

Guia de desenvolvimento do Amazon EMR no EKS

Amazon EMR



Copyright © 2025 Amazon Web Services, Inc. and/or its affiliates. All rights reserved.

Amazon EMR: Guia de desenvolvimento do Amazon EMR no EKS

Copyright © 2025 Amazon Web Services, Inc. and/or its affiliates. All rights reserved.

As marcas comerciais e imagens comerciais da Amazon não podem ser usadas no contexto de nenhum produto ou serviço que não seja da Amazon, nem de qualquer maneira que possa gerar confusão entre os clientes ou que deprecie ou desprestigie a Amazon. Todas as outras marcas comerciais que não pertencem à Amazon pertencem a seus respectivos proprietários, que podem ou não ser afiliados, patrocinados pela Amazon ou ter conexão com ela.

Table of Contents

O que é o Amazon EMR no EKS?	1
Arquitetura do Amazon EMR no EKS	2
Noções básicas de conceitos e terminologia do Amazon EMR no EKS	2
Namespace do Kubernetes	2
Cluster virtual	3
Execução de trabalho	3
Contêineres do Amazon EMR	3
O que acontece quando você envia trabalhos para um cluster virtual do Amazon EMR no	
EKS	4
Conceitos básicos do Amazon EMR no EKS	5
Execução de uma aplicação do Spark	6
Práticas recomendadas	12
Segurança	12
Envio de trabalho para o Pyspark	12
Armazenamento	12
Integração com o metastore	13
Depuração	13
Solução de problemas do Amazon EMR no EKS	13
Posicionamento de nó	13
Performance	13
Otimização de custo	14
Usando AWS Outposts	14
Personalização de imagens do Docker	15
Como personalizar imagens do Docker	15
Pré-requisitos	16
Etapa 1: recuperar uma imagem base do Amazon Elastic Container Registry (Amazon	
ECR)	16
Etapa 2: personalizar uma imagem base	17
Etapa 3: (opcional, mas recomendada) validar uma imagem personalizada	18
Etapa 4: publicar uma imagem personalizada	20
Etapa 5: enviar uma workload do Spark no Amazon EMR usando uma imagem	
personalizada	21
Personalização de imagens do Docker para endpoints interativos	23
Trabalho com imagens de múltiplas arquiteturas	25

Detalhes sobre como selecionar um URI de imagem base	27
Contas de registro do Amazon ECR	28
Considerações para personalizar imagens	29
Execução de trabalhos do Flink	31
Operador do Kubernetes para Flink	31
Configuração	32
Instalando o operador Flink Kubernetes	33
Execução de uma aplicação do Flink	35
Permissões do perfil de segurança para executar uma aplicação do Flink	39
Desinstalação do operador	41
Kubernetes nativo do Flink	42
Configuração	42
Começar	43
Requisitos de segurança	45
Personalização de imagens do Docker para o Flink e o FluentD	46
Pré-requisitos	46
Recuperar uma imagem base do Amazon Elastic Container Registry	46
Personalização de uma imagem base	47
Como publicar a imagem personalizada	48
Como enviar uma workload	49
Monitoramento	49
Uso do Amazon Managed Service for Prometheus	50
Uso da interface do usuário do Flink	51
Uso da configuração de monitoramento	52
Como o Flink oferece suporte à alta disponibilidade e resiliência do trabalho	57
Uso da alta disponibilidade	58
Otimização dos tempos de reinicialização	64
Desativação tranquila	71
Uso do Autoscaler	74
Ajuste automático dos parâmetros do escalador automático	76
Manutenção e solução de problemas para trabalhos do Flink no Amazon EMR no EKS	85
Manutenção de aplicações do Flink	85
Solução de problemas	86
Versões compatíveis	90
Execução de trabalhos do Spark	92
StartJobRun	93

Configuração	93
Envio de uma execução de trabalho com StartJobRun	122
Uso da classificação de envio de trabalho	124
Usando a classificação padrão do contêiner Amazon EMR	130
Operador do Spark	132
Configuração	133
Começar	133
Autoescalabilidade vertical	138
Desinstalar	143
Usando a configuração de monitoramento para monitorar o Spark	143
Segurança	151
spark-submit	160
Configuração	161
Começar	162
Segurança	163
Apache Livy	168
Configuração	169
Começar	170
Execução de uma aplicação do Spark	175
Desinstalação	177
Segurança	177
Propriedades de instalação	187
Solução de erros comuns de formato de variável de ambiente	193
Gerenciamento de execuções de trabalhos	194
Gerenciamento com a CLI	195
Execução de scripts do Spark SQL	201
Estados de execução de trabalho	204
Visualização de trabalhos no console	205
Erros comuns na execução de trabalhos	205
Uso de modelos de trabalho	211
Criação e uso de um modelo de trabalho para iniciar uma execução de trabalho	212
Definição de parâmetros de modelos de trabalhos	214
Controle do acesso aos modelos de trabalhos	216
Uso de modelos de pod	218
Cenários comuns	218
Habilitação de modelos de pod com o Amazon EMR no EKS	220

Campos do modelo de pod	222
Considerações sobre contêineres sidecar	225
Uso de políticas de repetição	227
Definição de uma política de repetição	227
Recuperação de um status da política	230
Monitoramento do trabalho	231
Descoberta de logs de drivers	231
Uso da alternância de log de eventos do Spark	231
Uso da alternância de log do contêiner do Spark	232
Uso da escalabilidade automática vertical	235
Configuração	235
Começar	238
Configuração	240
Monitoramento das recomendações	246
Desinstalação	248
Execução de workloads interativas	249
Visão geral dos endpoints interativos	249
Pré-requisitos para os endpoints interativos	251
AWS CLI	251
eksctl	251
Cluster do Amazon EKS	251
Concessão de acesso ao cluster	252
Ativação de perfis do IAM para contas de serviço	252
Criação de um perfil de execução de trabalho do IAM	252
Concessão de acesso para os usuários	253
Registro do cluster do Amazon EKS com o Amazon EMR	253
Load Balancer Controller	254
Criação de um endpoint interativo	254
Criação de um endpoint interativo	254
Especificação de parâmetros personalizados	255
	256
Parâmetros para um endpoint interativo	256
Definição de configurações para endpoints interativos	258
Monitorar trabalhos Spark	258
Modelos de pod personalizados	260
Implantação de um pod do JEG em um grupo de nós	260

Opções de configuração do JEG	264
Modificando parâmetros PySpark	265
Imagem de kernel personalizada	266
Monitoramento de endpoints interativos	267
Exemplos	270
Uso de cadernos Jupyter de hospedagem própria	270
Criar um grupo de segurança	271
Criação de um endpoint interativo	271
Obtenção do URL do servidor de gateway	272
Obtenção do token de autenticação	272
Implantação de um caderno	273
Limpeza	278
Como obter informações sobre endpoints interativos com comandos da CLI	278
	278
Listagem de endpoints interativos	280
Exclusão de um endpoint interativo	281
Upload dos dados	283
Pré-requisitos	283
Começar	283
Monitorar trabalhos	286
Monitore trabalhos com a Amazon CloudWatch Events	286
Automatize o Amazon EMR no EKS com eventos CloudWatch	287
Exemplo: configuração de uma regra que invoque o Lambda	288
Monitore o driver pod do trabalho com uma política de repetição usando o Amazon Events	
CloudWatch	289
Gerenciamento de clusters virtuais	290
Criação de um cluster virtual	290
Listagem de clusters virtuais	292
Descrição de um cluster virtual	292
Exclusão de um cluster virtual	292
Estados de um cluster virtual	292
Tutoriais	293
Uso do Delta Lake	293
Uso do Iceberg	294
Configurações de sessão do Spark para integração de catálogos	295
Usando PyFlink	296

Usando AWS Glue com Flink	297
Uso do Apache Hudi	300
Envio de um trabalho do Apache Hudi	300
Uso do RAPIDS para Spark	304
Uso do Spark no Redshift	308
Iniciar uma aplicação do Spark	309
Autenticação no Amazon Redshift	310
Leitura e gravação para o Amazon Redshift	312
Considerações	314
Uso do Volcano	315
Visão geral	315
Instalação	316
Envio: operador do Spark	317
Envio: spark-submit	319
Usando YuniKorn	320
Visão geral	320
Crie seu cluster do	320
Instalar YuniKorn	323
Envio: operador do Spark	323
Envio: spark-submit	326
Segurança	12
Práticas recomendadas	330
Aplicação do princípio de privilégio mínimo	330
Listagem de controle de acesso para endpoints	330
Obtenção das atualizações de segurança mais recentes para as imagens personalizadas	331
Limitação do acesso à credencial do pod	. 331
Isolamento de código de uma aplicação não confiável	331
Permissões de controle de acesso por perfil (RBAC)	331
Restrição do acesso às credenciais do perfil do IAM ou do perfil de instância do grupo de	
nós	332
Proteção de dados	333
Criptografia em repouso	334
Criptografia em trânsito	336
Gerenciamento de Identidade e Acesso	
Público	
Autenticação com identidades	338

Gerenciar o acesso usando políticas	342
Como o Amazon EMR no EKS funciona com o IAM	345
Uso de perfis vinculadas ao serviço	351
Políticas gerenciadas para o Amazon EMR no EKS	355
Uso de perfis de execução de trabalho com o Amazon EMR no EKS	357
Exemplos de políticas baseadas em identidade	359
Políticas para controle de acesso baseado em etiquetas	362
Solução de problemas	366
Usando o Amazon EMR no EKS com Lake Formation AWS	368
Como o Amazon EMR no EKS funciona com AWS o Lake Formation	368
Habilite o Lake Formation com o Amazon EMR no EKS	369
Considerações e limitações	378
Solução de problemas	380
Registro em log e monitoramento	382
Criptografia de logs	382
CloudTrail troncos	386
Concessão de Acesso do S3	388
Visão geral	388
Executar um cluster	389
Considerações	391
Validação de conformidade	391
Resiliência	391
Segurança da infraestrutura	391
Análise de configuração e vulnerabilidade	392
Endpoints da VPC de interface	392
Criação de uma política de endpoint da VPC para o Amazon EMR no EKS	393
Acesso entre contas	396
Pré-requisitos	396
Como acessar um bucket do Amazon S3 ou uma tabela do DynamoDB entre contas	397
Marcar recursos	402
Conceitos Básicos de Tags	402
Marcar com tag os recursos do	
Restrições de tag	404
Trabalhe com tags usando o AWS CLI e o Amazon EMR na API EKS	
Solução de problemas	
Falhas em trabalhos que usam PVC	

Verificação	406
Patch	407
Aplicação manual de patches	411
Falhas de escalabilidade automática vertical	413
Erro 403 Forbidden	413
Namespace não encontrado	413
Erro de credenciais do Docker	414
Falhas em operadores do Spark	414
Falha na instalação do chart do Helm	414
Exceção FileSystem sem suporte	415
Cotas e endpoints de serviço	417
Service endpoints	417
Cotas de serviço	420
Visualize cotas e solicite aumentos de cotas	421
Versões de liberação	422
Versões 7.9.0	423
Versões	423
Notas da versão	425
Alterações	427
emr-7.9.0-mais recente	427
emr-7.9.0-20250425	427
emr-7.9.0-flink-latest	427
emr-7.9.0-flink-20250425	428
Versões 7.8.0	428
Versões	428
Notas da versão	430
Alterações	431
emr-7.8.0-mais recente	432
emr-7.8.0-20250228	432
emr-7.8.0-flink-latest	432
emr-7.8.0-flink-20250228	432
Versões 7.7.0	433
Versões	433
Notas da versão	434
Alterações	436
emr-7.7.0-mais recente	436

emr-7.7.0-20250131	436
emr-7.7.0-flink-latest	437
emr-7.7.0-flink-20250131	437
Versões 7.6.0	437
Versões	437
Notas da versão	439
Atributos	441
Alterações	441
emr-7.6.0-mais recente	441
emr-7.6.0-20241213	441
emr-7.6.0-flink-latest	442
emr-7.6.0-flink-20241213	442
Versões 7.5.0	442
Versões	442
Notas da versão	443
Versões 7.4.0	443
Versões	443
Notas da versão	443
Versões 7.3.0	444
Versões	444
Notas da versão	445
Atributos	447
Alterações	447
emr-7.3.0-latest	448
emr-7.3.0-20240920	448
emr-7.3.0-flink-latest	448
emr-7.3.0-flink-29240920	448
Versões 7.2.0	449
Versões	449
Notas da versão	450
Atributos	452
emr-7.2.0-latest	453
emr-7.2.0-20240610	453
emr-7.2.0-flink-latest	453
emr-7.2.0-flink-20240610	453
Versões 7.1.0	454

Versões	454
Notas da versão	456
Atributos	457
emr-7.1.0-latest	457
emr-7.1.0-20240321	458
emr-7.1.0-flink-latest	458
emr-7.1.0-flink-20240321	458
Versões 7.0.0	458
Versões	459
Notas da versão	
Atributos	462
Alterações	462
emr-7.0.0-latest	462
emr-7.0.0-2024321	463
emr-7.0.0-20231211	463
emr-7.0.0-flink-latest	463
emr-7.0.0-flink-2024321	
emr-7.0.0-flink-20231211	464
Versões 6.15.0	464
Versões	464
Notas da versão	466
Atributos	
emr-6.15.0-latest	
emr-6.15.0-20240105	468
emr-6.15.0-20231109	468
emr-6.15.0-flink-latest	468
emr-6.15.0-flink-20240105	469
emr-6.15.0-flink-20231109	469
Versões 6.14.0	469
Versões	469
Notas da versão	471
Atributos	472
emr-6.14.0-latest	472
emr-6.14.0-20231005	472
Versões 6.13.0	473
Versões	473

