

このコンテンツは公開から3年以上経過しており内容が古い可能性があります
最新情報については[サービス別資料](#)もしくはサービスのドキュメントをご確認ください

AWS Data Pipeline

AWS Black Belt Tech Webinar 2015

Yuta Imai

Solutions Architect, Amazon Data Services Japan

Archived

アジェンダ

1. What is AWS Data Pipeline
2. When not to use
3. Case studies

アジェンダ

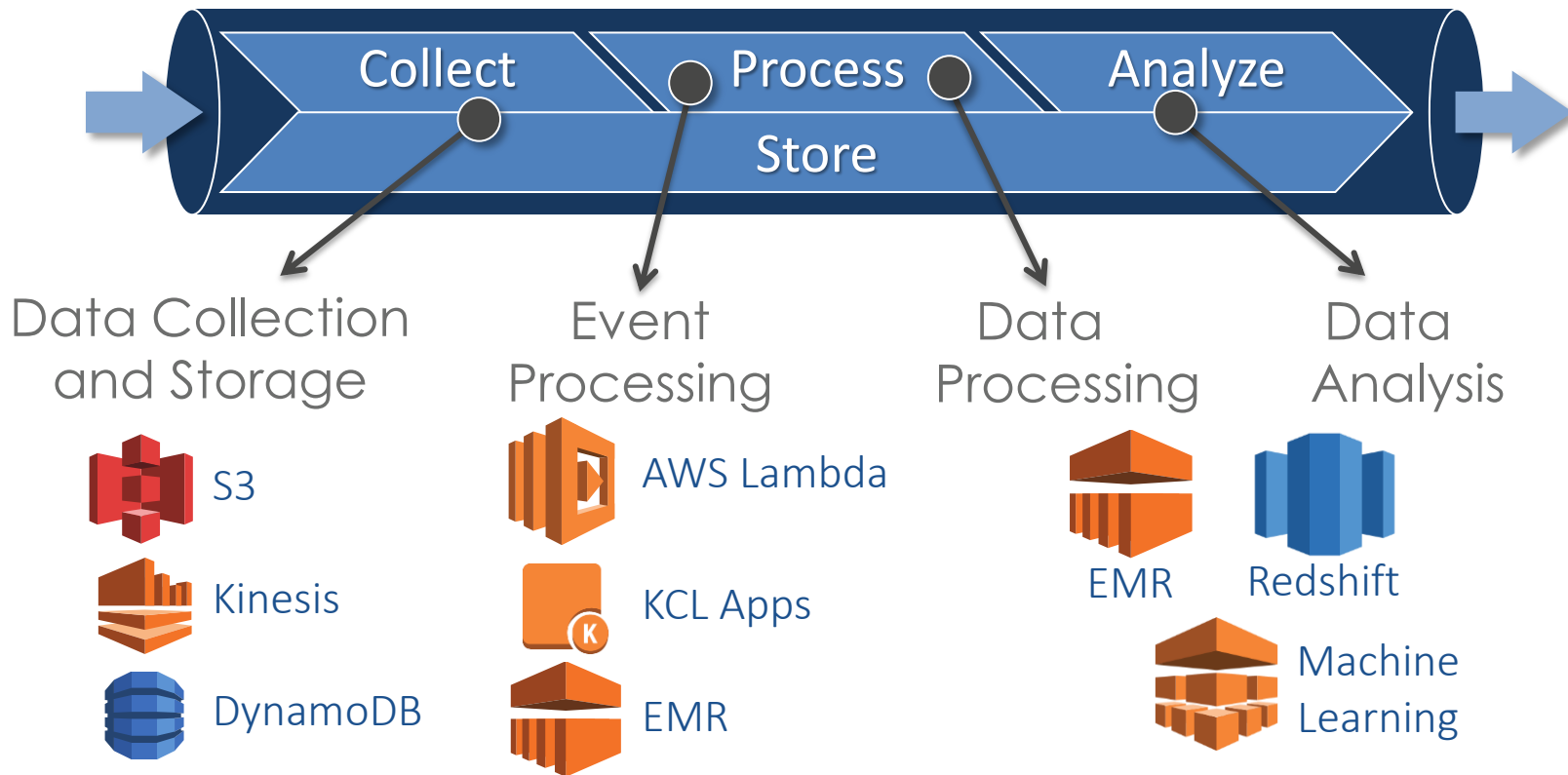
1. What is AWS Data Pipeline
2. When not to use
3. Case studies

AWS Data Pipelineでできること

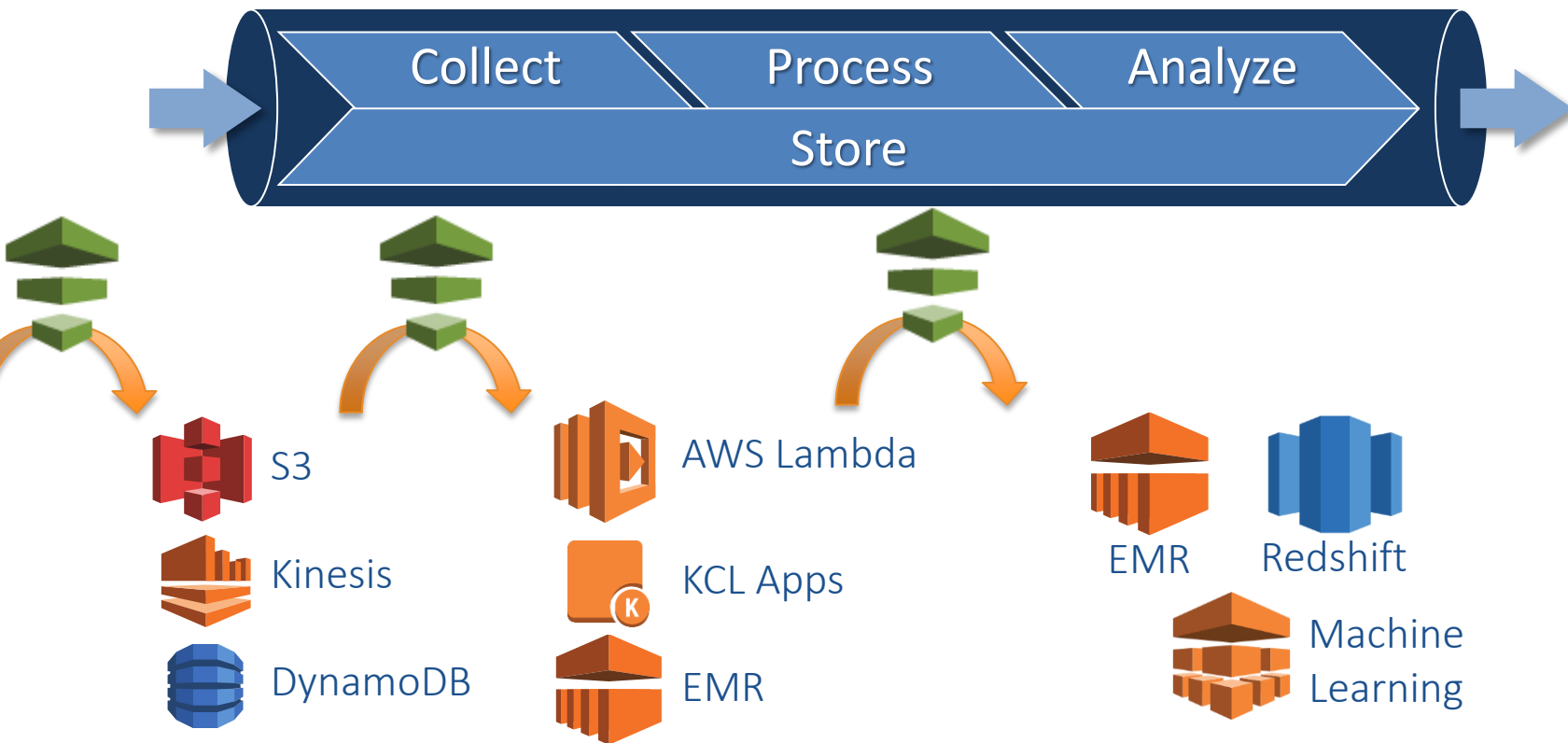
- さまざまなサーバーやデータベースからのデータの回収
- ETL
- Redshiftへのデータロード

