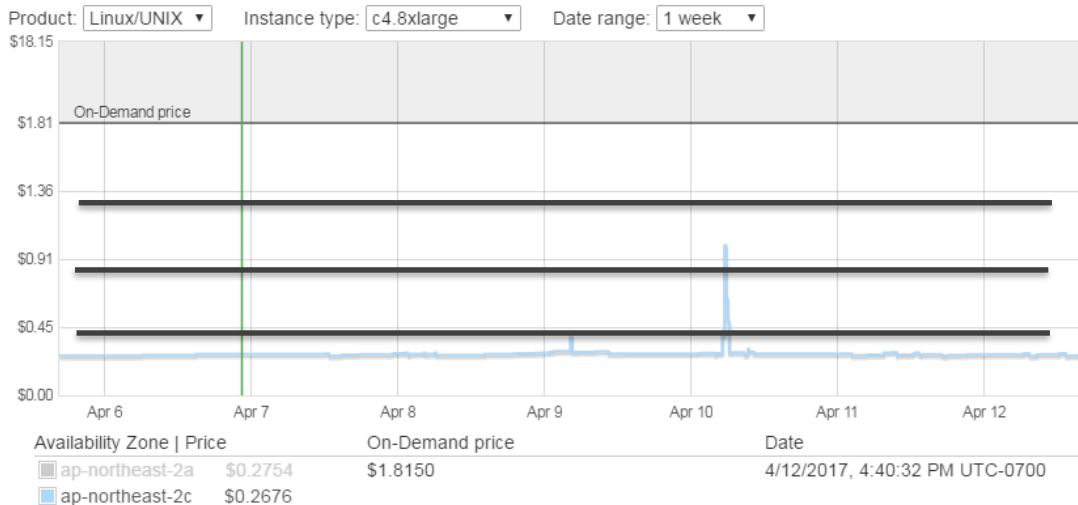


コンピュータの市場価格は需要と供給に基づいて変動



入札価格を超える支払いは発生しない

Spot Instance Pricing History



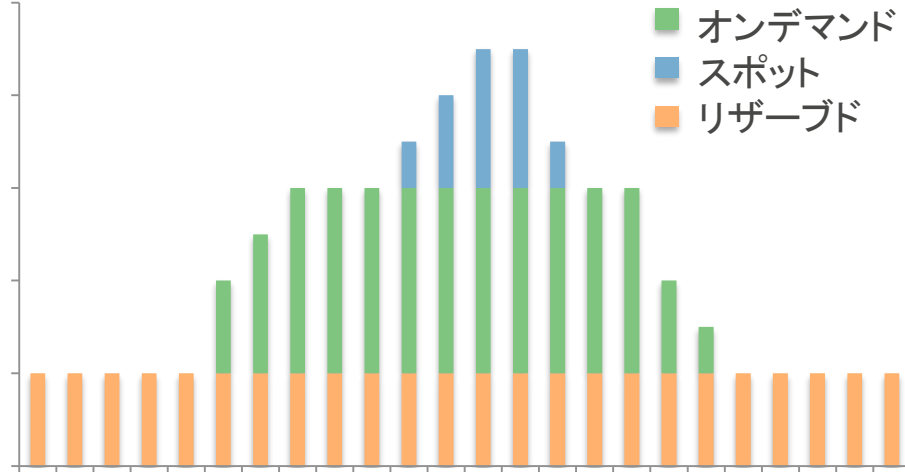
オンデマンドの75%
で入札
オンデマンドの25%
で入札

オンデマンドの
50%で入札

市場価格の85%
割引で支払い！

購入モデルを組み合わせて使用

1. 既知の定常的なワークロードにはリザーブドインスタンスを使用
2. 複数のAuto Scalingグループをセットアップ
3. スポット、オンデマンド、または両方を使ってスケーリング



単純なワークロードの要件はたいてい単純



Webサイト



ブログ



開発環境



プロトタイピング



ビルドサーバー



Amazon Lightsail : AWSでの取り組み を開始するための最も簡単な方法



VPS (Virtual Private Server)



