

Alibaba Cloud E-MapReduce

クイックスタート

Document Version20200420

目次

1	構成の選択.....	1
2	前提条件.....	4
3	クラスターの作成.....	6
4	E-MapReduce クイックスタート.....	12

1 構成の選択

E-MapReduce (EMR) を使用するには、適切な Hadoop クラスターを選択する必要があります。EMR の構成を選択する際、企業でのビッグデータの使用を考慮し、総データ量、サービスの信頼性および予算を見積もります。

ビッグデータシナリオ

EMR は以下の企業ビッグデータシナリオに適用されています：

- **バッチ操作** このシナリオは高いディスク I/O スループットと高ネットワークスループットを必要としますが、リアルタイムキャパシティ要件は低くなります。大量のデータを処理する必要があるがリアルタイム処理は必要ない場合、MapReduce、Pig および Spark を使用できます。このシナリオは高メモリキャパシティを必要としません。ゆえに、データシャッフリングを実行する際は CPU、メモリ、ネットワーク要求にフォーカスする必要があります。
- **アドホッククエリ** データサイエンティストおよびデータアナリストはアドホッククエリを使用してデータを収集します。このシナリオではリアルタイムクエリ、高いディスク I/O スループットとネットワークスループットが要求されます。このシナリオでは Impala や Presto が使用できます。このシナリオには高メモリ要件もあります。総データ量と同時クエリを考慮する必要があります。
- **ストリームコンピューティング、高ネットワークスループットおよび計算集約型シナリオ** これらのシナリオには、Flink、Spark Streaming または Storm を使用できます。
- **メッセージキュー** このシナリオは高いディスク I/O スループット、ネットワークスループットが必要で、大量のメモリを消費し、ストレージは HDFS に依存しません。したがって、EMR 内で Kafka クラスターが使用できます。EMR クラスターは Hadoop への影響を避けるために Kafka クラスターと Hadoop クラスターに分割されます。
- **コールドバックアップ** このシナリオは高い I/O スループットまたはコンピューティングスループットを要求せず、低コストです。EMR d1 インスタンスをコールドバックアップに使用することを推奨します。d1 インスタンスのストレージコストは 0.03 元 / 月 / GB です。

EMR インスタンス

3 種類のインスタンスで構成された EMR クラスター： マスターインスタンス、コアインスタンス、タスクインスタンス

EMR ストレージとしてウルトラディスク、SSD ディスクおよびローカルディスクを選択できます。異なるディスクのパフォーマンスは、SSD ディスク > ローカルディスク > ウルトラディスク、の順です。

EMR 基盤ストレージは OSS (標準 OSS のみ) と HDFS をサポートします。OSS は HDFS よりも高いデータ可用性を持っています。HDFS のデータ可用性がクラウドディスクまたはローカルディスクストレージの信頼性に依存するのに対し、OSS のデータ可用性は 99.99999999% です。

ストレージ価格は以下のとおりです：

- ローカルディスク付きインスタンス：0.003 USD/GB / 月
- OSS 標準ストレージ：0.02 USD/GB / 月
- ウルトラディスクストレージ：0.05 USD/GB / 月
- SSD ディスクストレージ：0.143 USD/GB / 月

EMR 設定の選択

- マスターインスタンス設定の選択
 - マスターインスタンスは、NameNode や ResourceManager のような Hadoop のマスタープロセスをデプロイするために使用されます。
 - 本番クラスターの高可用性を有効化することを推奨します。HDFS、YARN、Hive または HBase のような EMR コンポーネントに高可用性は利用可能です。クラスター作成時の「クラスター設定」ステップで高可用性を有効化することを推奨します。EMR クラスター作成時に高可用性が有効化されていない場合、後から有効化することはできません。
 - マスターインスタンスは HDFS メタデータとコンポーネントログファイルを格納するために使用されます。それらは低いディスク I/O 要求を持つ計算集約型です。HDFS メタデータはメモリに格納されます。最小推奨メモリサイズはファイル数に基づいて 16 GB です。
- コアインスタンス設定の選択
 - コアインスタンスはデータの格納、コンピューティングタスクおよび DataNode や NodeManager のようなプロセスの実行に使用されます。
 - HDFS (3 つのバックアップ) に格納された総データ量が 60 TB を超えると、ローカルディスク付きのインスタンス (ECS d1 インスタンスおよび ECS d1ne インスタンス) を使用することを推奨します。ローカルディスク容量は以下のとおり計算されます： $(\text{CPU cores 数}/2) * 5.5 \text{ TB} * \text{インスタンス数}$ たとえば、4 つの 8 コア d1 インスタンスを購入する場合、ローカルディスク容量は $8/2 * 5.5 * 4 = 88 \text{ TB}$ となります。HDFS は 3 つのバックアップが必要です。したがって、少なくともローカルディスクを使用する 3 つのインスタンスを購入する必要があります。データ信頼性とディスクリカバリのためには少なくとも 4 インスタンスを購入することを推奨します。
 - HDFS に格納された総データ量が 60 TB より少ない場合、ウルトラディスクまたは SSD ディスクを使用できます。