データ利用のためのマエシヨリ！！！！

AWS のビッグデータ・プラットフォーム



AWS のビッグデータ・プラットフォームとData Pipeline



前処理ってめんどくさい

こういうのやめたい
ですよね・・・

前処理は処理コードを書く
だけではなく下記が必要！

- スケジュール管理
- 状態管理

```
s3 = boto.s3.connect_to_region()
key = create_key_from_timestamp() # some_key_2014112500.tar.gz

if(s3.get_key(key)):

    result1 = do_first_step(key)

    if(result1):
        do_second_step(result1)
        ...

    else:
        print 'Error in first step...'
        sys.exit()

else:

    print 'error: no target object in s3 bucket. The key to be ...'
    sys.exit()
```

ETLってめんどくさい

データのマイグレーションが簡単じゃない

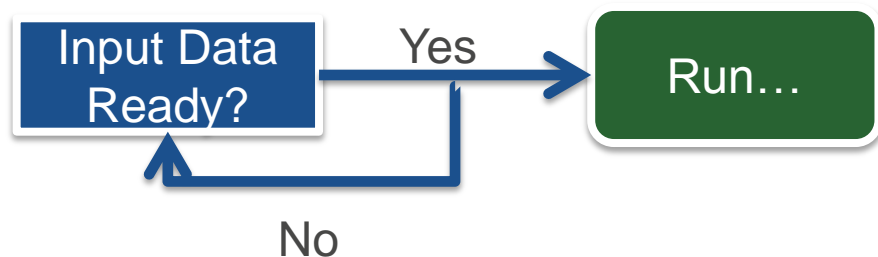
- 異なる場所、異なるフォーマット
 - S3, RDS, EMR, Redshift, DynamoDB
 - オンプレミス

依存管理が複雑

- 入力データが存在
- 前の処理が成功

異常処理が必要

- 失敗時のリトライ
- タイムアウト
- イベント通知



たとえば・・・

定期的に起動中のEC2のリストを取得してS3に保存

定期的に起動中のEC2のリストを取得してS3に保存

1. スクリプトを書く

2. EC2にそのスクリプトを配置してcronで動かす

- EC2が落ちたらどうする？
- cronのログはどこで管理する？
- EC2の時計がずれたらどうする

定期的に起動中のEC2のリストを取得してS3に保存

1. スクリプトを書く

—

~~2. EC2にそのスクリプトを配置してcronで動かす~~

~~--EC2が落ちたらどうする？~~

~~--cronのログはどこで管理する？~~

~~--EC2の時計がずれたらどうする~~

2. Data Pipelineで動かす！

定期的に起動中のEC2のリストを取得してS3に保存

Pipeline作成 - 1/2

Pipeline名

Create Pipeline

Name ListEc2

Description (optional)

Source

- Build using a template
- Import a definition
- Build using Architect

Schedule

Run

- once on pipeline activation
- on a schedule

Run every 15 minute(s)

Starting

- on pipeline activation
- 2015-09-08 05:20 UTC (Current time is 05:37 UTC)

Ending

- never
- after 1 occurrence(s)
- 2015-09-09 05:20 UTC (Current time is 05:37 UTC)

テンプレートや定義済みのJSONを使うかどうか

インターバル

定期的に起動中のEC2のリストを取得してS3に保存

Pipeline作成 - 2/2

The screenshot shows the 'Pipeline Configuration' page in the AWS Data Pipeline console. It is divided into three main sections: 'Logging', 'Security/Access', and 'Tags'. The 'Logging' section has 'Logging' set to 'Enabled' and 'S3 location for logs' set to 's3://YOUR_BUCKET/DIR/'. The 'Security/Access' section has 'IAM roles' set to 'Custom', with 'Pipeline role' set to 'DataPipelineDefaultRole' and 'EC2 instance role' set to 'general'. The 'Tags' section is empty. Three orange lines with arrows point from Japanese text annotations to the 'Logging', 'Pipeline role', and 'EC2 instance role' fields.

ログの吐き出し先

Logging Enabled Disabled Copy execution logs to S3. [More](#)

S3 location for logs
s3://YOUR_BUCKET/DIR/

Data Pipeline自身が利用するIAM Role

Security/Access

IAM roles Default Custom IAM Roles let you control permissions for AWS Data Pipeline and your EC2 applications. [More](#)

Pipeline role Control what AWS Data Pipeline can do with resources in your account. [More](#)

EC2 instance role Control what EC2 applications can do with resources in your account. [More](#)

Data Pipelineが起動するEC2等が利用するIAM Role

Tags

i Add up to 10 tags to your pipeline. These tags will be applied to the pipeline as well as any resources created by the pipeline. A tag consists of a case-sensitive key-value pair. [Learn more](#)

Key	Value (Optional)
<input type="text" value="Add key to create"/>	<input type="text"/>

Cancel

定期的に起動中のEC2のリストを取得してS3に保存

Architect

The screenshot displays the Amazon CloudFormation Architect console interface. At the top, the breadcrumb navigation shows 'Data Pipeline > List Pipelines > Architect: ListEC2 (df-0077777GWNLKPONAPAT) [Pending]'. Below this, a toolbar contains several action buttons: 'Add activity', 'Add data node', 'Save pipeline', 'Activate', 'Export', and 'View/Edit tags'. The main workspace is currently empty, featuring a central 'ORGANIZE LAYOUT' button with a directional icon. On the right side, a sidebar menu is visible with the following expandable sections: 'Activities', 'DataNodes', 'Schedules', 'Resources', 'Preconditions', 'Others', and 'Parameters'. At the bottom left of the console, the 'Errors/Warnings' section is partially visible.

定期的に起動中のEC2のリストを取得してS3に保存

Architect

The screenshot shows the AWS Architect console for a Data Pipeline named "ListEC2 (df-05812081F8894IBHLMD) [Active]". The pipeline consists of a "ListActivity" (ShellCommandActivity) connected to a "DefaultDataNode1". The ListActivity is configured with a "Name" of "ListActivity" and a "Type" of "ShellCommandActivity". The console interface includes buttons for "Add activity", "Add data node", "Save pipeline", "Activate", "Export", and "View/Edit tags". The right-hand panel shows the configuration for the "ListActivity", including fields for "Name", "Type", and "DataNodes".