Notas da versão		474
Atributos		476
emr-6.13.0-latest		476
emr-6.13.0-202308	14	476
Versões 6.12.0		477
Versões		477
Notas da versão		478
Atributos		479
emr-6.12.0-latest		479
emr-6.12.0-202403	21	480
emr-6.12.0-202307	01	480
Versões 6.11.0		480
Versões		480
Notas da versão		481
Atributos		482
emr-6.11.0-latest		483
emr-6.11.0-202309	05	483
emr-6.11.0-202305	09	483
Versões 6.10.0		484
emr-6.10.0-latest		486
emr-6.10.0-202309	05	487
emr-6.10.0-202306	24	487
emr-6.10.0-202304	21	487
emr-6.10.0-202304	03	487
emr-6.10.0-2023022	20	488
Versões 6.9.0		488
emr-6.9.0-latest		491
emr-6.9.0-2023090	5	491
emr-6.9.0-2023062	4	491
emr-6.9.0-2022110	8	492
Versões 6.8.0		492
emr-6.8.0-latest		496
emr-6.8.0-2023090	5	496
emr-6.8.0-2023062	4	496
emr-6.8.0-2022121	9	497
emr-6.8.0-2022080	2	497

Versões 6.7.0	497
emr-6.7.0-latest	. 499
emr-6.7.0-20240321	499
emr-6.7.0-20230624	500
emr-6.7.0-20221219	500
emr-6.7.0-20220630	500
Versões 6.6.0	500
emr-6.6.0-latest	. 502
emr-6.6.0-20240321	502
emr-6.6.0-20230624	502
emr-6.6.0-20221219	503
emr-6.6.0-20220411	503
Versões 6.5.0	503
emr-6.5.0-latest	. 504
emr-6.5.0-20240321	505
emr-6.5.0-20221219	505
emr-6.5.0-20220802	505
emr-6.5.0-20211119	505
Versões 6.4.0	506
emr-6.4.0-latest	. 507
emr-6.4.0-20240321	507
emr-6.4.0-20221219	508
emr-6.4.0-20210830	508
Versões 6.3.0	508
emr-6.3.0-latest	. 509
emr-6.3.0-20240321	510
emr-6.3.0-20220802	510
emr-6.3.0-20211008	510
emr-6.3.0-20210802	510
emr-6.3.0-20210429	511
Versões 6.2.0	511
emr-6.2.0-latest	. 512
emr-6.2.0-20240321	513
emr-6.2.0-20220802	513
emr-6.2.0-20211008	513
emr-6.2.0-20210802	513

	emr-6.2.0-20210615	514
	emr-6.2.0-20210129	514
	emr-6.2.0-20201218	514
	emr-6.2.0-20201201	514
	Versões 5.36.0	515
	emr-5.36.0-latest	516
	emr-5.36.0-20240321	516
	emr-5.36.0-20221219	517
	emr-5.36.0-20220620	517
	emr-5.36.0-20220525	517
•	Versões 5.35.0	517
	emr-5.35.0-latest	519
	emr-5.35.0-20240321	519
	emr-5.35.0-20221219	519
	emr-5.35.0-20220802	519
	emr-5.35.0-20220307	520
•	Versões 5.34	520
	emr-5.34.0-latest	521
	emr-5.34.0-20240321	521
	emr-5.34.0-20220802	522
	emr-5.34.0-20211208	522
,	Versões 5.33.0	522
	emr-5.33.0-latest	
	emr-5.33.0-20240321	
	emr-5.33.0-20221219	
	emr-5.33.0-20220802	524
	emr-5.33.0-20211008	
	emr-5.33.0-20210802	
	emr-5.33.0-20210615	
	emr-5.33.0-20210323	525
•	Versões 5.32.0	526
	emr-5.32.0-latest	
	emr-5.32.0-20240321	527
	emr-5.32.0-20220802	528
	emr-5.32.0-20211008	528
	emr-5.32.0-20210802	528

	emr-5.32.0-20210615	528
	emr-5.32.0-20210129	529
	emr-5.32.0-20201218	529
	emr-5.32.0-20201201	529
Histór	ico do documentos	530
	d>	αxxiii

O que é o Amazon EMR no EKS?

O Amazon EMR no EKS oferece uma opção de implantação para o Amazon EMR que permite executar estruturas de big data de código aberto no Amazon Elastic Kubernetes Service (Amazon EKS). Com essa opção de implantação, você pode se concentrar na execução de workloads de analytics enquanto o Amazon EMR no EKS cria, configura e gerencia contêineres para aplicações de código aberto.

Se você já usa o Amazon EMR, agora poderá executar aplicações baseadas no Amazon EMR com outros tipos de aplicações no mesmo cluster do Amazon EKS. Esta opção de implantação também melhora a utilização de recursos e simplifica o gerenciamento da infraestrutura em diversas zonas de disponibilidade. Se você já executa estruturas de big data no Amazon EKS, poderá usar o Amazon EMR para automatizar o provisionamento e o gerenciamento, e executar o Apache Spark mais rapidamente.

O Amazon EMR no EKS possibilita que sua equipe colabore com mais eficiência e processe grandes quantidades de dados com mais facilidade e economia:

- Você pode executar aplicações em um grupo comum de recursos sem a necessidade de provisionar infraestrutura. Você pode usar o <u>Amazon EMR Studio</u> e o AWS SDK ou AWS CLI desenvolver, enviar e diagnosticar aplicativos de análise executados em clusters EKS. É possível executar trabalhos programados no Amazon EMR no EKS usando o Apache Airflow autogerenciado ou o Amazon Managed Workflows for Apache Airflow (MWAA).
- As equipes de infraestrutura podem gerenciar centralmente uma plataforma de computação comum para consolidar as workloads do Amazon EMR com outras aplicações baseadas em contêiner. Você pode simplificar o gerenciamento da infraestrutura com ferramentas comuns do Amazon EKS e aproveitar um cluster compartilhado para workloads que precisam de versões diferentes de estruturas de código aberto. Também é possível reduzir a sobrecarga operacional com o gerenciamento automatizado de cluster do Kubernetes e com a aplicação de patches para o sistema operacional. Com a Amazon EC2 e AWS Fargate, você pode habilitar vários recursos computacionais para atender aos requisitos operacionais, financeiros ou de desempenho.

O diagrama a seguir mostra os dois diferentes modelos de implantação do Amazon EMR.

Tópicos

Arquitetura do Amazon EMR no EKS

1

- Noções básicas de conceitos e terminologia do Amazon EMR no EKS
- O que acontece quando você envia trabalhos para um cluster virtual do Amazon EMR no EKS

Arquitetura do Amazon EMR no EKS

O Amazon EMR no EKS realiza um acoplamento fraco das aplicações à infraestrutura na qual elas são executadas. Cada camada de infraestrutura fornece orquestração para a camada subsequente. Ao enviar um trabalho para o Amazon EMR, a definição do trabalho contém todos os parâmetros específicos da aplicação. O Amazon EMR usa esses parâmetros para instruir o Amazon EKS sobre quais pods e contêineres implantar. O Amazon EKS então coloca on-line os recursos de computação da Amazon EC2 e AWS Fargate necessários para executar o trabalho.

Com esse acoplamento fraco de serviços, você pode executar simultaneamente múltiplos trabalhos isolados com segurança. Também é possível realizar uma avalição comparativa do mesmo trabalho com diferentes back-ends de computação ou distribuir o trabalho em diversas zonas de disponibilidade para melhorar a disponibilidade.

O diagrama a seguir ilustra como o Amazon EMR no EKS funciona com AWS outros serviços.

Noções básicas de conceitos e terminologia do Amazon EMR no EKS

O Amazon EMR no EKS oferece uma opção de implantação para o Amazon EMR que permite executar estruturas de big data de código aberto no Amazon Elastic Kubernetes Service (Amazon EKS). Este tópico fornece contexto sobre algumas das terminologias comuns para isso, incluindo namespaces, clusters virtuais e execuções de trabalhos, que são unidades de trabalho enviadas para processamento.

Namespace do Kubernetes

O Amazon EKS usa namespaces do Kubernetes para dividir os recursos de cluster entre diversos usuários e aplicações. Esses namespaces são a base para ambientes multilocatários. Um namespace Kubernetes pode ter a Amazon EC2 ou AWS Fargate o provedor de computação. Essa flexibilidade oferece diferentes opções de performance e de custos para a execução de seus trabalhos.

Cluster virtual

Um cluster virtual corresponde a um namespace do Kubernetes no qual o Amazon EMR está registrado. O Amazon EMR usa clusters virtuais para executar trabalhos e hospedar endpoints. Diversos clusters virtuais podem ser apoiados pelo mesmo cluster físico. No entanto, cada cluster virtual é mapeado para um namespace em um cluster do EKS. Os clusters virtuais não criam quaisquer recursos ativos que contribuam para o seu faturamento ou que requeiram gerenciamento do ciclo de vida de forma externa ao serviço.

Execução de trabalho

A execução de um trabalho é uma unidade de trabalho, como um Spark jar, PySpark script ou consulta SparkSQL, que você envia para o Amazon EMR no EKS. Um trabalho pode ter várias execuções de trabalhos. Ao enviar uma execução de trabalho, você inclui as seguintes informações:

- Um cluster virtual no qual o trabalho deve ser executado.
- Um nome do trabalho para a identificação do trabalho.
- O perfil de execução, que é um perfil do IAM com escopo definido que executa o trabalho e permite especificar quais recursos podem ser acessados pelo trabalho.
- O rótulo de versão do Amazon EMR que especifica a versão das aplicações de código aberto a serem usadas.
- Os artefatos a serem usados ao enviar seu trabalho, como os parâmetros spark-submit.

Por padrão, os logs são carregados no servidor de histórico do Spark e podem ser acessados do AWS Management Console. Você também pode enviar registros de eventos, registros de execução e métricas para o Amazon S3 e a Amazon. CloudWatch

Contêineres do Amazon EMR

Os contêineres do Amazon EMR correspondem ao <u>nome da API do Amazon EMR no EKS</u>. O prefixo emr-containers é usado nos seguintes cenários:

- É o prefixo nos comandos da CLI para o Amazon EMR no EKS. Por exemplo, .aws emrcontainers start-job-run
- É o prefixo antes das ações de política do IAM para o Amazon EMR no EKS. Por exemplo, ."Action": ["emr-containers:StartJobRun"] Para obter mais informações, consulte Ações de políticas para o Amazon EMR no EKS.

Cluster virtual 3

• É o prefixo usado nos endpoints de serviço do Amazon EMR no EKS. Por exemplo, .emr-containers.us-east-1.amazonaws.com Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

O que acontece quando você envia trabalhos para um cluster virtual do Amazon EMR no EKS

O registro do Amazon EMR com um namespace do Kubernetes no Amazon EKS cria um cluster virtual. O Amazon EMR passará a poder executar workloads de analytics nesse namespace. Quando você usa o Amazon EMR no EKS para enviar trabalhos do Spark ao cluster virtual, o Amazon EMR no EKS solicita que o programador do Kubernetes no Amazon EKS programe pods.

As seguintes etapas e o diagrama abaixo ilustram o fluxo de trabalho do Amazon EMR no EKS:

- Uso de um cluster existente do Amazon EKS ou criação um cluster ao usar o utilitário de linha de comando eksctl ou o console do Amazon EKS.
- Criação de um cluster virtual ao registrar o Amazon EMR com um namespace em um cluster do EKS.
- Envie seu trabalho para o cluster virtual usando o AWS CLI ou SDK.

Para cada trabalho executado, o Amazon EMR no EKS cria um contêiner com uma imagem base do Amazon Linux 2, o Apache Spark e as dependências associadas. Cada execução de trabalho ocorre em um pod que faz download do contêiner e começa a executá-lo. O pod é encerrado após o término do trabalho. Se a imagem de contêiner tiver sido implantada anteriormente no nó, uma imagem em cache será usada e o download será ignorado. Os contêineres sidecar, como encaminhadores de log ou de métricas, podem ser implantados no pod. Após o término do trabalho, você ainda poderá depurá-lo usando a interface do usuário da aplicação do Spark no console do Amazon EMR.

Conceitos básicos do Amazon EMR no EKS

Este tópico ajuda você a começar a usar o Amazon EMR no EKS ao implantar uma aplicação do Spark em um cluster virtual. Inclui etapas para configurar as permissões corretas e iniciar um trabalho. Antes de começar, siga as etapas em Configuração do Amazon EMR no EKS. Isso ajuda você a obter ferramentas como a AWS CLI configurada antes de criar seu cluster virtual. Para outros modelos que podem ajudar você a começar, consulte nosso Guia de práticas recomendadas de contêineres do EMR em. GitHub

Você precisará das seguintes informações para as etapas de configuração:

 O ID do cluster virtual para o cluster do Amazon EKS e o namespace do Kubernetes registrado no Amazon EMR.

Important

Ao criar um cluster do EKS, certifique-se de usar m5.xlarge como tipo de instância ou qualquer outro tipo de instância com CPU e memória superiores. Usar um tipo de instância com CPU ou memória inferior ao m5.xlarge, pode levar à falha do trabalho devido à insuficiência de recursos disponíveis no cluster.

- O nome do perfil do IAM usado para a execução do trabalho.
- O rótulo de versão da versão do Amazon EMR (por exemplo, emr-6.4.0-latest).
- As metas de destino para o registro em log e o monitoramento:
 - Nome do grupo de CloudWatch log da Amazon e prefixo do fluxo de registros
 - A localização do Amazon S3 para armazenar logs de eventos e de contêineres.

Important

As tarefas do Amazon EMR no EKS usam a Amazon CloudWatch e o Amazon S3 como destinos para monitoramento e registro. Você pode monitorar o progresso do trabalho e solucionar falhas ao visualizar os logs de trabalho enviados para esses destinos. Para ativar o registro em log, a política do IAM associada ao perfil do IAM para a execução do trabalho deve ter as permissões necessárias para acessar os recursos das metas. Se a política do IAM não tiver as permissões necessárias, você deverá seguir as etapas descritas em Atualização da política de confiança do perfil de execução de trabalho Configurar uma

execução de trabalho para usar registros do Amazon S3 e Configurar uma execução de trabalho para CloudWatch usar registros antes de executar esse exemplo de trabalho.