永続ストレージ



ネットワーキング

事前に設定されたインスタンスイメージを選択

オペレーティング
システム



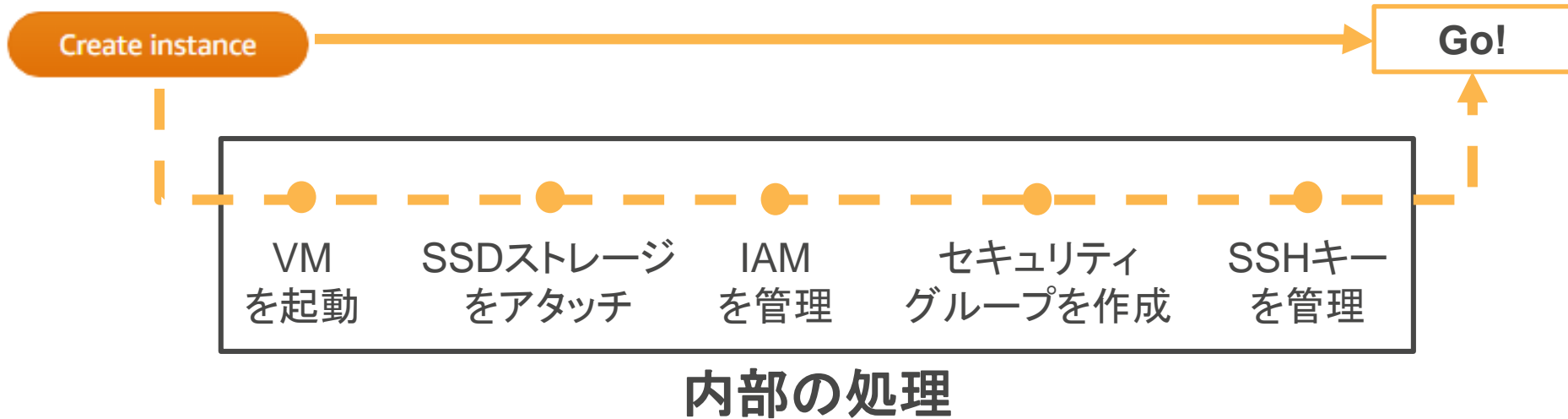
アプリケーション



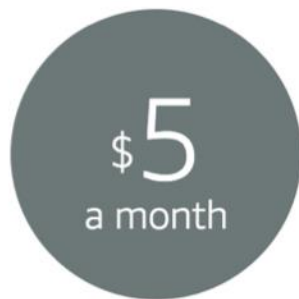
開発スタック



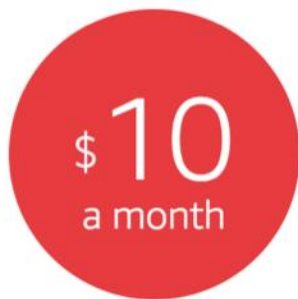
VPSインスタンスをワンクリックで起動



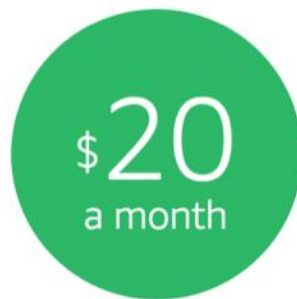
必要なものをどれでも予測可能な低価格で



512MB Memory
1 Core Processor
20GB SSD Disk
1TB Transfer



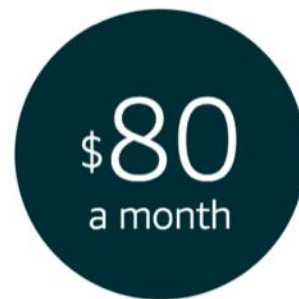
1GB Memory
1 Core Processor
30GB SSD Disk
2TB Transfer