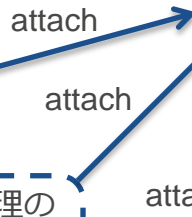
アクティビティ。どんな処理をどういったインターバルで、どのリソースの上で動かすかを定義する。

リソース。実際に処理を実行するEC2やEMRの定義の集合。

スケジュール。処理のインターバルの定義の集合。

アクション。タスク成功/失敗時の通知等の定義の集合

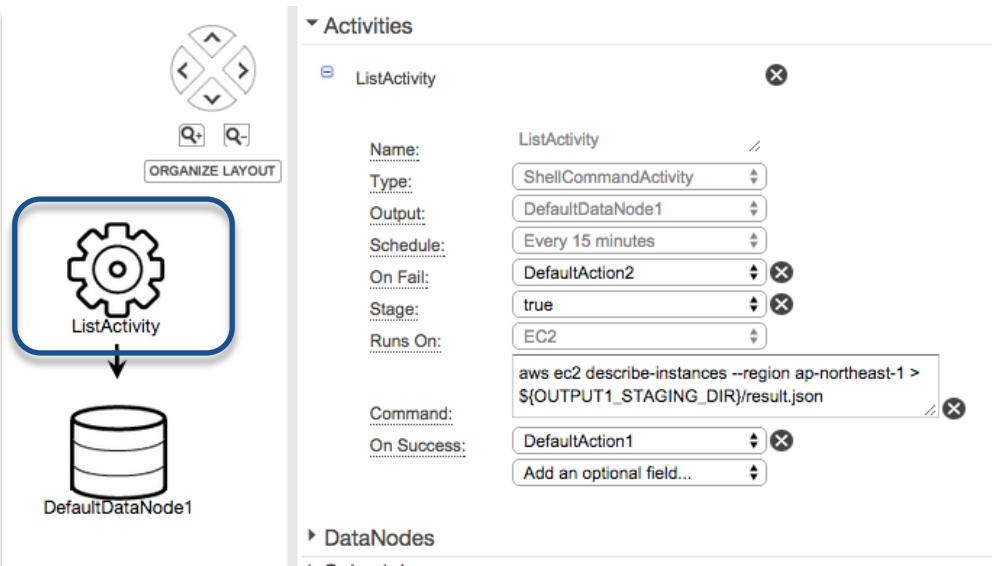
データノード。S3やRedshift、RDSを表す。例えばS3ならバケットやディレクトリ情報を持つ。



定期的に起動中のEC2のリストを取得してS3に保存

アクティビティ

Architect

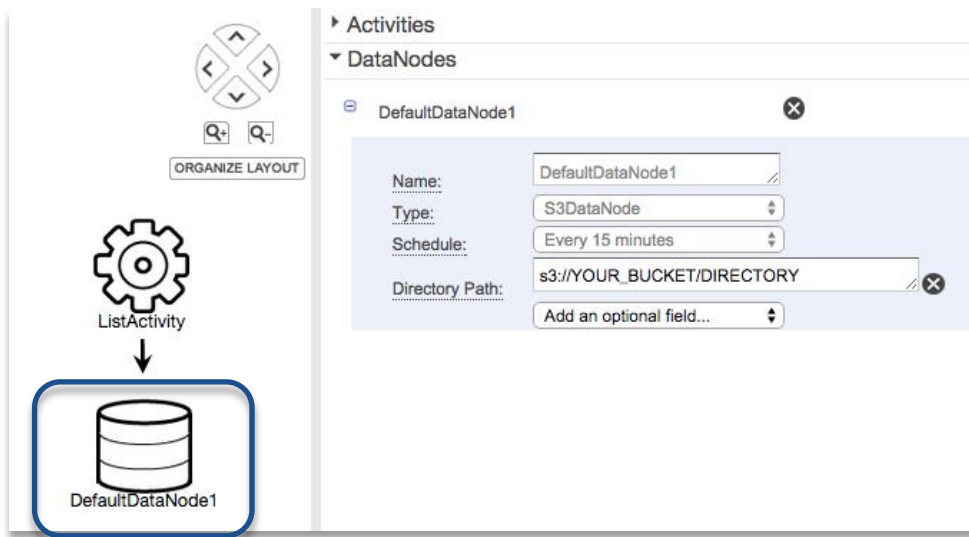


- Output: データの出力先
 - 結果を出力するS3バケット。詳細DataNodeで説明。
- On Success/On Fail
 - 成功/失敗時のアクション

- Type: ShellCommandActivity
 - 指定したシェルスクリプトを実行するアクティビティ
- Schedule: 15分毎
- Runs On
 - アクティビティを実行するEC2。詳細はResourceで説明。
- Command
 - 実行するコマンド
 - CommandではなくてScriptUriを使うと、S3に配置したスクリプトを指定可能
- Stage: True
 - `${INPUT1_STAGING_DIR}`, `${OUTPUT1_STAGING_DIR}`という変数が見えるようになる。S3からのデータ、S3へ送るデータが配置される/配置するディレクトリ。

定期的に起動中のEC2のリストを取得してS3に保存

Architect



データノード

- Type: S3
 - S3, Redshift, RDSなどの選択肢がある
- Directory Path:
 - S3のパス

定期的に起動中のEC2のリストを取得してS3に保存

Architect

The screenshot shows the configuration for a scheduled activity in the AWS IAM console. The activity is named "Every 15 minutes" and is of type "Schedule". The period is set to 15 minutes. The start at field is set to "FIRST_ACTIVATION_DATE_TI".

Activities

DataNodes

Schedules

Every 15 minutes

Name: Every 15 minutes

Type: Schedule

Period: 15 Minute(s)

Start At: FIRST_ACTIVATION_DATE_TI

Add an optional field...

Resources

Preconditions

Others

Parameters

スケジュール

- Type: Schedule
- Period
 - ここで実際のインターバルを設定する。
 - 最小インターバルは15分。
- Start At
 - スケジュールの開始日時

定期的に起動中のEC2のリストを取得してS3に保存

Architect

The screenshot shows the configuration for an EC2 resource in a CloudFormation stack. The resource is named 'EC2' and is of type 'Ec2Resource'. It is configured to run every 15 minutes. The configuration includes the following fields:

- Name: EC2
- Type: Ec2Resource
- Subnet Id: subnet-cabbba8c
- Schedule: Every 15 minutes
- Resource Role: general
- Role: DataPipelineDefaultRole
- Image Id: ami-1c1b9f1c
- Security Group Ids: sg-0a86016f
- Instance Type: t2.micro
- Key Pair: YOUR_KEY
- Terminate After: 10 Minute(s)

リソース

- Type: Ec2Resource
 - EC2もしくはEMRから選択
- Subnet Id
 - リソースを起動するSubnet
- Schedule
 - リソースを起動するインターバル
- Terminate After
 - 起動したリソースの削除時間の設定
- Image Id/Instance Type
 - 起動するリソースのAMIとインスタンスタイプ

定期的に起動中のEC2のリストを取得してS3に保存

Architect

The screenshot shows the AWS Architect console interface. On the left, a navigation pane lists categories: Activities, DataNodes, Schedules, Resources, Preconditions, and Others. Under 'Others', two actions are visible: DefaultAction1 and DefaultAction2. Each action has a list of configuration fields:

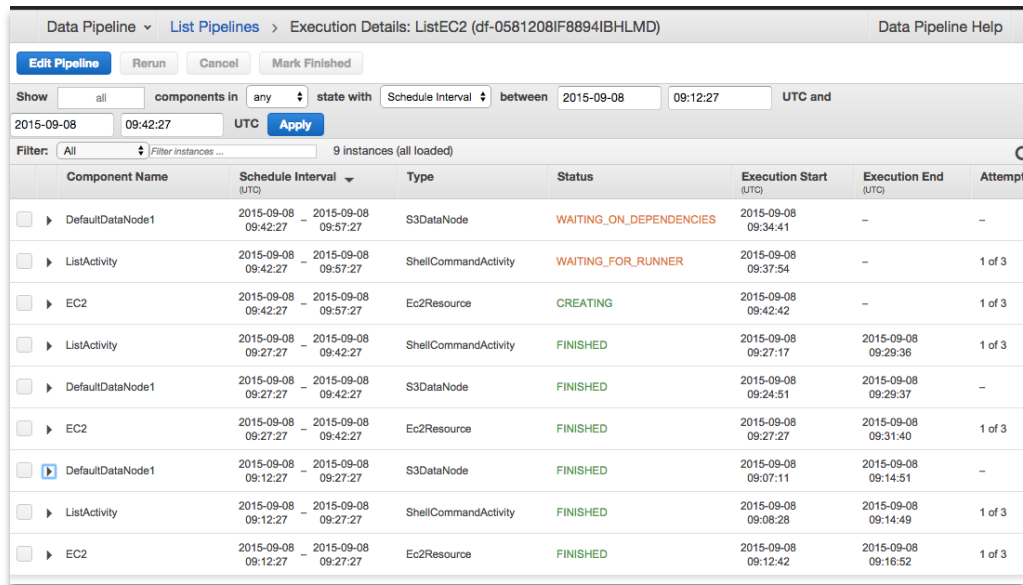
- DefaultAction1:**
 - Name: DefaultAction1
 - Type: SnsAlarm
 - Role: DataPipelineDefaultRole
 - Subject: Data Pipeline Test
 - Message: Succeeded!
 - Topic Arn: arn:aws:sns:ap-northeast-
- DefaultAction2:**
 - Name: DefaultAction2
 - Type: SnsAlarm
 - Role: DataPipelineDefaultRole
 - Subject: Data Pipeline Test
 - Message: Failed!

アクション

- Type: SnsAlarm
 - SnsAlarmもしくはTerminateから選択
- Role
 - アクションを実行するIAM Role
- Subject
 - メッセージのタイトル
- Message
 - メッセージの内容
- Topic Arn
 - SNSのトピックのARN(識別子)

定期的に起動中のEC2のリストを取得してS3に保存

Execution Details



Data Pipeline > List Pipelines > Execution Details: ListEC2 (df-0581208IF8894IBHLMD) Data Pipeline Help

Buttons: Edit Pipeline, Rerun, Cancel, Mark Finished

Show: all components in any state with Schedule Interval between 2015-09-08 09:42:27 UTC and 2015-09-08 09:12:27 UTC Apply

Filter: All 9 instances (all loaded)

Component Name	Schedule Interval (UTC)	Type	Status	Execution Start (UTC)	Execution End (UTC)	Attempts
DefaultDataNode1	2015-09-08 09:42:27 - 2015-09-08 09:57:27	S3DataNode	WAITING_ON_DEPENDENCIES	2015-09-08 09:34:41	-	-
ListActivity	2015-09-08 09:42:27 - 2015-09-08 09:57:27	ShellCommandActivity	WAITING_FOR_RUNNER	2015-09-08 09:37:54	-	1 of 3
EC2	2015-09-08 09:42:27 - 2015-09-08 09:57:27	Ec2Resource	CREATING	2015-09-08 09:42:42	-	1 of 3
ListActivity	2015-09-08 09:27:27 - 2015-09-08 09:42:27	ShellCommandActivity	FINISHED	2015-09-08 09:27:17	2015-09-08 09:29:36	1 of 3
DefaultDataNode1	2015-09-08 09:27:27 - 2015-09-08 09:42:27	S3DataNode	FINISHED	2015-09-08 09:24:51	2015-09-08 09:29:37	-
EC2	2015-09-08 09:27:27 - 2015-09-08 09:42:27	Ec2Resource	FINISHED	2015-09-08 09:27:27	2015-09-08 09:31:40	1 of 3
DefaultDataNode1	2015-09-08 09:12:27 - 2015-09-08 09:27:27	S3DataNode	FINISHED	2015-09-08 09:07:11	2015-09-08 09:14:51	-
ListActivity	2015-09-08 09:12:27 - 2015-09-08 09:27:27	ShellCommandActivity	FINISHED	2015-09-08 09:08:28	2015-09-08 09:14:49	1 of 3
EC2	2015-09-08 09:12:27 - 2015-09-08 09:27:27	Ec2Resource	FINISHED	2015-09-08 09:12:42	2015-09-08 09:16:52	1 of 3

- いわゆる実行状況や実行ログのが確認できる
- うまく動かないときはここからリソース起動やアクティビティのログを調査することができる。

定期的に起動中のEC2のリストを取得してS3に保存

Behind the scene

スケジュールで定義した「15分ごと」に・・・

リソースで定義にしたがってのEC2を起動。(Terminate Afterで定義した時間がくると削除される)

アクティビティで定義されたコマンドを実行

```
aws ec2 describe-instances --region ap-northeast-1 \  
> `${OUTPUT1_STAGING_DIR}/result.json
```

アップロード先はデータノードで定義

``${OUTPUT1_STAGING_DIR}`に配置されたファイルをS3へアップロード

定期的に起動中のEC2のリストを取得してS3に保存

1. スクリプトを書く

—

~~2. EC2にそのスクリプトを配置してcronで動かす~~

~~--EC2が落ちたらどうする？~~

~~--cronのログはどこで管理する？~~

~~--EC2の時計がずれたらどうする~~

2. Data Pipelineで動かす！

AWS Data Pipeline: 例えの使い道・

Pipelineのテンプレートより→

Getting Started

- Getting Started using ShellCommandActivity

AWS Command Line Interface (CLI) Templates

- Run AWS CLI command

DynamoDB Templates

- DynamoDB cross-regional table copy

- Export DynamoDB table to S3

- Import DynamoDB backup data from S3

Elastic MapReduce (EMR) Templates

- Run job on an Elastic MapReduce cluster

RDS Templates

- Full copy of RDS MySQL table to S3

- Incremental copy of RDS MySQL table to S3

- Load S3 data into RDS MySQL table

Redshift Templates

- Full copy of RDS MySQL table to Redshift

- Incremental copy of RDS MySQL table to Redshift

- Load AWS detailed billing report into Redshift

- Load data from S3 into Redshift

AWS Data Pipeline

AWS Data Pipelineとは

- サービス間のデータ統合・処理をスケジュールベースで自動化してくれるサービス
- 使い始めるにはワークフローを定義する
 1. データ：データソースや出力先の定義
 2. アクティビティ：処理の内容を定義
 3. スケジュールと依存関係：処理の依存関係とスケジュールを定義
 4. 通知：イベントの通知先を定義
- オンプレミス環境との連携も可能

Pipelineの定義

ワークフローを定義するための登場人物

- Data Node: データの場所、フォーマット
- Activity: データ処理のアクティビティ
- Schedule: 処理実行のスケジュール
- Resource: 処理や条件チェックを行うリソース
- Precondition: 処理実行の条件
- Action: 通知を送る方法

アクティビティ (Activities)



データ移動や処理の全体を管理

- 入出力、スケジュール、処理内容、リソース、通知アクション
- AWSとオンプレミスにて実行可能

サポートするアクティビティ一覧

- CopyActivity
- EmrActivity
- HiveActivity
- HiveCopyActivity
- PigActivity
- RedshiftCopyActivity
- SqlActivity
- ShellCommandActivity

ShellActivity1

Name:	ShellActivity1
Type:	ShellCommandActivity
Schedule:	Hourly S3 Upload
Script Uri:	s3://example/s3-upload.sh
Runs On:	EC2: uploader
On Success:	SNS: success report
	Add an optional field...