Execução de uma aplicação do Spark

Siga as etapas a seguir para executar uma aplicação simples do Spark no Amazon EMR no EKS. O arquivo da aplicação entryPoint para a aplicação em Python do Spark está localizado em s3://REGION.elasticmapreduce/emr-containers/samples/wordcount/scripts/wordcount.py. REGIONÉ a região na qual seu cluster virtual Amazon EMR no EKS reside, como. us-east-1

 Atualize a política do IAM para o perfil de execução de trabalho com as permissões necessárias, conforme demonstram as instruções de política a seguir.

JSON

```
"Version": "2012-10-17",
"Statement": [
    "Sid": "ReadFromLoggingAndInputScriptBuckets",
    "Effect": "Allow",
    "Action": [
      "s3:GetObject",
      "s3:ListBucket"
    ],
    "Resource": [
      "arn:aws:s3:::*.elasticmapreduce",
      "arn:aws:s3:::*.elasticmapreduce/*",
      "arn:aws:s3:::amzn-s3-demo-bucket",
      "arn:aws:s3:::amzn-s3-demo-bucket/*",
      "arn:aws:s3:::amzn-s3-demo-bucket-b",
      "arn:aws:s3:::amzn-s3-demo-bucket-b/*"
    ]
 },
    "Sid": "WriteToLoggingAndOutputDataBuckets",
    "Effect": "Allow",
    "Action": [
```

```
"s3:PutObject",
        "s3:DeleteObject"
      ],
      "Resource": [
        "arn:aws:s3:::amzn-s3-demo-bucket/*",
        "arn:aws:s3:::amzn-s3-demo-bucket-b/*"
     ]
   },
      "Sid": "DescribeAndCreateCloudwatchLogStream",
      "Effect": "Allow",
      "Action": [
        "logs:CreateLogStream",
        "logs:DescribeLogGroups",
        "logs:DescribeLogStreams"
      ],
      "Resource": [
        "arn:aws:logs:*:*:*"
      ]
   },
      "Sid": "WriteToCloudwatchLogs",
      "Effect": "Allow",
      "Action": [
        "logs:PutLogEvents"
      ],
      "Resource": [
        "arn:aws:logs:*:*:log-group:my_log_group_name:log-
stream:my_log_stream_prefix/*"
      ]
   }
 ]
}
```

- A primeira instrução ReadFromLoggingAndInputScriptBuckets nesta política concede acesso para ListBucket e GetObjects aos seguintes buckets do Amazon S3:
 - *REGION*. elasticmapreduce: o bucket em que o arquivo da aplicação entryPoint está localizado.
 - amzn-s3-demo-destination-bucket- um bucket que você define para seus dados de saída.

- amzn-s3-demo-logging-bucket- um bucket que você define para seus dados de registro.
- A segunda instrução WriteToLoggingAndOutputDataBuckets nesta política concede ao trabalho as permissões para gravar dados em seus buckets de saída e de registro em log, respectivamente.
- A terceira declaração DescribeAndCreateCloudwatchLogStream concede ao trabalho permissões para descrever e criar Amazon CloudWatch Logs.
- A quarta declaração WriteToCloudwatchLogs concede permissões para gravar registros em um grupo de CloudWatch registros da Amazon chamado my_log_group_name em um fluxo de registros chamadomy_log_stream_prefix.
- Para executar uma aplicação em Python do Spark, use o comando apresentado a seguir.
 Substitua todos os red italicized valores substituíveis por valores apropriados. REGIONÉ a região na qual seu cluster virtual Amazon EMR no EKS reside, como. us-east-1

```
aws emr-containers start-job-run \
--virtual-cluster-id cluster_id
--name sample-job-name \
--execution-role-arn execution-role-arn \
--release-label emr-6.4.0-latest \
--job-driver '{
  "sparkSubmitJobDriver": {
    "entryPoint": "s3://REGION.elasticmapreduce/emr-containers/samples/wordcount/
scripts/wordcount.py",
    "entryPointArguments": ["s3://amzn-s3-demo-destination-bucket/
wordcount_output"],
    "sparkSubmitParameters": "--conf spark.executor.instances=2 --
conf spark.executor.memory=2G --conf spark.executor.cores=2 --conf
 spark.driver.cores=1"
}'\
--configuration-overrides '{
  "monitoringConfiguration": {
    "cloudWatchMonitoringConfiguration": {
      "logGroupName": "my_log_group_name",
      "logStreamNamePrefix": "my_log_stream_prefix"
    "s3MonitoringConfiguration": {
       "logUri": "s3://amzn-s3-demo-logging-bucket"
    }
  }
```

}'

Os dados de saída deste trabalho estarão disponíveis em s3://amzn-s3-demo-destination-bucket/wordcount_output.

Você também pode criar um arquivo JSON com parâmetros especificados para sua execução de trabalho. Em seguida, execute o comando start-job-run com um caminho para o arquivo JSON. Para obter mais informações, consulte Envio de uma execução de trabalho com StartJobRun. Para obter mais detalhes sobre como configurar parâmetros de execução de trabalho, consulte Opções para a configuração de uma execução de trabalho.

3. Para executar uma aplicação em SQL do Spark, use o comando apresentado a seguir. Substitua todos *red italicized* os valores pelos valores apropriados. *REGION*É a região na qual seu cluster virtual Amazon EMR no EKS reside, como. *us-east-1*

```
aws emr-containers start-job-run \
--virtual-cluster-id cluster_id \
--name sample-job-name \
--execution-role-arn execution-role-arn \
--release-label emr-6.7.0-latest \
--job-driver '{
  "sparkSqlJobDriver": {
    "entryPoint": "s3://query-file.sql",
    "sparkSqlParameters": "--conf spark.executor.instances=2 --
conf spark.executor.memory=2G --conf spark.executor.cores=2 --conf
 spark.driver.cores=1"
  }
}'\
--configuration-overrides '{
  "monitoringConfiguration": {
    "cloudWatchMonitoringConfiguration": {
      "logGroupName": "my_log_group_name",
      "logStreamNamePrefix": "my_log_stream_prefix"
    "s3MonitoringConfiguration": {
       "logUri": "s3://amzn-s3-demo-logging-bucket"
    }
  }
}'
```

Um arquivo de consulta SQL de exemplo é mostrado abaixo. Você deve ter um armazenamento de arquivos externo, como o S3, no qual os dados das tabelas são armazenados.

```
CREATE DATABASE demo;
CREATE EXTERNAL TABLE IF NOT EXISTS demo.amazonreview( marketplace string,
  customer_id string, review_id string, product_id string, product_parent string,
  product_title string, star_rating integer, helpful_votes integer, total_votes
  integer, vine string, verified_purchase string, review_headline string,
  review_body string, review_date date, year integer) STORED AS PARQUET LOCATION
  's3://URI to parquet files';
SELECT count(*) FROM demo.amazonreview;
SELECT count(*) FROM demo.amazonreview WHERE star_rating = 3;
```

A saída desse trabalho estará disponível nos registros de stdout do driver no S3 ou CloudWatch, dependendo do monitoringConfiguration que estiver configurado.

4. Você também pode criar um arquivo JSON com parâmetros especificados para sua execução de trabalho. Em seguida, execute o comando start-job-run com um caminho para o arquivo JSON. Para obter mais informações, consulte Envio de uma execução de trabalho. Para obter mais detalhes sobre como configurar parâmetros de execução de trabalho, consulte Opções para a configuração de uma execução de trabalho.

Para monitorar o progresso do trabalho ou para depurar falhas, você pode inspecionar os logs enviados para o Amazon S3 CloudWatch, Logs ou ambos. Consulte o caminho do log no Amazon S3 em Configurar uma execução de trabalho para usar os registros do S3 e para os registros do Cloudwatch em Configurar uma execução de trabalho para usar os registros. CloudWatch Para ver os registros no CloudWatch Logs, siga as instruções abaixo.

- Abra o CloudWatch console em https://console.aws.amazon.com/cloudwatch/.
- No painel de navegação, selecione Logs. Em seguida, escolha Grupos de logs.
- Escolha o grupo de logs do Amazon EMR no EKS e, em seguida, visualize os eventos de log carregados.



▲ Important

Os trabalhos têm uma política de repetição padrão configurada. Para obter informações sobre como modificar ou desabilitar a configuração, consulte <u>Uso de políticas de repetição de</u> trabalho.

Links para o Amazon EMR nos guias de melhores práticas do EKS sobre GitHub

Desenvolvemos o Guia de práticas recomendadas do Amazon EMR no EKS com a colaboração da comunidade de código aberto para que fosse possível iterar rapidamente e fornecer recomendações sobre aspectos de criação e execução de um cluster virtual. Recomendamos usar o Guia de práticas recomendadas do Amazon EMR no EKS para as seções. Escolha os links em cada seção para acessar o GitHub site.

Segurança



Note

Para obter mais informações sobre segurança com o Amazon EMR no EKS, consulte Práticas recomendadas de segurança para o Amazon EMR no EKS.

Práticas recomendadas de criptografia: descreve como usar a criptografia para dados em repouso e em trânsito.

O gerenciamento da segurança da rede descreve como configurar grupos de segurança para pods do Amazon EMR no EKS ao se conectar a fontes de dados hospedadas em Serviços da AWS, como o Amazon RDS e o Amazon Redshift.

Usando o gerenciador de AWS segredos para armazenar segredos.

Envio de trabalho para o Pyspark

Envio de trabalho para o Pyspark: especifica diferentes tipos de empacotamentos para as aplicações do PySpark usando formatos de empacotamento como zip, egg, wheel e pex.

Armazenamento

Uso de volumes do EBS: descreve como usar o provisionamento estático e dinâmico para trabalhos que precisam de volumes do EBS.

12 Segurança

<u>Usando volumes do Amazon FSx for Luster:</u> como usar o provisionamento estático e dinâmico para trabalhos que precisam de volumes do Amazon FSx for Luster.

<u>Uso de volumes de armazenamento de instância:</u> descreve como usar volumes de armazenamento de instância para o processamento de trabalhos.

Integração com o metastore

<u>Uso do Hive Metastore:</u> oferece diferentes maneiras de usar o Hive Metastore.

<u>Usando o AWS Glue:</u> oferece maneiras diferentes de configurar o catálogo do AWS Glue.

Depuração

Uso da depuração do Spark: descreve como alterar o nível de log.

Conexão com a interface do usuário do Spark no pod do driver.

Como usar o servidor de histórico do Spark de hospedagem própria com o Amazon EMR no EKS.

Solução de problemas do Amazon EMR no EKS

Solução de problemas.

Posicionamento de nó

<u>Uso de seletores de nó do Kubernetes</u> para single-az e outros casos de uso.

Uso do posicionamento de nó do Fargate.

Performance

Uso da alocação dinâmica de recursos (DRA).

<u>Práticas recomendadas do EKS</u> para o plug-in Container Network Interface (CNI) da Amazon VPC, o Cluster Autoscaler e o CoreDNS.

Integração com o metastore 13

Otimização de custo

<u>Usando instâncias spot:</u> melhores práticas de instâncias EC2 spot da Amazon e como usar o recurso de desativação de nós do Spark.

Usando AWS Outposts

Executando o Amazon EMR no EKS usando AWS Outposts

Otimização de custo 14

Personalização de imagens do Docker para o Amazon EMR no EKS

É possível usar imagens do Docker personalizadas com o Amazon EMR no EKS. A personalização da imagem de runtime do Amazon EMR no EKS oferece os seguintes benefícios:

- Empacotamento das dependências e do ambiente de runtime da aplicação em um único contêiner imutável que promove a portabilidade e simplifica o gerenciamento de dependências para cada workload.
- Instalação e configuração de pacotes otimizados para as workloads. Esses pacotes podem não estar amplamente disponíveis na distribuição pública dos runtimes do Amazon EMR.
- Integração do Amazon EMR no EKS aos processos atuais de criação, teste e implantação estabelecidos em sua organização, incluindo desenvolvimento e testes locais.
- Aplicação dos processos de segurança estabelecidos, como a verificação de imagens, que atendem aos requisitos de conformidade e governança da sua organização.

Tópicos

- Como personalizar imagens do Docker
- Detalhes sobre como selecionar um URI de imagem base
- Considerações para personalizar imagens

Como personalizar imagens do Docker

Siga as etapas a seguir para personalizar as imagens do Docker para o Amazon EMR no EKS. As etapas mostram como obter uma imagem base, personalizá-la e publicá-la e enviar uma workload usando a imagem.

- Pré-requisitos
- Etapa 1: recuperar uma imagem base do Amazon Elastic Container Registry (Amazon ECR)
- Etapa 2: personalizar uma imagem base
- Etapa 3: (opcional, mas recomendada) validar uma imagem personalizada
- Etapa 4: publicar uma imagem personalizada
- Etapa 5: enviar uma workload do Spark no Amazon EMR usando uma imagem personalizada



Note

Outras opções que você pode considerar ao personalizar imagens do Docker são a personalização para endpoints interativos, feita para garantir que tenha as dependências necessárias, ou o uso de imagens de contêiner com várias arquiteturas:

- Personalização de imagens do Docker para endpoints interativos
- Trabalho com imagens de múltiplas arquiteturas

Pré-requisitos

- Conclua as etapas em Configuração do Amazon EMR no EKS para o Amazon EMR no EKS.
- Instale o Docker em seu ambiente. Para obter mais informações, consulte Get Docker.