- タスクインスタンス設定の選択

- タスクインスタンスはコアインスタンスの CPU とメモリが十分なコンピューティング容量を持たない場合に使用されます。タスクインスタンスはデータの格納や DataNode の実行を行いません。CPU とメモリ要件に基づいてインスタンスの数を見積もることができます。

EMR ライフサイクル

EMR は素早い[クラスターのスケールアップ](#)を可能にするオートスケーリングをサポートします。クラスターノードの設定を柔軟に調整でき[ECS インスタンスの設定をアップグレードまたはダウングレード](#)できます。

ゾーンの選択

EMR および事業システムを[同一ゾーン](#)および[同一リージョン](#)にデプロイすることを推奨します。

2 前提条件

E-MapReduce クラスターを作成する前に、次の前提条件が満たされていることを確認します。

1. Alibaba Cloud アカウントの作成

E-MapReduce クラスターを作成するためには、Alibaba Cloud エコシステムで内でユーザーを一意に識別するための Alibaba Cloud アカウントが必要です。このアカウントを使用して、E-MapReduce クラスターを作成し、[OSS \(Object Storage Service\)](#) や [RDS \(ApsaraDB for RDS\)](#) を含む、他の Alibaba Cloud サービスを有効化します。

Alibaba Cloud アカウントの作成についてのより詳しい情報は、[Alibaba Cloud へのサインアップ](#) をご参照ください。

2. AccessKey の作成 (任意)

E-MapReduce を使用するには、最低でも 1 つの AccessKey を作成する必要があります。次の手順に従って、AccessKey を作成します。

- a. [Alibaba Cloud ウェブサイト](#) にログインします。
- b. Alibaba Cloud コンソールへ移動します。
- c. プロフィール画像の上にカーソルを置き、**[AccessKeys]** をクリックします。



:

次のダイアログボックスが表示されたら、**[Continue to manage AccessKey]** ボタンをクリックします。

Security Tips



AccessKey of your cloud account is the secret key to access Alibaba Cloud APIs. Since the AccessKey has full permissions of your cloud account, please make sure you keep it well. To avoid the AccessKey being used by others to cause [Sensitive information leakage](#), do not release your AccessKey to any external channels (for example, Github). We strongly recommend you use the AccessKeys of RAM users in API calls, according to [Alibaba Cloud account security best practices](#).

Continue to manage AccessKey

Get Started with Sub Users's AccessKey

- d. **[Create AccessKey]** をクリックします。

- e. AccessKey が作成されます。

3. Alibaba Cloud OSS の有効化

E-MapReduce では、ジョブと実行ログが Alibaba Cloud OSS に保存されます。したがって、Alibaba Cloud OSS を有効化する必要があります。詳しい情報については、[OSSへのサインアップ](#)をご参照ください。作成すべき EMR クラスターと同じリージョンに OSS バケットを作成します。詳しい情報については、[バケットの作成](#)をご参照ください。

4. 高性能な構成のインスタンスの作成 (任意)

Pay-As-You-Go-based クラスター用に 8 個以上のコアを持つインスタンスを作成する必要がある場合は、Alibaba Cloud アカウントにログインし、アプリケーションのチケットを送信します。[サポートとサービス](#)

3 クラスターの作成

このセクションでは、E-MapReduce クラスターの作成および設定方法を説明します。

クラスター作成ページに移動します。

1. [\[Alibaba Cloud E-MapReduce コンソール\]](#) にログインします。
2. RAM 認証を完了します。詳細は、「[#unique_6](#)」をご参照ください。
3. クラスターを作成するリージョンを選択します。リージョンは、クラスター作成後に変更することはできません。
4. 右上の **[クラスター作成]** をクリックします。