データノード



Input / Outputデータの場所やタイプを定義

- S3パス
- SQL データベース
- DynamoDB
- Redshift

フォーマット指定は自由

- CSV Data Format
- カスタマイズ

S3: upload dir

Name:	S3: upload dir
Type:	S3DataNode
Schedule:	Hourly S3 Upload
Directory Path:	s3://mybucket/upload
Precondition:	S3: prefix exist
	Add an optional field...

データとテーブルのステー징

データにアクセスしやすいための仕組み

- データを自動的にリソースにコピー / テーブルを自動的に作成
 - データがリソースのローカルにあるように
- サポートデータノード: S3DataNode, SqlDataNode
- サポートアクティビティ
 - ShellCommandActivity
 - 初期値: off。Stage = true にセット
 - 変数: `${INPUTx_STAGING_DIR}` , `${OUTPUTx_STAGING_DIR}`
 - HiveActivity
 - 初期値: on.
 - 変数: `${inputx}`, `${outputx}`

Table作成は
不要

```
{  
  "id": "MyHiveActivity",  
  ...  
  "hiveScript": "INSERT OVERWRITE TABLE ${output1} select * from ${input1};"  
},
```

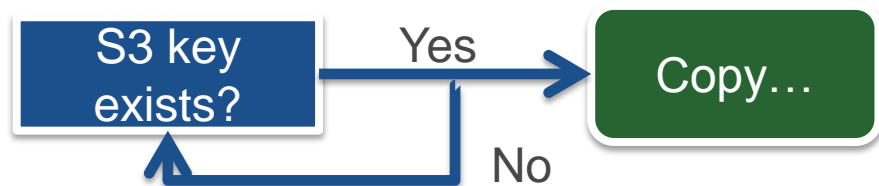
依存関係(Preconditions)

条件が成立した場合のみ後継タスクを実行

- DynamoDB tableが存在／データがある
- S3キーが存在
- S3プリフィックスが存在
- 独自のShellコマンド実行が成功
- 依存するpipelineタスクが成功

☐ S3: prefix exist

Name:	S3: prefix exist
Type:	S3PrefixNotEmpty
Role:	DataPipelineDefaultRole
S3 Prefix:	s3://mybucket/upload
	Add an optional field...

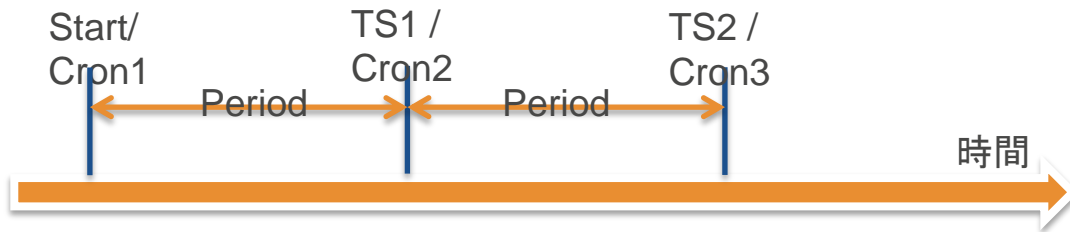


スケジュール



処理を実行するタイミング

- Cronスタイル: 指定した間隔のstart時点で起動
- Time Seriesスタイル: 指定した間隔のend時点で起動
 - ただし、EC2やEMRリソースは常にstart時点で作成
- 間隔: 15分、時、日、週など
 - 15min ~ 3year



Hourly S3 Upload

Name:	Hourly S3 Upload
Type:	Schedule
Start Date Time (in UTC):	2014-03-15 02:56:00
Period:	1 Hour(s)
End Date Time (in UTC):	2014-03-21 02:56:00
	Add an optional field...

スケジュール(2)

Backfillタスク

- 開始時間に過去を指定した場合、現在まで遡ってタスクを繰り返し実行
 - テストに便利
- 開始時間が1日以前の場合、タスク起動しない
 - CLIで `--force` 引数を指定して起動可能

タイムゾーン

- 初期値: UTC, “YYYY-MM-DDTHH:MM:SS” フォーマット
- 変数にタイムゾーン指定可能
 - `#{inTimeZone(myDateTime,'Asia/Tokyo')}`

リソース

タスクを実行するリソース

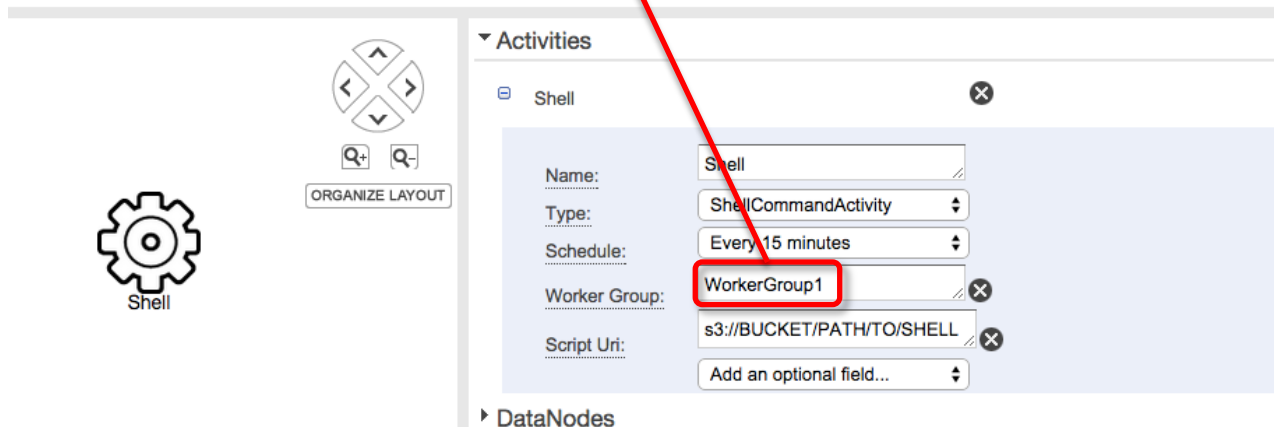
- Task Runnerと呼ばれるAWS Data Pipelineのエージェントプロセスが実体
- EC2: EC2-ClassicとEC2-VPC両方サポート
- EMR: タスクノードにspot instance利用可能
- Multi-regionのリソース管理が可能

EC2: upload server

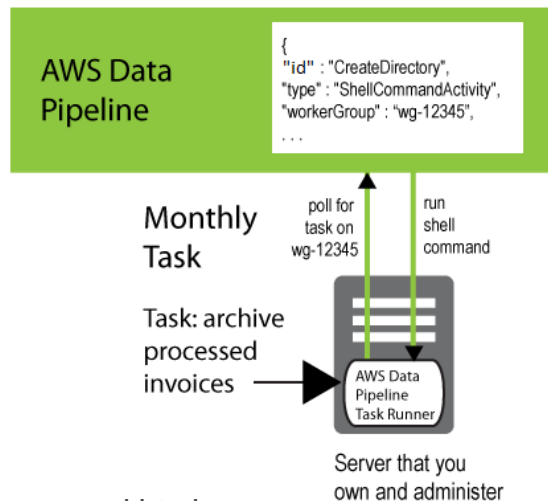
Name:	EC2: upload server
Type:	Ec2Resource
Terminate After:	5 Minute(s)
Schedule:	Hourly S3 Upload
Role:	DataPipelineDefaultRole
Log Uri:	s3://yifeng/logs/datapipeline/
Resource Role:	DataPipelineDefaultResourceRc
Image Id:	ami-1624987f
	Add an optional field...