Etapa 1: recuperar uma imagem base do Amazon Elastic Container Registry (Amazon ECR)

A imagem base contém o runtime do Amazon EMR e os conectores usados para acessar outros serviços da AWS. Para a versão 6.9.0 e posteriores do Amazon EMR, você pode obter as imagens base na galeria pública do Amazon ECR. Navegue pela galeria para encontrar o link da imagem e extraia-a para seu Workspace local. Por exemplo, para a versão 7.9.0 do Amazon EMR, o docker pull comando a seguir fornece a última imagem base padrão. Você pode substituir emr-7.9.0:latest por emr-7.9.0-spark-rapids:latest para recuperar a imagem que tem o acelerador RAPIDS da Nvidia. Você também pode substituir emr-7.9.0:latest por emr-7.9.0-java11:latest para recuperar a imagem com o runtime do Java 11.

```
docker pull public.ecr.aws/emr-on-eks/spark/emr-7.9.0:latest
```

Se você desejar recuperar a imagem base para uma versão 6.9.0 ou para versões anteriores do Amazon EMR, ou se preferir recuperá-la de contas de registro do Amazon ECR em cada região, use as seguintes etapas:

Escolha um URI de imagem base. O URI de imagem segue este formato, *ECR-registry*account.dkr.ecr.Region.amazonaws.com/spark/container-image-tag, como demonstra o exemplo a seguir.

Pré-requisitos

```
895885662937.dkr.ecr.us-west-2.amazonaws.com/spark/emr-6.6.0:latest
```

Para escolher uma imagem base em sua região, consulte <u>Detalhes sobre como selecionar um</u> URI de imagem base.

 Faça login no repositório do Amazon ECR no qual a imagem base está armazenada. Substitua 895885662937 e us-west-2 pela conta de registro do Amazon ECR e pela AWS região que você selecionou.

```
aws ecr get-login-password --region us-west-2 | docker login --username AWS -- password-stdin 895885662937.dkr.ecr.us-west-2.amazonaws.com
```

 Extraia a imagem base para seu Workspace local. emr-6.6.0:latestSubstitua pela tag de imagem do contêiner que você selecionou.

```
docker pull 895885662937.dkr.ecr.us-west-2.amazonaws.com/spark/emr-6.6.0:latest
```

Etapa 2: personalizar uma imagem base

Siga as etapas a seguir para personalizar a imagem base extraída do Amazon ECR.

- 1. Crie um novo Dockerfile em seu Workspace local.
- Edite o Dockerfile que você acabou de criar e adicione o conteúdo a seguir. Este Dockerfile usa a imagem de contêiner que você extraiu de 895885662937.dkr.ecr.uswest-2.amazonaws.com/spark/emr-6.6.0:latest.

```
FROM 895885662937.dkr.ecr.us-west-2.amazonaws.com/spark/emr-6.6.0:latest USER root ### Add customization commands here #### USER hadoop:hadoop
```

 Adicione comandos no Dockerfile para personalizar a imagem base. Por exemplo, adicione um comando para instalar bibliotecas Python, como demonstra o Dockerfile a seguir.

```
FROM 895885662937.dkr.ecr.us-west-2.amazonaws.com/spark/emr-6.6.0:latest USER root
RUN pip3 install --upgrade boto3 pandas numpy // For python 3
USER hadoop:hadoop
```

4. No mesmo diretório em que o Dockerfile foi criado, execute o comando apresentado a seguir para criar a imagem do Docker. Forneça um nome para a imagem do Docker, por exemplo, emr6.6_custom.

```
docker build -t emr6.6_custom .
```

Etapa 3: (opcional, mas recomendada) validar uma imagem personalizada

Recomendamos testar a compatibilidade da sua imagem personalizada antes de publicá-la. Você pode usar a <u>CLI de imagem personalizada do Amazon EMR no EKS</u> para verificar se sua imagem tem as estruturas de arquivos necessárias e as configurações corretas para a execução no Amazon EMR no EKS.



A CLI de imagem personalizada do Amazon EMR no EKS não pode confirmar se sua imagem está isenta de erros. Tenha cuidado ao remover dependências das imagens base.

Siga as etapas apresentadas a seguir para validar a imagem personalizada.

- Faça download e instale a CLI de imagem personalizada do Amazon EMR no EKS. Para obter mais informações, consulte Amazon EMR on EKS custom image CLI Installation Guide.
- 2. Execute o comando apresentado a seguir para testar a instalação.

```
emr-on-eks-custom-image --version
```

Confira a seguir um exemplo da saída.

```
Amazon EMR on EKS Custom Image CLI
Version: x.xx
```

3. Execute o comando apresentado a seguir para validar a imagem personalizada.

```
emr-on-eks-custom-image validate-image -i image_name -r release_version [-
t image_type]
```

- -i especifica o URI da imagem local que precisa ser validado. Pode ser o URI da imagem, qualquer nome ou etiqueta que você definiu para a imagem.
- -r especifica a versão de liberação exata para a imagem base, por exemplo, emr-6.6.0latest.
- -t especifica o tipo de imagem. Se for uma imagem do Spark, insira spark. O valor padrão é spark. A versão atual da CLI de imagem personalizada do Amazon EMR no EKS oferece suporte somente para imagens de runtime do Spark.

Se você executar o comando com êxito e a imagem personalizada atender a todas as configurações e estruturas de arquivos necessárias, a saída retornada exibirá os resultados de todos os testes, como demonstra o exemplo a seguir.

```
Amazon EMR on EKS Custom Image Test
Version: x.xx
... Checking if docker cli is installed
... Checking Image Manifest
[INFO] Image ID: xxx
[INFO] Created On: 2021-05-17T20:50:07.986662904Z
[INFO] Default User Set to hadoop:hadoop : PASS
[INFO] Working Directory Set to /home/hadoop : PASS
[INFO] Entrypoint Set to /usr/bin/entrypoint.sh : PASS
[INFO] SPARK_HOME is set with value: /usr/lib/spark : PASS
[INFO] JAVA_HOME is set with value: /etc/alternatives/jre : PASS
[INFO] File Structure Test for spark-jars in /usr/lib/spark/jars: PASS
[INFO] File Structure Test for hadoop-files in /usr/lib/hadoop: PASS
[INFO] File Structure Test for hadoop-jars in /usr/lib/hadoop/lib: PASS
[INFO] File Structure Test for bin-files in /usr/bin: PASS
... Start Running Sample Spark Job
[INFO] Sample Spark Job Test with local:///usr/lib/spark/examples/jars/spark-
examples.jar : PASS
Overall Custom Image Validation Succeeded.
```

Se a imagem personalizada não atender às configurações ou estruturas de arquivos necessárias, você obterá mensagens de erro. A saída retornada fornece informações sobre configurações ou estruturas de arquivos incorretas.

Etapa 4: publicar uma imagem personalizada

Publique a nova imagem do Docker no registro do Amazon ECR.

 Execute o comando apresentado a seguir para criar um repositório do Amazon ECR para o armazenamento da imagem do Docker. Forneça um nome para seu repositório, por exemplo, emr6.6_custom_repo. Substitua us-west-2 pela sua região.

```
aws ecr create-repository \
    --repository-name emr6.6_custom_repo \
    --image-scanning-configuration scanOnPush=true \
    --region us-west-2
```

Para obter mais informações, consulte Criar um repositório no Guia do usuário do Amazon ECR.

2. Execute o comando apresentado a seguir para realizar a autenticação em seu registro padrão.

```
aws ecr get-login-password --region us-west-2 | docker login --username AWS -- password-stdin aws_account_id.dkr.ecr.us-west-2.amazonaws.com
```

Para obter mais informações, consulte <u>Autenticar-se no registro padrão</u> no Guia do usuário do Amazon ECR.

3. Faça a marcação e publique uma imagem no repositório do Amazon ECR que você criou.

Marque a imagem.

```
docker tag emr6.6_custom aws_account_id.dkr.ecr.us-
west-2.amazonaws.com/emr6.6_custom_repo
```

Envie a imagem.

```
docker push aws_account_id.dkr.ecr.us-west-2.amazonaws.com/emr6.6_custom_repo
```

Para obter mais informações, consulte <u>Enviar uma imagem ao Amazon ECR</u> no Guia do usuário do Amazon ECR.

Etapa 5: enviar uma workload do Spark no Amazon EMR usando uma imagem personalizada

Depois que uma imagem personalizada for criada e publicada, será possível enviar um trabalho do Amazon EMR no EKS usando uma imagem personalizada.

Primeiro, crie um start-job-run-request arquivo.json e especifique o spark.kubernetes.container.image parâmetro para referenciar a imagem personalizada, como demonstra o exemplo de arquivo JSON a seguir.

Note

Você pode usar o esquema local:// para se referir aos arquivos disponíveis na imagem personalizada, conforme mostrado com o argumento entryPoint no trecho de código JSON abaixo. Você também pode usar o esquema local:// para se referir às dependências da aplicação. Todos os arquivos e as dependências referenciados usando o esquema local:// já devem estar presentes no caminho especificado na imagem personalizada.

```
{
    "name": "spark-custom-image",
    "virtualClusterId": "virtual-cluster-id",
    "executionRoleArn": "execution-role-arn",
    "releaseLabel": "emr-6.6.0-latest",
    "jobDriver": {
      "sparkSubmitJobDriver": {
        "entryPoint": "local:///usr/lib/spark/examples/jars/spark-examples.jar",
        "entryPointArguments": [
                  "10"
              ],
         "sparkSubmitParameters": "--class org.apache.spark.examples.SparkPi --conf
 spark.kubernetes.container.image=123456789012.dkr.ecr.us-west-2.amazonaws.com/
emr6.6_custom_repo"
       }
    }
}
```

Você também pode fazer referência à imagem personalizada com propriedades applicationConfiguration, como demonstra o exemplo a seguir.

```
{
    "name": "spark-custom-image",
    "virtualClusterId": "virtual-cluster-id",
    "executionRoleArn": "execution-role-arn",
    "releaseLabel": "emr-6.6.0-latest",
    "jobDriver": {
      "sparkSubmitJobDriver": {
        "entryPoint": "local:///usr/lib/spark/examples/jars/spark-examples.jar",
        "entryPointArguments": [
                  "10"
              ],
         "sparkSubmitParameters": "--class org.apache.spark.examples.SparkPi"
       }
    },
    "configurationOverrides": {
        "applicationConfiguration": [
            {
                "classification": "spark-defaults",
                "properties": {
                    "spark.kubernetes.container.image": "123456789012.dkr.ecr.us-
west-2.amazonaws.com/emr6.6_custom_repo"
                }
            }
        ]
    }
}
```

Em seguida, execute o comando start-job-run para enviar o trabalho.

```
aws emr-containers start-job-run --cli-input-json file://./start-job-run-request.json
```

Nos exemplos de JSON acima, *emr-6.6.0-latest* substitua pela sua versão de lançamento do Amazon EMR. Recomendamos fortemente usar a versão de liberação -latest para garantir que a versão selecionada contenha as atualizações de segurança mais recentes. Para obter mais informações sobre as versões de liberação do Amazon EMR e suas etiquetas de imagem, consulte Detalhes sobre como selecionar um URI de imagem base.



Note

É possível usar spark.kubernetes.driver.container.image e spark.kubernetes.executor.container.image para especificar uma imagem diferente para pods de drivers e de executores.

Personalização de imagens do Docker para endpoints interativos

Você também pode personalizar as imagens do Docker para endpoints interativos com a finalidade de executar imagens base de kernel personalizadas. Isso ajuda a garantir que você tenha as dependências necessárias ao executar workloads interativas do EMR Studio.

Siga as etapas 1 a 4 descritas acima para personalizar uma imagem do Docker. Para versões 6.9.0 e posteriores do Amazon EMR, é possível obter o URI da imagem base na galeria pública do Amazon ECR. Para versões anteriores ao Amazon EMR 6.9.0, você pode obter a imagem nas contas de registro do Amazon ECR em cada Região da AWS, e a única diferença é o URI da imagem base em seu Dockerfile. O URI da imagem base segue o formato:

ECR-registry-account.dkr.ecr.Region.amazonaws.com/notebook-spark/container-imagetag

Você precisa usar notebook-spark no URI da imagem base, em vez de spark. A imagem base contém o runtime do Spark e os kernels do caderno que são executados com ele. Para obter mais informações sobre como selecionar as regiões e as etiquetas de imagem de contêiner, consulte Detalhes sobre como selecionar um URI de imagem base.



Note

Atualmente, somente substituições de imagens básicas são suportadas e a introdução de kernels completamente novos de outros tipos além dos fornecidos pelas imagens AWS básicas não é suportada.

2. Crie um endpoint interativo que possa ser usado com a imagem personalizada.

Primeiro, crie um arquivo JSON chamado custom-image-managed-endpoint.json com o conteúdo apresentado a seguir.

```
{
    "name": "endpoint-name",
    "virtualClusterId": "virtual-cluster-id",
    "type": "JUPYTER_ENTERPRISE_GATEWAY",
    "releaseLabel": "emr-6.6.0-latest",
    "executionRoleArn": "execution-role-arn",
    "certificateArn": "certificate-arn",
    "configurationOverrides": {
        "applicationConfiguration": [
            {
                "classification": "jupyter-kernel-overrides",
                "configurations": [
                    {
                         "classification": "python3",
                         "properties": {
                             "container-image": "123456789012.dkr.ecr.us-
west-2.amazonaws.com/custom-notebook-python:latest"
                        }
                    },
                    {
                         "classification": "spark-python-kubernetes",
                         "properties": {
                             "container-image": "123456789012.dkr.ecr.us-
west-2.amazonaws.com/custom-notebook-spark:latest"
                        }
                    }
                ]
            }
        ]
    }
}
```

Em seguida, crie um endpoint interativo usando as configurações especificadas no arquivo JSON, como demonstra o exemplo a seguir.

```
aws emr-containers create-managed-endpoint --cli-input-json custom-image-managed-endpoint.json
```

Para obter mais informações, consulte Criação de um endpoint interativo para o cluster virtual.

3. Conecte-se ao endpoint interativo usando o EMR Studio. Para obter mais informações, consulte Connecting from Studio.

Trabalho com imagens de múltiplas arquiteturas

O Amazon EMR no EKS oferece suporte a imagens de contêiner de múltiplas arquiteturas para o Amazon Elastic Container Registry (Amazon ECR). Para obter mais informações, consulte Introducing multi-architecture container images for Amazon ECR.