2GB Memory
1 Core Processor
40GB SSD Disk
3TB Transfer

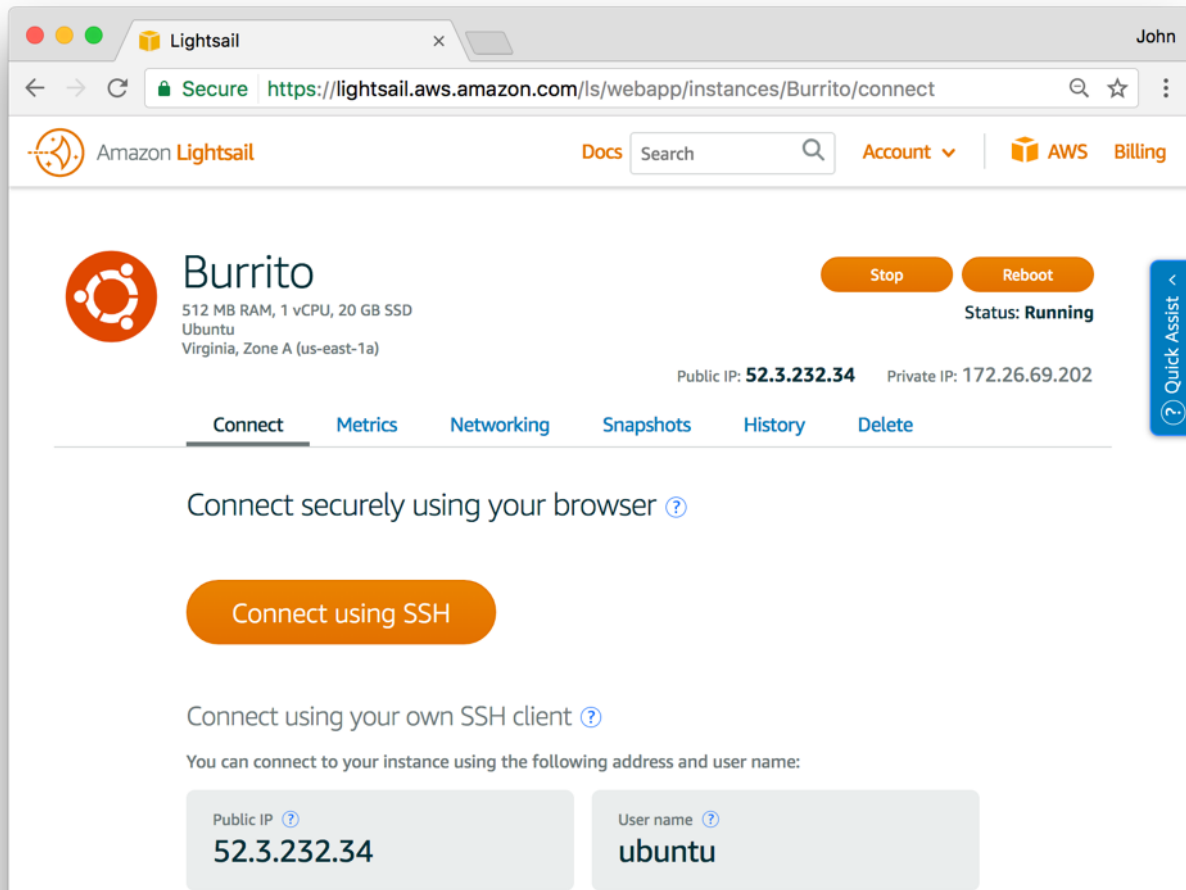


4GB Memory
2 Core Processor
60GB SSD Disk
4TB Transfer



8GB Memory
2 Core Processor
80GB SSD Disk
5TB Transfer

使いやすいインターフェイス



The screenshot displays the Amazon Lightsail console interface for a specific instance named 'Burrito'. The browser window shows the URL 'https://lightsail.aws.amazon.com/ls/webapp/instances/Burrito/connect'. The page header includes the 'Amazon Lightsail' logo, navigation links for 'Docs', 'Search', 'Account', 'AWS', and 'Billing', and a user profile 'John'.

The main content area features the 'Burrito' instance details: an Ubuntu logo, the instance name 'Burrito', and specifications '512 MB RAM, 1 vCPU, 20 GB SSD'. The status is 'Running'. It also shows the public IP '52.3.232.34' and private IP '172.26.69.202'. Action buttons for 'Stop' and 'Reboot' are visible.

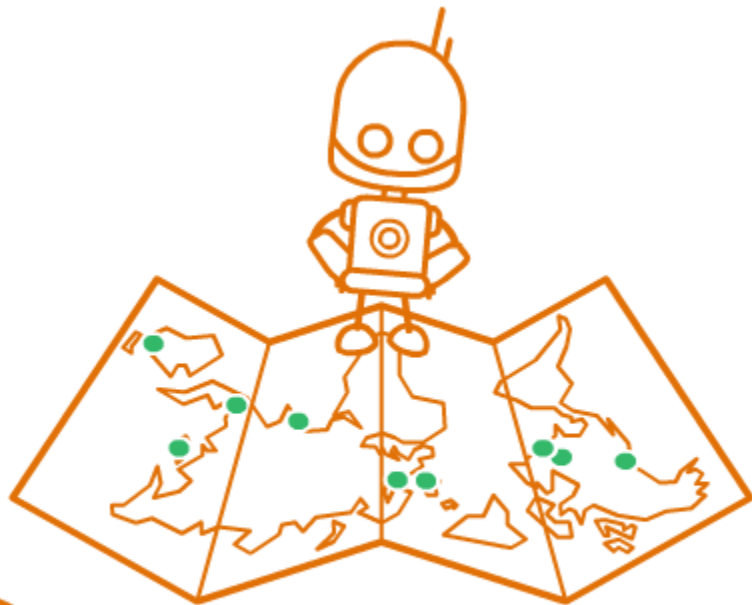
Below the instance details is a horizontal menu with tabs: 'Connect', 'Metrics', 'Networking', 'Snapshots', 'History', and 'Delete'. The 'Connect' tab is selected.

The 'Connect' section provides instructions: 'Connect securely using your browser' and 'Connect using SSH'. A large orange button labeled 'Connect using SSH' is prominently displayed.

Further down, it says 'Connect using your own SSH client' and provides the connection details: 'You can connect to your instance using the following address and user name:'.

At the bottom, two input fields are shown: 'Public IP' with the value '52.3.232.34' and 'User name' with the value 'ubuntu'.

世界中をライトセーリング！



Amazon**Lightsail**

まとめ

AWSのインスタンスロードマップを決定するのは顧客

- AWSは実行すべき顧客のジョブとパフォーマンスの定義を理解している
- この作業をより効果的に行うために新しいハードウェアとソフトウェアを調査している
- パフォーマンスの改善を可能にする次世代のインスタンスと、顧客の新しいニーズを解決するための新しいインスタンスファミリとインスタンス機能を提供する

AWSでは、新しいインスタンス、機能、コンピュートモデルを使ってすべてのワークロードをカバーするために、コンピュートサービスのイノベーションに取り組んでいる

AWS

S U M M I T