Task Runnerは 既存EC2やオンプレミスサーバーでも動く

```
java -jar TaskRunner-1.0.jar --config ~/credentials.json --workerGroup=WorkerGroup1 --  
region=MyRegion --logUri=s3://mybucket/foldername
```



The screenshot shows the AWS Data Pipeline console interface. On the left, there is a navigation menu with a gear icon labeled 'Shell' and an 'ORGANIZE LAYOUT' button. The main area displays the configuration for a 'Shell' activity under the 'Activities' section. The configuration fields are: Name: Shell, Type: ShellCommandActivity, Schedule: Every 15 minutes, Worker Group: WorkerGroup1 (highlighted with a red box), and Script Uri: s3://BUCKET/PATH/TO/SHELL. A red arrow points from the 'WorkerGroup1' field in the console to the 'WorkerGroup' parameter in the command line above.

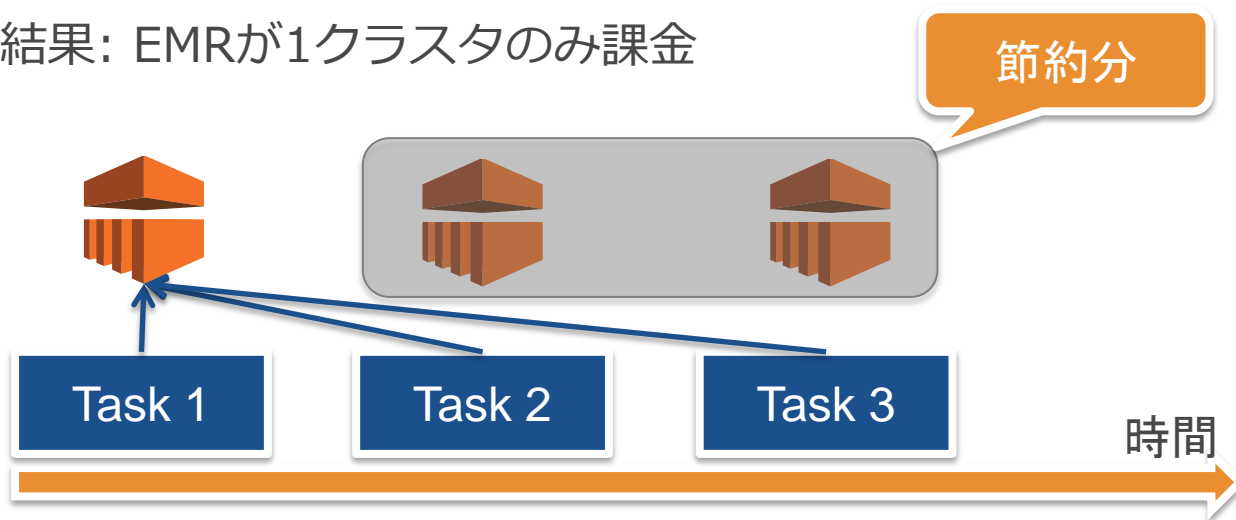


<http://docs.aws.amazon.com/datapipeline/latest/DeveloperGuide/dp-how-task-runner-user-managed.html>

リソース(2)

アクティビティとリソースのスケジュールを別々に指定可能

- リソースを最大限に利用
 - Activityスケジュール: 20分
 - Resourceスケジュール: 1時間
 - 結果: EMRが1クラスタのみ課金



イベントと異常管理

イベントが発生するタイミング

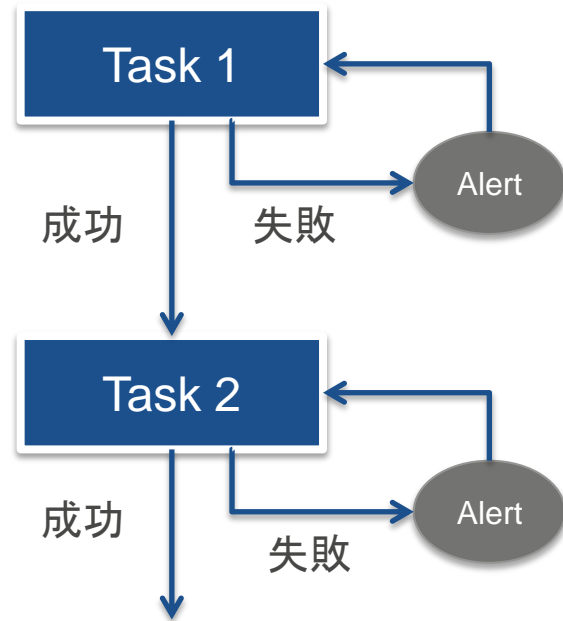
- 成功時
- 失敗時
- 遅れが発生した場合

設定可能なアクション

- SNS通知
- リソースを削除

失敗時自動リトライ

- 初回実行含め、1~6回の実行回数を設定可能
- 初期値は3回



ファイル出力

GUIの操作結果をJSONに出力可能

- Pipelineのバージョン管理
- 出力したファイルは

```
./datapipeline --create pipeline_name --put pipeline_file  
--activate --force
```

- インポート前にバリデーションで
きる

```
./datapipeline --validate my-pipeline.json --credential  
credetials.json --force --id df-0123456789ABCD
```

```
{  
  "objects": [  
    {  
      "id": "ActivityId_YYbJV",  
      "schedule": {  
        "ref": "ScheduleId_X8kbH"  
      },  
      "scriptUri":  
"s3://mybucket/myscript.sh",  
      "name": "ShellActivity1",  
      "runsOn": {  
        "ref": "ResourceId_5nJlh"  
      },  
      ...  
    }  
  ]  
}
```

料金

- 無料使用枠あり
- アクティビティ、または依存関係の従量課金
- 実際に利用したリソース(EC2, S3など)の課金

	高頻度 (>1日1回)	低頻度
On AWS	\$1.00	\$0.60
オンプレミス	\$2.50	\$1.50
実行しないpipeline	\$1.00	

* 2014/3/19時点、東京リージョンのアクティビティ / 依存関係の月額単価

AWS Data Pipelineがもたらしてくれるもの

- 分離

- データと処理リソースの分離
- 処理リソースと処理ロジックの分離
- 処理ロジックとスケジュールの分離

- 統合

- 分散環境での整合性
- 一環したエラー処理や処理のやりなおし

AWS Data Pipelineがもたらしてくれるもの

例えばこんなパイプラインなら・・・

1.1.1.1, /login, 20140226000101, ...
192.168..., /home, 20140226011226, ...
1.1.1.2, /home, 20140226011331, ...