As imagens personalizadas do Amazon EMR no EKS oferecem suporte tanto para instâncias AWS baseadas em Graviton EC2 quanto para instâncias. non-Graviton-based EC2 As imagens baseadas em Graviton são armazenadas nos mesmos repositórios de imagens no Amazon ECR que as imagens. non-Graviton-based

Por exemplo, para inspecionar a lista de manifestos do Docker para imagens 6.6.0, execute o comando a seguir.

```
docker manifest inspect 895885662937.dkr.ecr.us-west-2.amazonaws.com/spark/
emr-6.6.0:latest
```

Confira a saída a seguir. A arquitetura arm64 é para instâncias baseadas no Graviton. A arquitetura amd64 é para instâncias não baseadas no Graviton.

```
{
   "schemaVersion": 2,
   "mediaType": "application/vnd.docker.distribution.manifest.list.v2+json",
   "manifests": [
      {
         "mediaType": "application/vnd.docker.distribution.manifest.v2+json",
         "size": 1805,
         "digest":
 "xxx123:6b971cb47d11011ab3d45fff925e9442914b4977ae0f9fbcdcf5cfa99a7593f0",
         "platform": {
            "architecture": "arm64",
            "os": "linux"
         }
      },
         "mediaType": "application/vnd.docker.distribution.manifest.v2+json",
         "size": 1805,
```

```
"digest":
 "xxx123:6f2375582c9c57fa9838c1d3a626f1b4fc281e287d2963a72dfe0bd81117e52f",
         "platform": {
            "architecture": "amd64",
            "os": "linux"
         }
      }
   ]
}
```

Execute as seguintes etapas para criar imagens de múltiplas arquiteturas:

Crie um Dockerfile com o conteúdo apresentado a seguir para que você possa extrair a imagem arm64.

```
FROM --platform=arm64 895885662937.dkr.ecr.us-west-2.amazonaws.com/spark/
emr-6.6.0:latest
USER root
RUN pip3 install boto3 // install customizations here
USER hadoop:hadoop
```

2. Siga as instruções em Introducing multi-architecture container images for Amazon ECR para desenvolver uma imagem com múltiplas arquiteturas.



Note

Você deve criar imagens arm64 em instâncias arm64. De forma semelhante, você deve criar imagens amd64 em instâncias amd64.

Também é possível desenvolver imagens com múltiplas arquiteturas sem a necessidade de criar cada tipo de instância específico com o comando buildx do Docker. Para obter mais informações, consulte Leverage multi-CPU architecture support.

Após criar uma imagem com múltiplas arquiteturas, você poderá enviar um trabalho com o mesmo parâmetro spark.kubernetes.container.image e direcioná-lo para a imagem. Em um cluster heterogêneo com instâncias AWS baseadas em Graviton e non-Graviton-based EC2 instâncias, a instância determina a imagem correta da arquitetura com base na arquitetura da instância que extrai a imagem.

Detalhes sobre como selecionar um URI de imagem base



Note

Para versões 6.9.0 e posteriores do Amazon EMR, é possível recuperar a imagem base da galeria pública do Amazon ECR, portanto, não é necessário criar o URI da imagem base conforme as instruções nesta página. Para encontrar a etiqueta da imagem de contêiner para sua imagem base, consulte a página de notas de versão da versão correspondente do Amazon EMR no EKS.

As imagens base do Docker que você pode selecionar são armazenadas no Amazon Elastic Container Registry (Amazon ECR). O URI da imagem segue este formato: ECR-registryaccount.dkr.ecr.Region.amazonaws.com/spark/container-image-tag,como demonstra o seguinte exemplo a seguir.

```
895885662937.dkr.ecr.us-west-2.amazonaws.com/spark/emr-7.9.0:latest
```

O URI da imagem para endpoints interativos segue este formato: *ECR-registry*account.dkr.ecr.Region.amazonaws.com/notebook-spark/container-image-tag, como demonstra o exemplo a seguir. Você precisa usar notebook-spark no URI da imagem base, em vez de spark.

```
895885662937.dkr.ecr.us-west-2.amazonaws.com/notebook-spark/emr-7.9.0:latest
```

De forma semelhante, para imagens python3 que não sejam do Spark para endpoints interativos, o URI da imagem é ECR-registry-account.dkr.ecr.Region.amazonaws.com/notebookpython/container-image-tag. O seguinte URI de exemplo está formatado corretamente:

```
895885662937.dkr.ecr.us-west-2.amazonaws.com/notebook-python/emr-7.9.0:latest
```

Para encontrar a etiqueta da imagem de contêiner para sua imagem base, consulte a página de notas de versão da versão correspondente do Amazon EMR no EKS.

Contas de registro do Amazon ECR por região

Para evitar a alta latência da rede, extraia uma imagem base da imagem mais próxima Região da AWS. Selecione a conta de registro do Amazon ECR que corresponde à região da qual você extrairá a imagem com base na tabela a seguir.

Regiões	Contas de registro do Amazon ECR	
ap-east-1	736135916053	
ap-northeast-1	059004520145	
ap-northeast-2	996579266876	
ap-northeast-3	705689932349	
ap-southeast-3	946962994502	
ap-south-1	235914868574	
ap-south-2	691480105545	
ap-southeast-1	671219180197	
ap-southeast-2	038297999601	
ca-central-1	351826393999	
eu-central-1	107292555468	
eu-central-2	314408114945	
eu-north-1	830386416364	
eu-west-1	483788554619	
eu-west-2	118780647275	
eu-west-3	307523725174	

Regiões	Contas de registro do Amazon ECR
eu-south-1	238014973495
eu-south-2	350796622945
il-central-1	395734710648
me-south-1	008085056818
me-central-1	818935616732
sa-east-1	052806832358
us-gov-west-1	299385240661
us-gov-east-1	299393998622
us-east-1	755674844232
us-east-2	711395599931
us-west-1	608033475327
us-west-2	895885662937
af-south-1	358491847878
cn-north-1	068337069695
cn-northwest-1	068420816659

Considerações para personalizar imagens

Ao personalizar imagens do Docker, é possível escolher o runtime exato para seu trabalho em um nível granular. Considere estas práticas recomendadas ao usar esse recurso. Isso inclui considerações sobre segurança, configuração e montagem de uma imagem:

- A segurança é uma responsabilidade compartilhada entre você AWS e você. Você é responsável
 pela aplicação de patches de segurança nos binários adicionados à imagem. Siga as <u>Práticas</u>
 recomendadas de segurança para o Amazon EMR no EKS, especialmente o que é colocado
 em <u>Obtenção das atualizações de segurança mais recentes para as imagens personalizadas</u> e
 Aplicação do princípio de privilégio mínimo.
- Ao personalizar uma imagem base, você deve alterar o usuário do Docker para hadoop: hadoop com a finalidade de que os trabalhos não sejam executados com o usuário raiz.
- O Amazon EMR no EKS monta arquivos sobre as configurações para as imagens, como o spark-defaults.conf, no runtime. Para substituir esses arquivos de configuração, recomendamos usar o parâmetro applicationOverrides durante o envio do trabalho e não modificar diretamente os arquivos na imagem personalizada.
- O Amazon EMR no EKS monta determinadas pastas no runtime. Quaisquer modificações feitas nessas pastas não estarão disponíveis no contêiner. Se você desejar adicionar uma aplicação ou suas dependências para as imagens personalizadas, recomendamos escolher um diretório que não faça parte dos seguintes caminhos definidos previamente:
 - /var/log/fluentd
 - /var/log/spark/user
 - /var/log/spark/apps
 - /mnt
 - /tmp
 - /home/hadoop
- Você pode fazer upload de sua imagem personalizada para qualquer repositório compatível com Docker, como o Amazon ECR, o Docker Hub ou um repositório empresarial privado. Para obter mais informações sobre como configurar a autenticação de cluster do Amazon EKS com o repositório do Docker selecionado, consulte Pull an Image from a Private Registry.

Execução de trabalhos do Flink com o Amazon EMR no EKS

As versões 6.13.0 e superiores do Amazon EMR oferecem suporte ao Amazon EMR no EKS com o Apache Flink ou ao operador do Kubernetes para Flink, como um modelo de envio de trabalhos para o Amazon EMR no EKS. Com o Amazon EMR no EKS com o Apache Flink, você pode implantar e gerenciar aplicações do Flink com o runtime da versão do Amazon EMR em seus próprios clusters do Amazon EKS. Depois de implantar o operador do Kubernetes para Flink em seu cluster do Amazon EKS, você poderá enviar aplicações do Flink diretamente com o operador. O operador gerencia o ciclo de vida das aplicações do Flink.

Tópicos

- Configuração e uso do operador do Kubernetes do Flink
- Uso do Kubernetes nativo do Flink
- Personalização de imagens do Docker para o Flink e o FluentD
- Monitoramento do operador do Kubernetes para Flink e dos trabalhos do Flink
- Como o Flink oferece suporte à alta disponibilidade e resiliência do trabalho
- Uso do Autoscaler para aplicações do Flink
- Manutenção e solução de problemas para trabalhos do Flink no Amazon EMR no EKS
- Versões compatíveis para Amazon EMR no EKS com Apache Flink

Configuração e uso do operador do Kubernetes do Flink

As páginas apresentadas a seguir descrevem como configurar e usar o operador do Kubernetes para Flink com a finalidade de executar trabalhos do Flink com o Amazon EMR no EKS. Os tópicos disponíveis incluem pré-requisitos obrigatórios, como configurar seu ambiente e executar uma aplicação do Flink no Amazon EMR no EKS.

Tópicos

- Configuração do operador do Kubernetes para Flink para o Amazon EMR no EKS
- Desinstalação do operador do Kubernetes do Flink no Amazon EMR no EKS
- Execução de uma aplicação do Flink
- Permissões do perfil de segurança para executar uma aplicação do Flink

Desinstalação do operador do Kubernetes para Flink para o Amazon EMR no EKS

Configuração do operador do Kubernetes para Flink para o Amazon EMR no EKS

Conclua as tarefas apresentadas a seguir para se preparar antes de instalar o operador do Kubernetes para Flink no Amazon EKS. Se você já se inscreveu na Amazon Web Services (AWS) e usou o Amazon EKS, está com quase tudo pronto para usar o Amazon EMR no EKS. Conclua as tarefas apresentadas a seguir para se preparar para usar o operador do Flink no Amazon EKS. Se você já completou algum dos pré-requisitos, pode ignorá-los e passar para os próximos.

- <u>Instale ou atualize para a versão mais recente do AWS CLI</u> Se você já instalou o AWS CLI, confirme se você tem a versão mais recente.
- <u>Instale kubectl e eksctl</u>: o eksctl é uma ferramenta de linha de comando usada para se comunicar com o Amazon EKS.
- <u>Instale o Helm</u>: o gerenciador de pacotes Helm para o Kubernetes ajuda a instalar e gerenciar aplicações em seu cluster do Kubernetes.
- Comece a usar o Amazon EKS (eksctl): siga as etapas para criar um cluster do Kubernetes com nós no Amazon EKS.
- Escolha um rótulo de lançamento do Amazon EMR (versão 6.13.0 ou superiores): o operador do Kubernetes do Flink é compatível com as versões 6.13.0 e superiores do Amazon EMR.
- Habilite perfis do IAM para contas de serviço (IRSA) no cluster do Amazon EKS.
- Crie um perfil de execução de trabalho.
- Atualize a política de confiança do perfil de execução de trabalho.
- Crie um perfil de execução de operador. Esta etapa é opcional. É possível usar o mesmo perfil
 para trabalhos e operadores do Flink. Se desejar ter um perfil do IAM diferente para o operador,
 você poderá criar um perfil separado.
- Atualize a política de confiança do perfil de execução do operador. Você deve adicionar explicitamente uma entrada de política de confiança para os perfis que deseja usar para a conta de serviço do operador do Kubernetes para Flink no Amazon EMR. Você pode seguir este formato de exemplo:

JSON

{

Configuração 32

```
"Version": "2012-10-17",
  "Statement": [
      "Effect": "Allow",
      "Action": [
        "sts:AssumeRoleWithWebIdentity"
      ],
      "Resource": [
        11 * 11
      ],
      "Condition": {
        "StringLike": {
          "aws:userid": "system:serviceaccount:emr:emr-containers-sa-flink-
operator"
        }
      "Sid": "AllowSTSAssumerolewithwebidentity"
    }
  ]
}
```

Desinstalação do operador do Kubernetes do Flink no Amazon EMR no EKS

Este tópico ajuda você a começar a usar o operador do Kubernetes do Flink no Amazon EKS ao preparar uma implantação do Flink.

Como instalar o operador do Kubernetes

Use as etapas a seguir para instalar o operador do Kubernetes para Apache Flink.

- 1. Caso ainda não tenha feito, conclua as etapas em the section called "Configuração".
- 2. Instale o *cert-manager* (uma vez por cluster do Amazon EKS) para permitir a adição do componente webhook.

```
kubectl apply -f https://github.com/cert-manager/cert-manager/releases/download/
v1.12.0/cert-manager.yaml
```

3. Instale o chart do Helm.

```
export VERSION=7.9.0 # The Amazon EMR release version
```

```
export NAMESPACE=The Kubernetes namespace to deploy the operator
helm install flink-kubernetes-operator \
oci://public.ecr.aws/emr-on-eks/flink-kubernetes-operator \
--version $VERSION \
--namespace $NAMESPACE
```

Resultado do exemplo:

NAME: flink-kubernetes-operator

LAST DEPLOYED: Tue May 31 17:38:56 2022

NAMESPACE: \$NAMESPACE

STATUS: deployed

REVISION: 1

TEST SUITE: None

4. Aguarde a conclusão da implantação e verifique a instalação do chart.

```
kubectl wait deployment flink-kubernetes-operator --namespace $NAMESPACE --for
condition=Available=True --timeout=30s
```

5. Você deverá visualizar a mensagem a seguir quando a implantação for concluída.

```
deployment.apps/flink-kubernetes-operator condition met
```

6. Use o comando apresentado a seguir para visualizar o operador implantado.

```
helm list --namespace $NAMESPACE
```

O exemplo a seguir mostra um exemplo de saída, em que a versão da aplicação x.y.z-amznn corresponderia à versão do operador do Flink para seu Amazon EMR na versão EKS. Para obter mais informações, consulte <u>Versões compatíveis para Amazon EMR no EKS com Apache</u> Flink.

```
NAME
STATUS
CHART

NAMESPACE
REVISION
UPDATED
APP VERSION

flink-kubernetes-operator
NAMESPACE 1
2023-02-22 16:43:45.24148
-0500 EST deployed flink-kubernetes-operator-emr-7.9.0 x.y.z-amzn-n
```

Execução de uma aplicação do Flink

Com as versões 6.13.0 e superiores do Amazon EMR, você pode executar uma aplicação do Flink usando o operador do Kubernetes para Flink no modo Aplicação do Amazon EMR no EKS. Com as versões 6.15.0 e superiores do Amazon EMR, você também pode executar uma aplicação do Flink no modo Sessão. Esta página descreve ambos os métodos para executar uma aplicação do Flink com o Amazon EMR no EKS.