クラスターの作成

操作手順は下記のとおりです。

- ソフトウェアの設定
- ハードウェアの設定
- 基本設定

手順 1：ソフトウェアの設定

説明

- **EMR** バージョン：デフォルトでは最新のバージョンが選択されます。
- クラスタータイプ：現在、E-MapReduce (EMR) は以下のクラスタータイプを提供しています：
 - Hadoop クラスター これらのクラスターは、Hadoop、Hive、Spark、Spark Streaming、Flink、Storm、Presto、Impala、Oozie、または Pig といった複数のエコシステムコンポーネントを提供しています。Hadoop、Hive および Spark はセミホスティングサービスで分散型大規模データストレージとコンピューティングに使用されます。Spark Streaming、Flink および Storm はストリームコンピューティングを提供します。Presto と Impala はインタラクティブなクエリを実現するために使用されます。コンポーネントの詳細については、クラスターとサービス管理ページに表示されている「サービスリスト」をご参照ください。
 - Kafka クラスター これらのクラスターは高スループットと高スケーラビリティを備えたセミホスティング分散メッセージシステムを提供します。Kafka クラスターはクラスターの安定稼動を維持する包括的なサービス監視システムを提供します。Kafka クラスターはよりプロフェッショナルで、信頼性があり、安全です。これらのクラスターをデプロイした

り維持する必要はありません。これらのクラスターは一般的にログ収集や監視データ集約のようなシナリオで使用されます。オフラインデータ処理とストリームコンピューティング、およびリアルタイムデータ分析もサポートされています。

- Druid クラスター これらのクラスターはセミホスティングおよびリアルタイムインタラクティブ分析サービスを提供します。Druid クラスターはミリ秒単位での大量のデータ照会および複数の方法によるデータ書き込みをサポートしています。Druid クラスターは EMR Hadoop、EMR Spark、OSS、RDS のような他のサービスと共に柔軟で安定したリアルタイム照会を提供します。
- データサイエンスクラスター これらのクラスターは一般的にビッグデータおよび AI シナリオに適用されます。データサイエンスクラスターは Hive および Spark のオフラインビッグデータ ETL と TensorFlow モデルトレーニングを提供します。CPU + GPU の混成コンピューティングフレームワークと NVIDIA CPU にサポートされるディープラーニングアルゴリズムを選択しコンピューティングタスクを効率的に実行することができます。
- 必要サービス： デフォルト設定を選択します。管理ページ上で後からサービスの追加、有効化および無効化ができます。
- ハイセキュリティモード： クラスターに Kerberos 認証機能を有効化するかを示します。一般ユーザー用に作成されたクラスターは一般的にこのモードを必須としないため、このモードはデフォルトで無効化されています。
- カスタム設定の有効化： JSON ファイルを指定しソフトウェア設定を変更できます。手順の詳細については、[\[ソフトウェア設定\]](#) をご参照ください。

手順 2： ハードウェア設定

説明

- 課金設定
 - 課金設定： クラスターのテスト時は従量課金を選択できます。すべてのテストに合格すると、サブスクリプションベースのクラスターを作成し使用できます。

- ネットワーク設定

- ゾーン：一般的に、デフォルトゾーンが使用されます。
- ネットワークタイプ：VPC がデフォルトで選択されます。VPC を選択しなかった場合は [VPC コンソール](#) に移動して作成します。
- ゾーン：同一リージョン内の地理的なエリア これらのゾーンは VPC を通じて相互接続されています。
- **VPC**：指定のリージョンで作成された VPC を選択します。利用可能な VPC が存在しない場合、**[VPC/VSwitch の作成]** をクリックし、現在のゾーンで作成します。
- **VSwitch**：現在のゾーン内の指定された VPC 用に VSwitch を選択します。利用可能な VSwitch が存在しない場合、VPC コンソールに移動し現在のゾーンで作成します。
- セキュリティグループ名：クラスターを初めて作成した場合は、デフォルトではセキュリティグループは利用できません。名前を入力しセキュリティグループを作成する必要があります。セキュリティグループを既に作成している場合はセキュリティグループを選択します。