USER	PATH	TIMESTAMP
USER1	/login	2014-02-26 00:00:01
USER2	/home	2014-02-26 01:13:31

Webサーバー



fluentd

ログ (オリジナル)

S3



EMR



S3
処理済みデータ



Redshift

ETL済みデータ



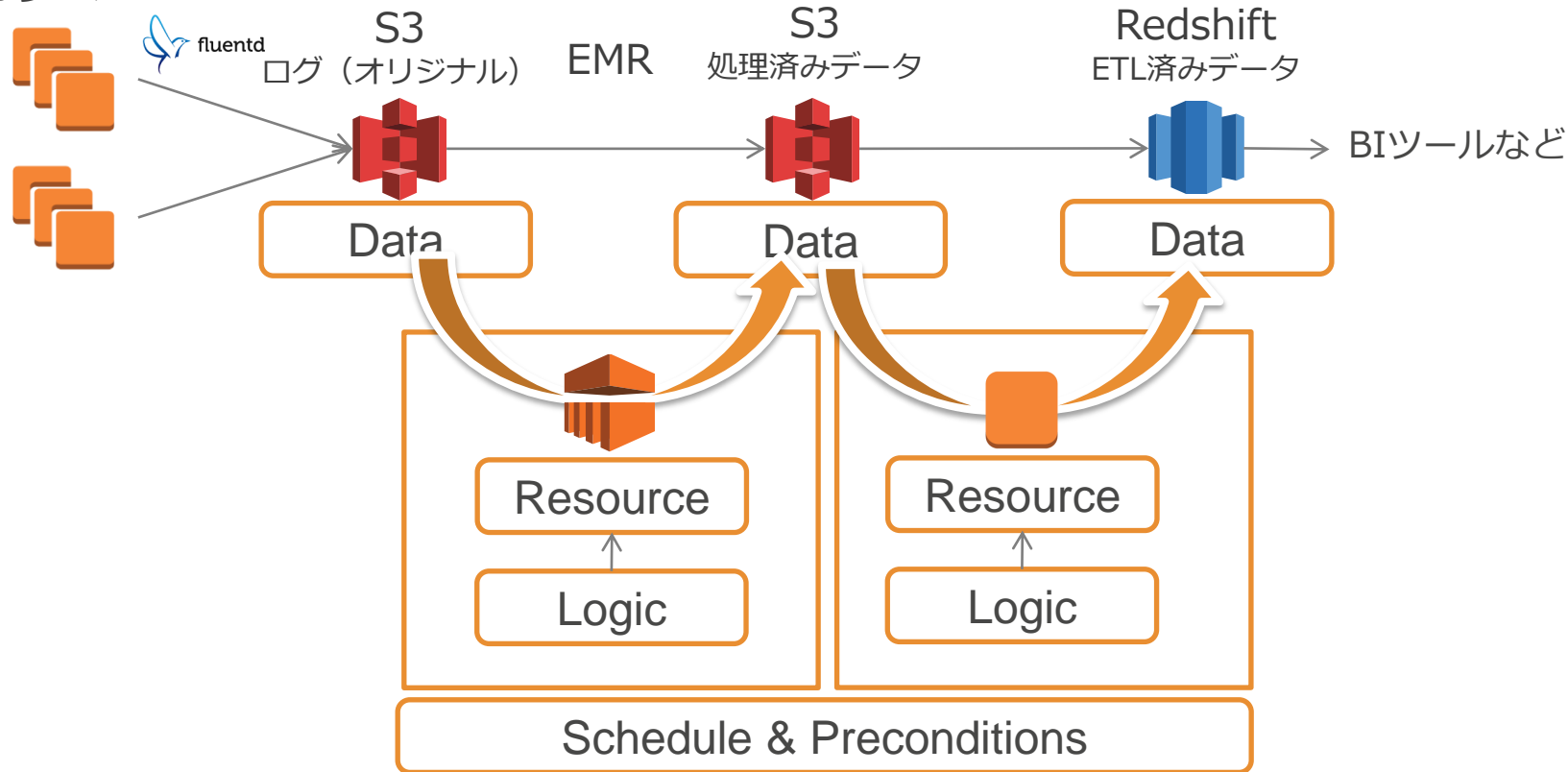
BIツールなど

1.1.1.1, /login, 20140226000101, ...
~~192.168..., /home, 20140226011226, ...~~
1.1.1.2, /home, 20140226011331, ...

AWS Data Pipelineがもたらしてくれるもの

例えばこんなパイプラインなら・・・

Webサーバー



アジェンダ

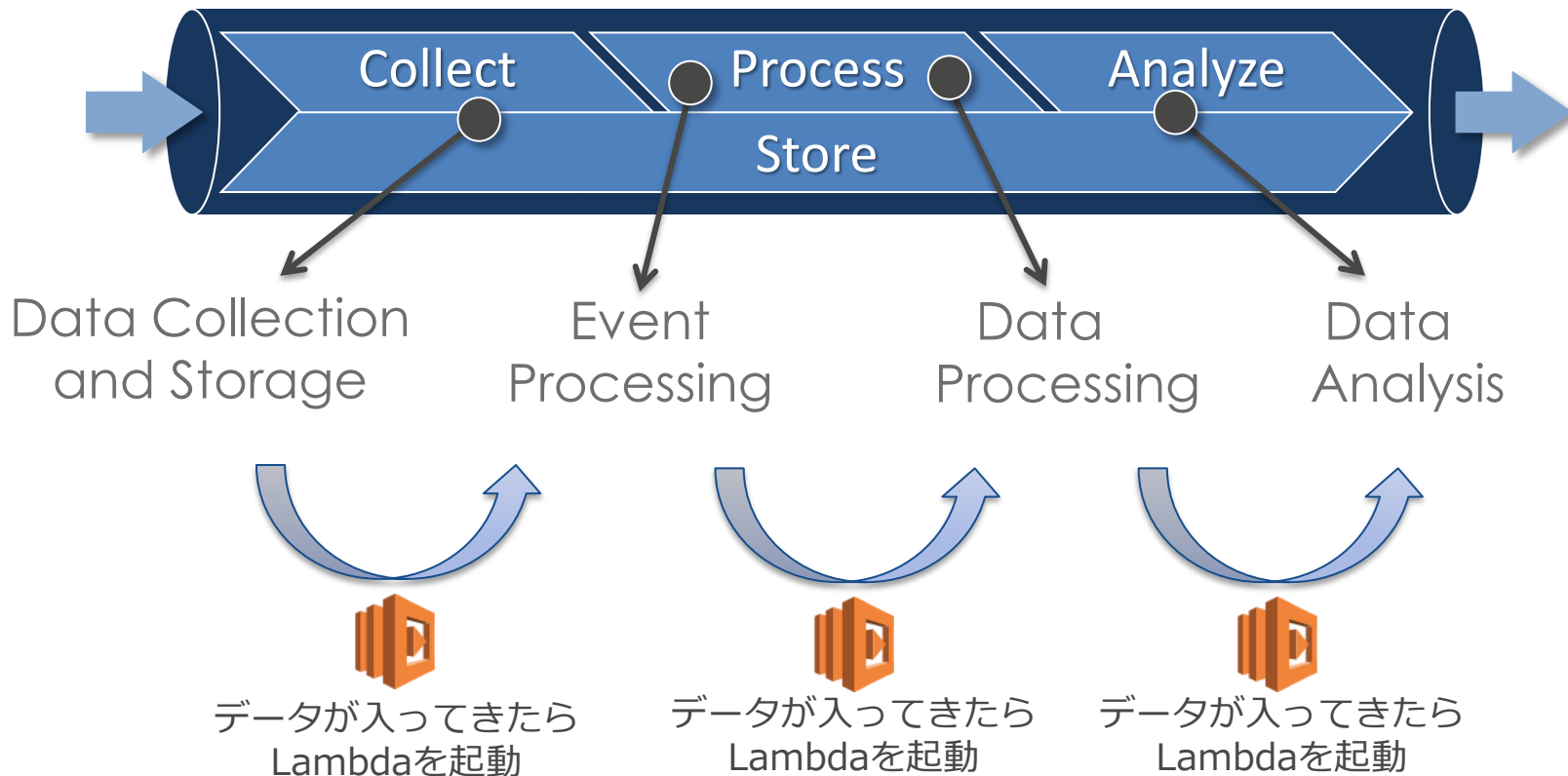
1. What is AWS Data Pipeline
2. When not to use
3. Case studies