Note

Você deve ter um bucket do Amazon S3 criado para armazenar os metadados de alta disponibilidade ao enviar o trabalho do Flink. Se não desejar esse atributo, você poderá desativá-lo. Por padrão, ele é habilitado.

Pré-requisito: antes de executar uma aplicação do Flink com o operador do Kubernetes para Flink, conclua as etapas em the section called "Configuração" e the section called "Como instalar o operador do Kubernetes".

Application mode

Com as versões 6.13.0 e superiores do Amazon EMR, você pode executar uma aplicação do Flink usando o operador do Kubernetes para Flink no modo Aplicação do Amazon EMR no EKS.

1. Crie um arquivo basic-example-app-cluster. yaml de definição FlinkDeployment como no exemplo a seguir. Se você ativou e usou um dos opt-in Regiões da AWS, certifiquese de descomentar e configurar a configuração. fs.s3a.endpoint.region

```
apiVersion: flink.apache.org/v1beta1
kind: FlinkDeployment
metadata:
  name: basic-example-app-cluster
spec:
  flinkConfiguration:
    taskmanager.numberOfTaskSlots: "2"
    #fs.s3a.endpoint.region: OPT_IN_AWS_REGION_NAME
    state.checkpoints.dir: CHECKPOINT_S3_STORAGE_PATH
    state.savepoints.dir: SAVEPOINT_S3_STORAGE_PATH
  flinkVersion: v1_17
  executionRoleArn: JOB_EXECUTION_ROLE_ARN
```

```
emrReleaseLabel: "emr-6.13.0-flink-latest" # 6.13 or higher
jobManager:
  storageDir: HIGH_AVAILABILITY_STORAGE_PATH
  resource:
    memory: "2048m"
    cpu: 1
taskManager:
  resource:
    memory: "2048m"
    cpu: 1
job:
  # if you have your job jar in S3 bucket you can use that path as well
  jarURI: local:///opt/flink/examples/streaming/StateMachineExample.jar
  parallelism: 2
  upgradeMode: savepoint
  savepointTriggerNonce: 0
monitoringConfiguration:
  cloudWatchMonitoringConfiguration:
     logGroupName: LOG_GROUP_NAME
```

2. Envie a implantação do Flink com o comando apresentado a seguir. Isso também criará um objeto FlinkDeployment chamado basic-example-app-cluster.

```
kubectl create -f basic-example-app-cluster.yaml -n <NAMESPACE>
```

3. Acesse a interface do usuário do Flink.

```
kubectl port-forward deployments/basic-example-app-cluster 8081 -n NAMESPACE
```

- 4. Abra localhost: 8081 para visualizar os trabalhos do Flink localmente.
- Limpe o trabalho. Lembre-se de limpar os artefatos do S3 que foram criados para esse trabalho, como pontos de verificação, alta disponibilidade, metadados de pontos de salvamento e registros. CloudWatch

Para obter mais informações sobre o envio de inscrições para o Flink por meio do operador Flink Kubernetes, consulte exemplos de operadores do Flink Kubernetes na pasta em. apache/flink-kubernetes-operator GitHub

Session mode

Com as versões 6.15.0 e superiores do Amazon EMR, você pode executar uma aplicação do Flink usando o operador do Kubernetes para Flink no modo Sessão no Amazon EMR no EKS.

 Crie um arquivo de definição FlinkDeployment chamado basic-exampleapp-cluster.yaml como no exemplo a seguir. Se você ativou e usou um dos opt-in Regiões da AWS, certifique-se de descomentar e configurar a configuração. fs.s3a.endpoint.region

```
apiVersion: flink.apache.org/v1beta1
kind: FlinkDeployment
metadata:
  name: basic-example-session-cluster
spec:
  flinkConfiguration:
    taskmanager.numberOfTaskSlots: "2"
    #fs.s3a.endpoint.region: OPT_IN_AWS_REGION_NAME
    state.checkpoints.dir: CHECKPOINT_S3_STORAGE_PATH
    state.savepoints.dir: SAVEPOINT_S3_STORAGE_PATH
  flinkVersion: v1_17
  executionRoleArn: JOB_EXECUTION_ROLE_ARN
  emrReleaseLabel: "emr-6.15.0-flink-latest"
  jobManager:
    storageDir: HIGH_AVAILABILITY_S3_STORAGE_PATH
    resource:
      memory: "2048m"
      cpu: 1
  taskManager:
    resource:
      memory: "2048m"
      cpu: 1
  monitoringConfiguration:
    s3MonitoringConfiguration:
       logUri:
    cloudWatchMonitoringConfiguration:
       logGroupName: LOG_GROUP_NAME
```

2. Envie a implantação do Flink com o comando apresentado a seguir. Isso também criará um objeto FlinkDeployment chamado basic-example-session-cluster.

```
kubectl create -f basic-example-app-cluster.yaml -n NAMESPACE
```

3. Use o seguinte comando para confirmar se o cluster da sessão LIFECYCLE é STABLE:

```
kubectl get flinkdeployments.flink.apache.org basic-example-session-cluster -
n NAMESPACE
```

A saída deve ser semelhante ao seguinte exemplo:

```
NAME

basic-example-session-cluster

JOB STATUS LIFECYCLE STATE

STABLE
```

4. Crie um arquivo basic-session-job.yaml de recursos de definição personalizado do FlinkSessionJob com o seguinte conteúdo de exemplo:

```
apiVersion: flink.apache.org/v1beta1
kind: FlinkSessionJob
metadata:
  name: basic-session-job
spec:
  deploymentName: basic-session-deployment
  job:
    # If you have your job jar in an S3 bucket you can use that path.
    # To use jar in S3 bucket, set
    # OPERATOR_EXECUTION_ROLE_ARN (--set emrContainers.operatorExecutionRoleArn=
$OPERATOR_EXECUTION_ROLE_ARN)
    # when you install Spark operator
    jarURI: https://repo1.maven.org/maven2/org/apache/flink/flink-examples-
streaming_2.12/1.16.1/flink-examples-streaming_2.12-1.16.1-TopSpeedWindowing.jar
    parallelism: 2
    upgradeMode: stateless
```

5. Envie o trabalho de sessão do Flink com o comando apresentado a seguir. Isso criará um objeto do FlinkSessionJob basic-session-job.

```
kubectl apply -f basic-session-job.yaml -n $NAMESPACE
```

6. Use o seguinte comando para confirmar se o cluster da sessão LIFECYCLE é STABLE e o JOB STATUS é RUNNING:

```
\label{lem:kubectl} \mbox{ kubectl get flinkdeployments.flink.apache.org basic-example-session-cluster - } \mbox{ } \m
```

A saída deve ser semelhante ao seguinte exemplo:

```
NAME JOB STATUS LIFECYCLE STATE basic-example-session-cluster RUNNING STABLE
```

7. Acesse a interface do usuário do Flink.

```
kubectl port-forward deployments/basic-example-session-cluster 8081 -n NAMESPACE
```

- 8. Abra localhost: 8081 para visualizar os trabalhos do Flink localmente.
- 9. Limpe o trabalho. Lembre-se de limpar os artefatos do S3 que foram criados para esse trabalho, como pontos de verificação, alta disponibilidade, metadados de pontos de salvamento e registros. CloudWatch

Permissões do perfil de segurança para executar uma aplicação do Flink

Este tópico descreve os perfis de segurança para implantar e executar uma aplicação do Flink. Há dois perfis necessários para gerenciar uma implantação e criar e gerenciar trabalhos: o perfil de operador e o perfil de trabalho. Este tópico os apresenta e lista as permissões.

Controle de acesso baseado em perfis

Para implantar o operador e executar os trabalhos do Flink, é necessário criar dois perfis do Kubernetes: um perfil de operador e um perfil de trabalho. Por padrão, o Amazon EMR cria os dois perfis ao instalar o operador.

Perfil de operador

Usamos a função de operador para gerenciar flinkdeployments a criação e o JobManager gerenciamento de cada tarefa do Flink e de outros recursos, como serviços.

O nome padrão do perfil de operador é emr-containers-sa-flink-operator e ele requer as permissões apresentadas a seguir.

rules:

- apiGroups:
 - _ ""

resources:

- pods
- services

- events - configmaps - secrets - serviceaccounts verbs: _ '*' - apiGroups: - rbac.authorization.k8s.io resources: - roles - rolebindings verbs: _ '*' - apiGroups: - apps resources: - deployments - deployments/finalizers - replicasets verbs: _ '*' - apiGroups: - extensions resources: - deployments - ingresses verbs: _ '*' - apiGroups: - flink.apache.org resources: - flinkdeployments - flinkdeployments/status - flinksessionjobs - flinksessionjobs/status verbs: _ '*' - apiGroups: - networking.k8s.io resources: - ingresses verbs: _ '*' - apiGroups:

```
- coordination.k8s.io
resources:
- leases
verbs:
- '*'
```

Perfil de trabalho

JobManager Ele usa a função de trabalho para criar e gerenciar TaskManagers e ConfigMaps para cada trabalho.

```
rules:
    apiGroups:
    ""
    resources:
    pods
    configmaps
    verbs:
        '*'
    apiGroups:
        apps
    resources:
        deployments
        deployments/finalizers
    verbs:
        '*'
```

Desinstalação do operador do Kubernetes para Flink para o Amazon EMR no EKS

Siga estas etapas para desinstalar o operador do Kubernetes para Flink.

1. Exclua o operador.

```
helm uninstall flink-kubernetes-operator -n <NAMESPACE>
```

2. Exclua os recursos do Kubernetes que o Helm não desinstala.

```
kubectl delete serviceaccounts, roles, rolebindings -l emr-
containers.amazonaws.com/component=flink.operator --namespace <namespace>
```

Desinstalação do operador 41

kubectl delete crd flinkdeployments.flink.apache.org
flinksessionjobs.flink.apache.org

3. (Opcional) Exclua o cert-manager.

```
kubectl delete -f https://github.com/jetstack/cert-manager/releases/download/
v1.12.0/cert-manager.yaml
```

Uso do Kubernetes nativo do Flink

As versões 6.13.0 e superiores do Amazon EMR oferecem suporte ao Kubernetes nativo para Flink como uma ferramenta de linha de comando que você pode usar para enviar e executar aplicações do Flink para um cluster do Amazon EMR no EKS.

Tópicos

- Configuração do Kubernetes nativo para Flink para o Amazon EMR no EKS
- Conceitos básicos do Kubernetes nativo para Flink para o Amazon EMR no EKS
- Requisitos de segurança da conta JobManager de serviço Flink para Kubernetes nativo

Configuração do Kubernetes nativo para Flink para o Amazon EMR no EKS

Conclua as tarefas apresentadas a seguir para se preparar antes de executar uma aplicação com a CLI do Flink no Amazon EMR no EKS. Se você já se inscreveu na Amazon Web Services (AWS) e usou o Amazon EKS, está com quase tudo pronto para usar o Amazon EMR no EKS. Se você já completou algum dos pré-requisitos, pode ignorá-los e passar para os próximos.

- Instale ou atualize para a versão mais recente do AWS CLI Se você já instalou o AWS CLI, confirme se você tem a versão mais recente.
- Comece a usar o Amazon EKS (eksctl): siga as etapas para criar um cluster do Kubernetes com nós no Amazon EKS.
- <u>Selecione um URI de imagem base do Amazon EMR</u> (versão 6.13.0 ou superiores): o comando do Kubernetes para Flink é compatível com as versões 6.13.0 e superiores do Amazon EMR.
- Confirme se a conta JobManager de serviço tem as permissões apropriadas para criar e assistir TaskManager pods. Para obter mais informações, consulte <u>Requisitos de segurança da conta</u> JobManager de serviço do Flink para o Kubernetes nativo.
- Configure seu perfil de credenciais locais da AWS.

Kubernetes nativo do Flink 42

 <u>Crie ou atualize um arquivo kubeconfig para um cluster do Amazon EKS</u> no qual você deseja executar as aplicações do Flink.

Conceitos básicos do Kubernetes nativo para Flink para o Amazon EMR no EKS

Essas etapas mostram como configurar, definir uma conta de serviço e executar uma aplicação do Flink. O Kubernetes nativo para Flink é usado na implantação do Flink em um cluster do Kubernetes em execução.