- クラスター設定

- 高可用性：この機能が有効になると、リソースマネージャーとネームノードの高可用性を保証する 2 つのマスターノードが提供されます。HBase クラスターはデフォルトで高可用性をサポートしています。HBase クラスターは 2 つのマスターノードの 1 つとしてコアノードを使用する必要があります。高可用性機能が有効化されると、HBase クラスターはより安全で信頼性のある高可用性をサポートする 1 つのマスターノードのみが必要です。

高可用性をサポートするクラスターを作成する必要がある場合、テスト中に高可用性を有効化します。

- マスターノード：リソースマネージャーおよびネームノードのようなプロセスをデプロイします。
 - マスターインスタンスタイプ：ニーズに合わせてインスタンス仕様を選択します。詳細については、「[インスタンスタイプファミリー](#)」をご参照ください。
 - システムディスクタイプ：ニーズに合わせてウルトラディスクまたは SSD ディスクを選択します。
 - システムディスクサイズ：ニーズに合わせてディスクをリサイズできます。推奨最小ディスクサイズは 120 GB です。
 - データディスクタイプ：ニーズに合わせてウルトラディスクまたは SSD ディスクを選択します。
 - データディスクサイズ：ディスクをリサイズできます。推奨最小ディスクサイズは 80 GB です。
 - マスターインスタンス：デフォルトのマスターインスタンス数は 1 です。
- コアインスタンス：すべてのクラスターデータを格納します。ニーズに合わせてインスタンスのスケールアップができます。
 - コアインスタンスタイプ：ニーズに合わせてインスタンス仕様を選択できます。詳細については、「[インスタンスタイプファミリー](#)」をご参照ください。
 - システムディスクタイプ：ニーズに合わせてウルトラディスクまたは SSD ディスクを選択します。
 - システムディスクサイズ：ニーズに合わせてディスクをリサイズできます。推奨最小ディスクサイズは 80 GB です。
 - データディスクサイズ：ニーズに合わせてウルトラディスクまたは SSD ディスクを選択します。
 - データディスクサイズ：ニーズに合わせてディスクをリサイズできます。推奨最小ディスクサイズは 80 GB です。
 - コアインスタンス：デフォルトのコアインスタンス数は 2 です。ニーズに合わせてコアインスタンスの数を調整できます。
- タスクインスタンスグループ：タスクインスタンスグループにはデータは格納されません。タスクインスタンスグループはクラスターのコンピューティング容量を調整するために使用されます。この機能はデフォルトで無効になっています。ニーズに合わせて有効化できます。

手順 3：基本設定

説明

- 基本情報
 - クラスター名：クラスターの名前を入力します。1～64 文字の長さで設定でき、漢字、大文字アルファベット、小文字アルファベット、数字、ハイフン (-) およびアンダースコア (_) を使用できます。
- 実行中ログ
 - 実行中ログ：この機能を有効化または無効化できます。この機能はデフォルトで無効になっています。この機能を有効化すると、ログを保存する OSS パスを指定しなければなりません。すべての実行中ログは指定の OSS ディレクトリに格納されます。この機能を有効化するには、まず OSS をアクティベートする必要があります。アップロードされたデータはスペース使用量に基づいて課金されます。デバッグおよびトラブルシューティングに役立つこの機能を有効化することを推奨します。
 - ログパス：ログを保存する OSS パスを指定します。
 - 統合メタデータベース：Hive はクラスターから独立した統合メタデータベースを使用します。メタ情報はクラスターリリース後に削除されません。この機能を無効化することを推奨します。
- 権限設定：一般的にデフォルト設定が使用されます。
- ログイン設定
 - リモートログイン：セキュリティグループのポート 22 を開放するかどうかを示します。この機能はデフォルトで無効になっています。
 - パスワード：マスターインスタンスにログインするためのパスワードを設定します。ログインパスワードは 8～30 文字の長さで大文字、小文字、数字、エクスクラメーションマーク (!)、アットマーク (@)、シャープ (#)、ドルマーク (\$)、パーセント (%)、アンドマーク (&) およびアスタリスク (*) のような特殊文字を含めることができます。
- ブートストラップ操作 (任意)：クラスターが Hadoop を開始する前にカスタムスクリプトを実行するように設定することができます。詳細については、「ブートストラップアクション」をご参照ください。

構成リスト

構成リストで構成アイテムと料金を確認します。

作成の確認

設定が終わりすべての設定が有効であることを確認した後、**【作成】** ボタンがハイライトされます。情報を確認し、**【作成】** をクリックしクラスターを作成します。