AWS Data Pipelineはスケジュールベースのサービス

- イベントドリブンな処理を行いたければAWS Lambdaを使う

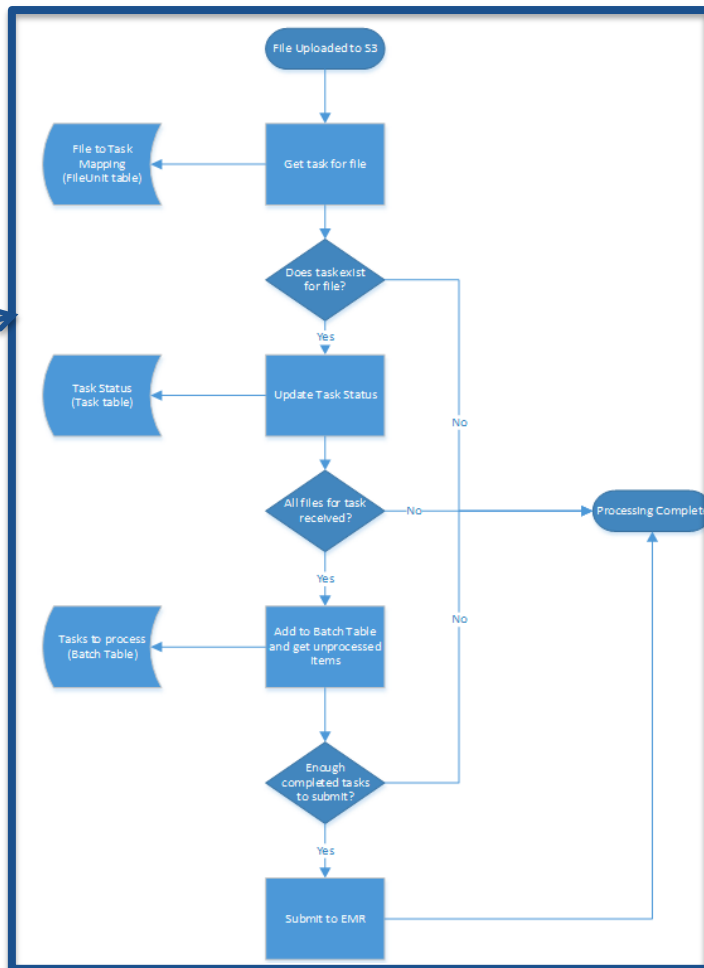
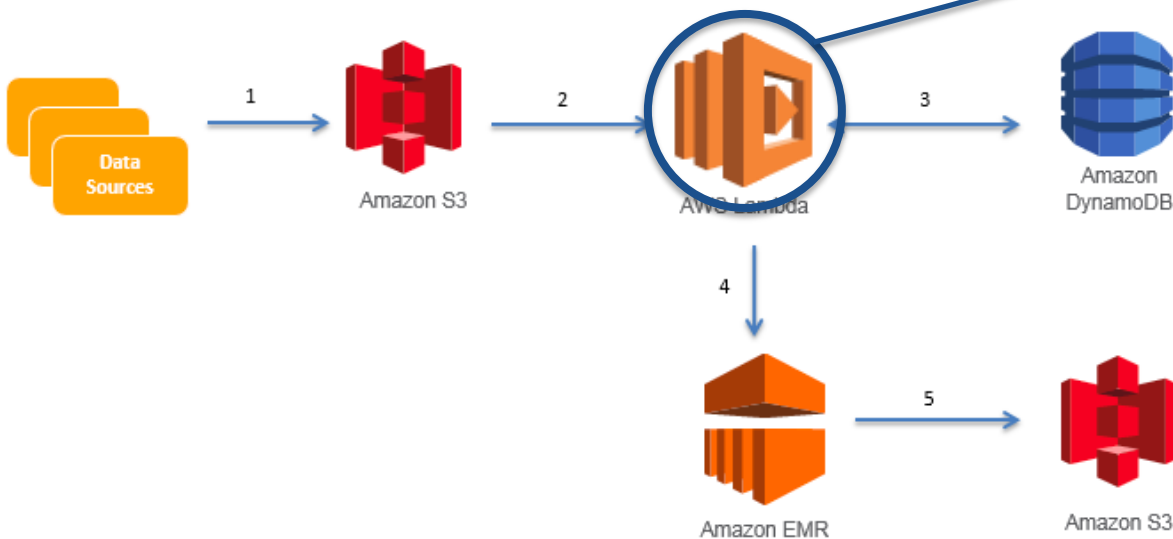


AWS Lambda as a pipeline glue



AWS Lambdaの事例: Expedia

Pipelineのオーケストレーション



<http://blogs.aws.amazon.com/bigdata/post/Tx1R28PXR3NAO11/How-Expedia-Implemented-Near-Real-time-Analysis-of-Interdependent-Datasets>

アジェンダ

1. What is AWS Data Pipeline
2. When not to use
3. Case studies

AWS Data Pipelineがフィットするケース

- 1日に1回などの頻度で行われる処理
 - 複数の拠点やデータセンターからのデータをS3へ回収
 - そしてRedshiftへロード
 - 定期的なクローリングの実行とS3への結果アップロード
 - そしてRedshiftへロード
 - データベースのバックアップ
 - RDSのようにマネージドサービスではないデータベースも、Data Pipelineを使えばバックアップを自動化できる

定義済みのアクティビティからユースケースを考えてみる

- CopyActivity
 - MysqlDataNodeとS3DataNode間、もしくはそれら同士のコピーを実施してくれる
 - RDBのバックアップ
- SqlActivity
 - データベースにSQLを発行してくれるアクティビティ。
 - 定期的にローデータテーブルを集計して結果をサマリテーブルに入れるなど。
- EmrActivity
 - EMR上でJARを指定して処理を実行させるためのアクティビティ。
 - 実行するJAR次第でなんでもできてしまうのである意味ShellCommandActivity並に自由度が高い。
 - 素直に使うなら大量のデータに対してのETLや集計、分析など。
- HiveActivity
 - EMR上でHiveを実行させるためのアクティビティ。
 - 生データが格納されているS3バケットからデータを取り出して整形したりフィルタして別バケットに格納、など。
 - S3とDynamoDB間のデータやりとりにも使えるので、DynamoDBのデータインポート/エクスポートにも。

定義済みのアクティビティからユースケースを考えてみる

- HiveCopyActivity
 - S3とDynamoDBのデータ連携に特化したアクティビティ。
- PigActivity
 - EMR上でPigを実行させるためのアクティビティ。
- RedshiftCopyActivity
 - S3からRedshiftにデータを取り込むためのアクティビティ。
- ShellCommandActivity
 - シェルスクリプトを実行するアクティビティ。なんでもできる！

まとめ

AWS Data Pipelineがもたらしてくれるもの: バッチ処理を管理しやすくしてくれる！

- 分離

- データと処理リソースの分離
- 処理リソースと処理ロジックの分離
- 処理ロジックとスケジュールの分離

- 統合

- 分散環境での整合性
- 一環したエラー処理や処理のやりなおし

まとめ: バッチ処理の管理省力化

- もう少し極端に言うところのこと -

1. スクリプトを書く

—

~~2. EC2にそのスクリプトを配置してcronで動かす~~

~~--EC2が落ちたらどうする？~~

~~--cronのログはどこで管理する？~~

~~--EC2の時計がずれたらどうする~~

2. Data Pipelineで動かす！

CopyActivity, EmrActivityなど、いろいろプリセットもある！

AWS Data Pipelineが面倒見てくれる！