Como configurar e executar uma aplicação do Flink

As versões 6.13.0 e posteriores do Amazon EMR oferecem suporte ao Kubernetes nativo para Flink para a execução de aplicações do Flink em um cluster do Amazon EKS. Para executar uma aplicação do Flink, siga estas etapas:

- 1. Antes de executar uma aplicação do Flink com o comando do Kubernetes nativo para Flink, conclua as etapas em the section called "Configuração".
- 2. Download e instalação do Flink.
- 3. Defina os valores para as variáveis de ambiente a seguir.

```
#Export the FLINK_HOME environment variable to your local installation of Flink export FLINK_HOME=/usr/local/bin/flink #Will vary depending on your installation export NAMESPACE=flink export CLUSTER_ID=flink-application-cluster export IMAGE=<123456789012.dkr.ecr.sample-Região da AWS-.amazonaws.com/flink/emr-6.13.0-flink:latest> export FLINK_SERVICE_ACCOUNT=emr-containers-sa-flink export FLINK_CLUSTER_ROLE_BINDING=emr-containers-crb-flink
```

4. Crie uma conta de serviço para gerenciar os recursos do Kubernetes.

```
kubectl create serviceaccount $FLINK_SERVICE_ACCOUNT -n $NAMESPACE
kubectl create clusterrolebinding $FLINK_CLUSTER_ROLE_BINDING --clusterrole=edit --
serviceaccount=$NAMESPACE:$FLINK_SERVICE_ACCOUNT
```

5. Execute o comando run-application da CLI.

```
$FLINK_HOME/bin/flink run-application \
```

Começar 43

```
--target kubernetes-application \
    -Dkubernetes.namespace=$NAMESPACE \
    -Dkubernetes.cluster-id=$CLUSTER_ID \
    -Dkubernetes.container.image.ref=$IMAGE \
    -Dkubernetes.service-account=$FLINK_SERVICE_ACCOUNT \
   local:///opt/flink/examples/streaming/Iteration.jar
2022-12-29 21:13:06,947 INFO org.apache.flink.kubernetes.utils.KubernetesUtils
           [] - Kubernetes deployment requires a fixed port. Configuration
blob.server.port will be set to 6124
2022-12-29 21:13:06,948 INFO org.apache.flink.kubernetes.utils.KubernetesUtils
            [] - Kubernetes deployment requires a fixed port. Configuration
taskmanager.rpc.port will be set to 6122
2022-12-29 21:13:07,861 WARN
org.apache.flink.kubernetes.KubernetesClusterDescriptor [] - Please note that
Flink client operations(e.g. cancel, list, stop, savepoint, etc.) won't work from
outside the Kubernetes cluster since 'kubernetes.rest-service.exposed.type' has
been set to ClusterIP.
2022-12-29 21:13:07,868 INFO
org.apache.flink.kubernetes.KubernetesClusterDescriptor
                                                              [] - Create flink
 application cluster flink-application-cluster successfully, JobManager Web
Interface: http://flink-application-cluster-rest.flink:8081
```

Examine os recursos do Kubernetes criados.

```
kubectl get all -n <namespace>
NAME READY STATUS RESTARTS AGE
pod/flink-application-cluster-546687cb47-w2p2z 1/1 Running 0 3m37s
pod/flink-application-cluster-taskmanager-1-1 1/1 Running 0 3m24s

NAME TYPE CLUSTER-IP EXTERNAL-IP PORT(S) AGE
service/flink-application-cluster ClusterIP None <none> 6123/TCP,6124/TCP 3m38s
service/flink-application-cluster-rest ClusterIP 10.100.132.158 <none> 8081/TCP
3m38s

NAME READY UP-TO-DATE AVAILABLE AGE
deployment.apps/flink-application-cluster 1/1 1 1 3m38s

NAME DESIRED CURRENT READY AGE
replicaset.apps/flink-application-cluster-546687cb47 1 1 1 3m38s
```

7. Encaminhe a porta para 8081.

Começar 44

kubectl port-forward service/flink-application-cluster-rest 8081 -n <namespace>
Forwarding from 127.0.0.1:8081 -> 8081

- 8. Acesse localmente a interface do usuário do Flink.
- 9. Exclua a aplicação do Flink.

```
kubectl delete deployment.apps/flink-application-cluster -n <namespace>
deployment.apps "flink-application-cluster" deleted
```

Para obter mais informações sobre o envio de aplicações para o Flink, consulte <u>Native Kubernetes</u> na documentação do Apache Flink.

Requisitos de segurança da conta JobManager de serviço Flink para Kubernetes nativo

O JobManager pod Flink usa uma conta de serviço do Kubernetes para acessar o servidor da API Kubernetes para criar e assistir pods. TaskManager A conta JobManager de serviço deve ter as permissões apropriadas para create/delete TaskManager os pods e permitir que o TaskManager líder do ConfigMaps to watch recupere o endereço de JobManager e ResourceManager no seu cluster.

As regras apresentadas a seguir se aplicam a esta conta de serviço.

```
rules:
- apiGroups:
  _ ""
  resources:
  - pods
  verbs:
  _ "*"
- apiGroups:
  _ ""
  resources:
  - services
  verbs:
  _ "*"
- apiGroups:
  _ ""
  resources:
```

Requisitos de segurança 45

```
- configmaps
verbs:
- "*"
- apiGroups:
- "apps"
resources:
- deployments
verbs:
- "*"
```

Personalização de imagens do Docker para o Flink e o FluentD

Siga as etapas a seguir para personalizar as imagens do Docker para o Amazon EMR no EKS com imagens do Apache Flink ou do FluentD. Isso inclui orientação técnica para obter uma imagem base, personalizá-la, publicá-la e enviar uma workload.

Tópicos

- Pré-requisitos
- Etapa 1: recuperar uma imagem base do Amazon Elastic Container Registry
- Etapa 2: personalizar uma imagem base
- Etapa 3: publicar uma imagem personalizada
- Etapa 4: enviar uma workload do Flink no Amazon EMR usando uma imagem personalizada

Pré-requisitos

Antes de personalizar a imagem do Docker, certifique-se de ter atendido aos seguintes prérequisitos:

- · Concluir as etapas em Setting up the Flink Kubernetes operator for Amazon EMR on EKS.
- Instalar o Docker em seu ambiente. Para obter mais informações, consulte Get Docker.

Etapa 1: recuperar uma imagem base do Amazon Elastic Container Registry

A imagem base contém o runtime do Amazon EMR e os conectores necessários para acessar outros Serviços da AWS. Se estiver usando o Amazon EMR no EKS com a versão 6.14.0 ou posterior do

Flink, você pode obter as imagens base na galeria pública do Amazon ECR. Navegue pela galeria para encontrar o link da imagem e extraia-a para seu Workspace local. Por exemplo, para a versão 6.14.0 do Amazon EMR, o comando docker pull a seguir retorna a imagem base padrão mais recente. Substitua emr-6.14.0:latest pela versão de lançamento desejada.

```
docker pull public.ecr.aws/emr-on-eks/flink/emr-6.14.0-flink:latest
```

Estes são os links para a imagem da galeria do Flink e a imagem da galeria do FluentD:

- emr-on-eks/flink/emr-6.14.0-flink
- emr-on-eks/fluente/emr-6.14.0 (

Etapa 2: personalizar uma imagem base

As etapas a seguir descrevem como personalizar a imagem base extraída do Amazon ECR.

- 1. Crie um novo Dockerfile em seu Workspace local.
- 2. Edite o Dockerfile e adicione o conteúdo a seguir. Este Dockerfile usa a imagem de contêiner que você extraiu de public.ecr.aws/emr-on-eks/flink/emr-7.10.0-flink:latest.

```
FROM public.ecr.aws/emr-on-eks/flink/emr-7.10.0-flink:latest
USER root
### Add customization commands here ####
USER hadoop:hadoop
```

Use a configuração a seguir se estiver utilizando o Fluentd.

```
FROM public.ecr.aws/emr-on-eks/fluentd/emr-7.10.0:latest
USER root
### Add customization commands here ####
USER hadoop:hadoop
```

 Adicione comandos no Dockerfile para personalizar a imagem base. O comando a seguir demonstra como instalar bibliotecas Python.

```
FROM public.ecr.aws/emr-on-eks/flink/emr-7.10.0-flink:latest
USER root
```

```
RUN pip3 install --upgrade boto3 pandas numpy // For python 3
USER hadoop:hadoop
```

4. No mesmo diretório em que o DockerFile foi criado, execute o comando a seguir para criar a imagem do Docker. O campo que você fornece após o sinalizador -t é o nome personalizado da imagem.

```
docker build -t <YOUR_ACCOUNT_ID>.dkr.ecr.<YOUR_ECR_REGION>.amazonaws.com/
<ECR_REPO>:<ECR_TAG>
```

Etapa 3: publicar uma imagem personalizada

Agora você pode publicar a nova imagem do Docker no Amazon ECR Registry.

 Execute o comando a seguir para criar um repositório do Amazon ECR e armazenar a imagem do Docker. Forneça um nome para o repositório, como emr_custom_repo.. Para obter mais informações, consulte <u>Create a repository</u> no Guia do usuário do Amazon Elastic Container Registry.

```
aws ecr create-repository \
     --repository-name emr_custom_repo \
     --image-scanning-configuration scanOnPush=true \
     --region <AWS_REGION>
```

Execute o comando apresentado a seguir para realizar a autenticação em seu registro padrão.
 Para obter mais informações, consulte <u>Authenticate to your default registry</u> no Guia do usuário do Amazon Elastic Container Registry.

```
aws ecr get-login-password --region <a href="https://docker.login.news.com">AWS_REGION> | docker login --username AWS --password-stdin <a href="https://docker.login.news.com">AWS_ACCOUNT_ID> .dkr.ecr.</a><a href="https://docker.login.news.com">YOUR_ECR_REGION> .amazonaws.com</a><a href="https://docker.login.news.com">news.com</a><a href="https://docker.login.news.com">news.com</a
```

 Envie a imagem. Para obter mais informações, consulte <u>Push an image to Amazon ECR</u> no Guia do usuário do Amazon Elastic Container Registry.

```
docker push <YOUR_ACCOUNT_ID>.dkr.ecr.<YOUR_ECR_REGION>.amazonaws.com/
<ECR_REPO>:<ECR_TAG>
```

Etapa 4: enviar uma workload do Flink no Amazon EMR usando uma imagem personalizada

Faça as alterações a seguir na especificação FlinkDeployment para usar uma imagem personalizada. Para fazer isso, insira sua própria imagem na linha spec.image da especificação de implantação.

```
apiVersion: flink.apache.org/v1beta1
  kind: FlinkDeployment
  metadata:
    name: basic-example
  spec:
    flinkVersion: v1_18
    image: <YOUR_ACCOUNT_ID>.dkr.ecr.<YOUR_ECR_REGION>.amazonaws.com/
<ECR_REPO>:<ECR_TAG>
    imagePullPolicy: Always
  flinkConfiguration:
    taskmanager.numberOfTaskSlots: "1"
```

Para usar uma imagem personalizada no trabalho do FluentD, insira sua própria imagem na linha monitoringConfiguration.image da especificação de implantação.

```
monitoringConfiguration:
    image: <YOUR_ACCOUNT_ID>.dkr.ecr.<YOUR_ECR_REGION>.amazonaws.com/
<ECR_REPO>:<ECR_TAG>
    cloudWatchMonitoringConfiguration:
    logGroupName: flink-log-group
    logStreamNamePrefix: custom-fluentd
```

Monitoramento do operador do Kubernetes para Flink e dos trabalhos do Flink

Esta seção descreve diversas maneiras para monitorar trabalhos do Flink com o Amazon EMR no EKS. Isso inclui a integração do Flink com o Amazon Managed Service for Prometheus, o uso do Flink Web Dashboard, que fornece o status e as métricas do trabalho, ou o uso de uma configuração de monitoramento para enviar dados de log ao Amazon S3 e ao Amazon CloudWatch.

Tópicos

Uso do Amazon Managed Service para Prometheus no monitoramento de trabalhos do Flink

Como enviar uma workload 49

- Uso da interface do usuário do Flink para monitorar trabalhos do Flink
- <u>Uso da configuração de monitoramento para monitorar o operador do Kubernetes do Flink e os</u> trabalhos do Flink

Uso do Amazon Managed Service para Prometheus no monitoramento de trabalhos do Flink

Você pode integrar o Apache Flink ao Amazon Managed Service for Prometheus (portal de gerenciamento). O Amazon Managed Service for Prometheus oferece suporte à ingestão de métricas de servidores do Amazon Managed Service for Prometheus em clusters em execução no Amazon EKS. O Amazon Managed Service for Prometheus funciona em conjunto com um servidor do Prometheus que já está em execução no cluster do Amazon EKS. A execução da integração do Amazon Managed Service for Prometheus com o operador do Flink do Amazon EMR implantará e configurará automaticamente um servidor do Prometheus para a integração com o Amazon Managed Service for Prometheus.

- 1. <u>Crie um Workspace do Amazon Managed Service for Prometheus</u>. Este Workspace serve como um endpoint de ingestão. Você precisará do URL de gravação remota posteriormente.
- 2. Configure perfis do IAM para as contas de serviço.

Para esse método de integração, use perfis do IAM para as contas de serviço no cluster do Amazon EKS em que o servidor do Prometheus está em execução. Esses perfis também são chamados de perfis de serviço.

Se você ainda não tiver os perfis, configure perfis de serviço para a ingestão de métricas de clusters do Amazon EKS.

Antes de continuar, crie um perfil do IAM chamado amp-iamproxy-ingest-role.

3. Instale o operador do Flink do Amazon EMR com o Amazon Managed Service for Prometheus.

Agora que você tem um Workspace do Amazon Managed Service for Prometheus, um perfil do IAM dedicado para o Amazon Managed Service for Prometheus e as permissões necessárias, é possível instalar o operador do Flink do Amazon EMR.

Criar um arquivo enable-amp. yaml Esse arquivo permite usar uma configuração personalizada para substituir as configurações do Amazon Managed Service for Prometheus. Use seus próprios perfis.

```
kube-prometheus-stack:
    prometheus:
    serviceAccount:
        create: true
        name: "amp-iamproxy-ingest-service-account"
        annotations:
            eks.amazonaws.com/role-arn: "arn:aws:iam::<AWS_ACCOUNT_ID>:role/amp-
iamproxy-ingest-role"
    remoteWrite:
        - url: <AMAZON_MANAGED_PROMETHEUS_REMOTE_WRITE_URL>
        sigv4:
            region: <AWS_REGION>
        queueConfig:
            maxSamplesPerSend: 1000
            maxShards: 200
            capacity: 2500
```

Use o comando <u>Helm Install --set</u> para transferir as substituições para o chart flink-kubernetes-operator.

```
helm upgrade -n <namespace> flink-kubernetes-operator \
  oci://public.ecr.aws/emr-on-eks/flink-kubernetes-operator \
  --set prometheus.enabled=true
  -f enable-amp.yaml
```

Esse comando instala automaticamente um relator do Prometheus no operador na porta 9999. Qualquer FlinkDeployment futura também expõe uma porta para metrics em 9249.