- ・ 課金方法が従量課金の場合、クラスターはすぐに作成されます。クラスターリストページへ誘導されます。初期化中のクラスターをクラスターリストで確認できます。クラスター作成が完了するまで数分かかります。クラスターが作成されると、アイドル状態に変わります。
- ・ 課金方法がサブスクリプションの場合、オーダーが生成されます。クラスターは支払い完了後に作成されます。

作成の失敗

作成に失敗すると、**CREATE_FAILED** とクラスターリスト上に表示されます。ポインターを赤いエクスクラメーションマーク (!) 上に動かし、理由を表示します。

これらのクラスターにコンピューティングリソースは作成されていないため、作成に失敗したクラスターを扱う必要はありません。これらのクラスターは3日後に自動的に非表示になります。

4 E-MapReduce クイックスタート

このセクションでは E-MapReduce (EMR) 上のクラスターとジョブ、およびそれらの使用方法を説明します。たとえば、Spark ジョブを作成し、それをクラスター上で実行して Pi (π) を計算し、コンソール上で結果を確認できます。



注記：すべての [前提条件](#) が満たされていることを確認することをご確認ください。

1. クラスターの作成

- a. [Alibaba Cloud E-MapReduce コンソール](#)で、**[クラスター]** タブをクリックし、クラスターリストページへ移動します。右上の **[プロジェクトの作成]** をクリックします。
- b. ソフトウェアの設定
 - A. 最新の EMR バージョンを選択します。例：**EMR 3.13.0**。
 - B. デフォルトのソフトウェア設定を選択します。
- c. ハードウェアの設定
 - A. **[重量課金]** を選択します。
 - B. セキュリティグループが作成されていない場合、名前を入力してから作成します。
 - C. 4 コア、8 GB メモリのマスターインスタンスを選択します。
 - D. 4 コア、8 GB メモリの 2 つのコアインスタンスを選択します。
 - E. 残りの設定はすべてデフォルトの設定を使用します。
- d. 基本設定
 - A. クラスター名を入力します。
 - B. ジョブのログを保存するパスを指定します。実行中のログ機能が有効になっていることを確認してください。クラスターが作成されたリージョンで、[OSS バケット](#)を作成します。
 - C. クラスターにログインするためのパスワードを入力します。
- e. **[OK]** をクリックしてクラスターを作成します。

2. ジョブを作成します。

- a. **[データプラットフォーム]** タブをクリックし、プロジェクトリストページへ移動します。右上の **[新規プロジェクト]** をクリックします。
- b. **[新規プロジェクト]** ダイアログボックス内で、プロジェクト名および説明を入力し、**[作成]** をクリックします。
- c. 指定のプロジェクトの右にある **[ワークフローのデザイン]** をクリックし、**[ジョブの編集]** ページへ移動します。
- d. ジョブの編集ページの左側で、操作するフォルダを右クリックし **[新規ジョブ]** をクリックします。
- e. 名前と説明を入力します。
- f. ジョブタイプに **Spark** を選択します。
- g. **[OK]** をクリックします。



注：

フォルダを右クリックし、サブフォルダーの作成を選択し、フォルダをリネームしたりフォルダを削除することもできます。

- h. 以下のようにパラメーターを入力します：

```
--class org.apache.spark.examples.SparkPi --master yarn-client --driver-memory 512m --num-executors 1 --executor-memory 1g --executor-cores 2 /usr/lib/spark-current/examples/jars/spark-examples_2.11-2.1.1.jar 10
```



注：

/usr/lib/spark-current/examples/jars/spark-examples_2.11-2.1.1.jar JAR ファイル名はクラスター上の Spark のバージョンにより定義されます。たとえば、Spark のバージョンが 2.1.1 の場合、JAR ファイル名は spark-examples_2.11-2.1.1.jar となります。Spark のバージョンが 2.2.0 の場合、JAR ファイル名は spark-examples_2.11-2.2.0.jar です。

- i. **[実行]** をクリックします。

3. ジョブログエントリを表示し、結果を確認します。

ジョブを実行した後、ページ下部にある **[ログ]** タブをクリックし実行中のジョブのログを確認します。 **[詳細表示]** をクリックし、詳細ページへ移動します。このページでは、ジョブ投入ログや YARN コンテナログなどの詳細を確認できます。