- As métricas do operador do Flink aparecem no Prometheus sob o rótulo flink_k8soperator_.
- As métricas Task Manager do Flink aparecem no Prometheus sob o rótulo flink_taskmanager_.
- As métricas Job Manager do Flink aparecem no Prometheus sob o rótulo flink_jobmanager_.

Uso da interface do usuário do Flink para monitorar trabalhos do Flink

Para monitorar a integridade e a performance de uma aplicação do Flink em execução, use o painel Web do Flink. Esse painel fornece informações sobre o status do trabalho, o número TaskManagers e as métricas e registros do trabalho. Ele também permite visualizar e modificar a configuração do trabalho do Flink e interagir com o cluster do Flink para enviar ou cancelar trabalhos.

Para acessar o painel Web do Flink para uma aplicação do Flink em execução no Kubernetes:

 Use o kubectl port-forward comando para encaminhar uma porta local para a porta na qual o Flink Web Dashboard está sendo executado nos pods do TaskManager aplicativo Flink. Por padrão, esta porta é 8081. deployment-nameSubstitua pelo nome da implantação do aplicativo Flink acima.

```
kubectl get deployments -n namespace
```

Resultado do exemplo:

```
kubectl port-forward deployments/deployment-name 8081 -n namespace
```

2. Se você quiser usar uma porta diferente localmente, use o parâmetro:8081*local-port*.

```
kubectl port-forward -n flink deployments/basic-example 8080:8081
```

3. Em um navegador da Web, vá até http://localhost:8081 (ou http://localhost:local-port, se você usou uma porta local personalizada) para acessar o painel Web do Flink. Esse painel mostra informações sobre o aplicativo Flink em execução, como o status do trabalho, o número e as métricas e registros do trabalho. TaskManagers

Uso da configuração de monitoramento para monitorar o operador do Kubernetes do Flink e os trabalhos do Flink

A configuração de monitoramento permite que você configure facilmente o arquivamento de registros do seu aplicativo Flink e dos registros do operador no S3 and/or CloudWatch (você pode escolher um ou ambos). Isso adiciona um sidecar FluentD aos JobManager seus pods TaskManager e, posteriormente, encaminha os registros desses componentes para os coletores configurados.



Note

Você deve configurar perfis do IAM para a conta de serviço do seu operador do Flink e seu trabalho do Flink (contas de serviço) para poder usar esse recurso, pois ele requer interação com outros Serviços da AWS. Você deve configurar isso usando IRSA em Configuração do operador do Kubernetes para Flink para o Amazon EMR no EKS.

Logs da aplicação do Flink

Você pode definir essa configuração da maneira apresentada a seguir.

```
apiVersion: flink.apache.org/v1beta1
kind: FlinkDeployment
metadata:
  name: basic-example
spec:
  image: FLINK IMAGE TAG
  imagePullPolicy: Always
  flinkVersion: v1_17
  flinkConfiguration:
    taskmanager.numberOfTaskSlots: "2"
  executionRoleArn: JOB EXECUTION ROLE
  jobManager:
    resource:
      memory: "2048m"
      cpu: 1
  taskManager:
    resource:
      memory: "2048m"
      cpu: 1
  job:
    jarURI: local:///opt/flink/examples/streaming/StateMachineExample.jar
  monitoringConfiguration:
    s3MonitoringConfiguration:
      logUri: S3 BUCKET
    cloudWatchMonitoringConfiguration:
      logGroupName: LOG GROUP NAME
      logStreamNamePrefix: LOG GROUP STREAM PREFIX
    sideCarResources:
      limits:
        cpuLimit: 500m
```

```
memoryLimit: 250Mi
containerLogRotationConfiguration:
   rotationSize: 2GB
   maxFilesToKeep: 10
```

A seguir, são apresentadas as opções de configuração.

- s3MonitoringConfiguration: a chave de configuração para configurar o encaminhamento para o S3.
 - logUri (obrigatório): o caminho do bucket do S3 em que você deseja armazenar seus logs.
 - O caminho no S3 depois que os logs forem carregados será semelhante ao apresentado a seguir.
 - Nenhuma alternância de log habilitada:

```
s3://${logUri}/${POD NAME}/STDOUT or STDERR.gz
```

 A alternância de log está habilitada. Você pode usar um arquivo rotacionado e um arquivo atual (um que não tenha carimbo de data).

```
s3://${logUri}/${POD NAME}/STDOUT or STDERR.gz
```

O formato a seguir é um número incremental.

```
s3://${logUri}/${POD NAME}/stdout_YYYYMMDD_index.gz
```

As permissões do IAM apresentadas a seguir são obrigatórias para usar este encaminhador.

 cloudWatchMonitoringConfiguration— chave de configuração para a qual configurar o encaminhamento. CloudWatch

- logGroupName(obrigatório) nome do grupo de CloudWatch registros para o qual você deseja enviar registros (cria automaticamente o grupo se ele não existir).
- logStreamNamePrefix (opcional): o nome do fluxo de logs para o qual você deseja enviar os logs. O valor padrão é uma string vazia. O formato é o seguinte:

```
${logStreamNamePrefix}/${POD NAME}/STDOUT or STDERR
```

As permissões do IAM apresentadas a seguir são obrigatórias para usar este encaminhador.

```
{
    "Effect": "Allow",
    "Action": [
        "logs:CreateLogStream",
        "logs:CreateLogGroup",
        "logs:PutLogEvents"
],
    "Resource": [
        "arn:aws:logs:REGION:ACCOUNT-ID:log-group:{YOUR_LOG_GROUP_NAME}:*",
        "arn:aws:logs:REGION:ACCOUNT-ID:log-group:{YOUR_LOG_GROUP_NAME}"
]
}
```

- sideCarResources (opcional): a chave de configuração para definir limites de recursos no contêiner sidecar do Fluent Bit iniciado.
 - memoryLimit (opcional): o valor padrão é 512 Mi. Ajuste de acordo com suas necessidades.
 - cpuLimit (opcional): esta opção não tem um padrão. Ajuste de acordo com suas necessidades.
- containerLogRotationConfiguration (opcional): controla o comportamento de alternância de log do contêiner. Ele é habilitado por padrão.
 - rotationSize (obrigatório): especifica o tamanho do arquivo para a alternância de log. O intervalo de valores possíveis é de 2 KB a 2 GB. A parcela numérica da unidade do parâmetro rotationSize é transferida como um número inteiro. Como não há suporte para valores decimais, você pode especificar um tamanho de rotação de 1,5 GB, por exemplo, com o valor 1.500 MB. O padrão é 2 GB.
 - maxFilesToKeep (obrigatório): especifica o número máximo de arquivos a serem retidos no contêiner após a alternância ter ocorrido. O valor mínimo é de 1 e o valor máximo é de 50. O padrão é 10.

Logs do operador do Flink

Também podemos habilitar o arquivamento de logs para o operador ao usar as opções apresentadas a seguir no arquivo values.yaml da instalação do chart do Helm. Você pode ativar o S3 ou ambos. CloudWatch

```
monitoringConfiguration:
    s3MonitoringConfiguration:
    logUri: "$3-BUCKET"
    totalFileSize: "1G"
    uploadTimeout: "1m"
    cloudWatchMonitoringConfiguration:
    logGroupName: "flink-log-group"
    logStreamNamePrefix: "example-job-prefix-test-2"
    sideCarResources:
    limits:
        cpuLimit: 1
        memoryLimit: 800Mi
    memoryBufferLimit: 700M
```

A seguir, estão apresentadas as opções de configuração disponíveis em monitoringConfiguration.

- s3MonitoringConfiguration: defina esta opção para realizar o arquivamento no S3.
- logUri (obrigatório): o caminho do bucket do S3 em que você deseja armazenar seus logs.
- A seguir, estão apresentados os formatos de como os caminhos do bucket do S3 podem parecer depois que os logs são carregados.
 - Nenhuma alternância de log habilitada.

```
s3://${logUri}/${POD NAME}/OPERATOR or WEBHOOK/STDOUT or STDERR.gz
```

 A alternância de log está habilitada. Você pode usar um arquivo rotacionado e um arquivo atual (um que não tenha carimbo de data).

```
s3://${logUri}/${POD NAME}/OPERATOR or WEBHOOK/STDOUT or STDERR.gz
```

O índice de formato a seguir é um número incremental.

```
s3://${logUri}/${POD NAME}/OPERATOR or WEBHOOK/stdout_YYYYMMDD_index.gz
```

- cloudWatchMonitoringConfiguration— a chave de configuração para a qual configurar o encaminhamento. CloudWatch
 - logGroupName(obrigatório) nome do grupo de CloudWatch registros para o qual você deseja enviar registros. O grupo é criado automaticamente, se não existir.
 - logStreamNamePrefix (opcional): o nome do fluxo de logs para o qual você deseja enviar os logs. O valor padrão é uma string vazia. O formato em CloudWatch é o seguinte:

```
${logStreamNamePrefix}/${POD NAME}/STDOUT or STDERR
```

- sideCarResources (opcional): a chave de configuração para definir limites de recursos no contêiner sidecar do Fluent Bit iniciado.
 - memoryLimit (opcional): o limite de memória. Ajuste de acordo com suas necessidades. O padrão é de 512Mi.
 - cpuLimit: o limite de CPU. Ajuste de acordo com suas necessidades. Nenhum valor padrão.
- containerLogRotationConfiguration (opcional): controla o comportamento de alternância de log do contêiner. Ele é habilitado por padrão.
 - rotationSize (obrigatório): especifica o tamanho do arquivo para a alternância de log. O intervalo de valores possíveis é de 2 KB a 2 GB. A parcela numérica da unidade do parâmetro rotationSize é transferida como um número inteiro. Como não há suporte para valores decimais, você pode especificar um tamanho de rotação de 1,5 GB, por exemplo, com o valor 1.500 MB. O padrão é 2 GB.
 - maxFilesToKeep (obrigatório): especifica o número máximo de arquivos a serem retidos no contêiner após a alternância ter ocorrido. O valor mínimo é de 1 e o valor máximo é de 50. O padrão é 10.

Como o Flink oferece suporte à alta disponibilidade e resiliência do trabalho

As seções a seguir descrevem como o Flink torna os trabalhos mais confiáveis e altamente disponíveis. Ele faz isso por meio de capacidades integradas, como a alta disponibilidade do Flink e vários recursos de recuperação, caso ocorram falhas.

Tópicos

Uso da alta disponibilidade (HA) para operadores e aplicações do Flink

- Otimização dos tempos de reinicialização de trabalhos do Flink para operações de recuperação e ajuste de escala de tarefas com o Amazon EMR no EKS
- Desativação tranquila de instâncias spot com Flink no Amazon EMR no EKS

Uso da alta disponibilidade (HA) para operadores e aplicações do Flink

Este tópico mostra como configurar a alta disponibilidade e descreve como ela funciona em alguns casos de uso diferentes. Isso inclui o uso do gerenciador de trabalhos e do Kubernetes nativo para Flink.

Alta disponibilidade do operador do Flink

Habilitamos a alta disponibilidade do operador do Flink para que possamos fazer failover para um operador do Flink em espera a fim de minimizar o tempo de inatividade no loop de controle do operador se ocorrerem falhas. A alta disponibilidade é habilitada por padrão e o número padrão de réplicas iniciais para o operador é dois. É possível configurar o campo de réplicas em seu arquivo values. yaml para o chart do Helm.

Os seguintes campos são personalizáveis:

- replicas (opcional, o padrão é dois): definir esse número como maior que um cria outros operadores em espera e permite uma recuperação mais rápida do trabalho.
- highAvailabilityEnabled (opcional, o padrão é "true"): controla se você deseja habilitar a HA. Especificar esse parâmetro como "true" habilita o suporte à implantação multi-AZ e define os parâmetros flink-conf.yaml corretos.

Você pode desativar a HA para seu operador ao definir a configuração apresentada a seguir em seu arquivo values.yaml.

```
imagePullSecrets: []

replicas: 1

# set this to false if you don't want HA
highAvailabilityEnabled: false
...
```

Implantação multi-AZ

Criamos os pods do operador em várias zonas de disponibilidade. Esta é uma restrição leve, e seus pods do operador serão programados na mesma AZ, se você não tiver recursos suficientes em uma AZ diferente.

Determinação da réplica líder

Se o HA estiver habilitado, as réplicas usarão uma concessão para determinar qual delas JMs é a líder e usarão uma locação K8s para a eleição do líder. Você pode descrever a concessão e consultar o campo .Spec.Holder Identity para determinar o líder atual.

```
kubectl describe lease <Helm Install Release Name>-<NAMESPACE>-lease -n <NAMESPACE> |
grep "Holder Identity"
```

Interação entre o Flink e o S3

Configuração de credenciais de acesso

Certifique-se de ter configurado o IRSA com as permissões do IAM apropriadas para acessar o bucket do S3.

Busca por trabalhos em JARs do modo de aplicação do S3

O operador do Flink também oferece suporte à busca de aplicações do S3 em JARs. Você acabou de fornecer a localização S3 para o JARuri em sua FlinkDeployment especificação.

Você também pode usar esse recurso para baixar outros artefatos, como PyFlink scripts. O script do Python resultante é descartado no caminho /opt/flink/usrlib/.

O exemplo a seguir demonstra como usar esse recurso para um PyFlink trabalho. Observe os campos jarURI e args.

```
apiVersion: flink.apache.org/v1beta1
kind: FlinkDeployment
metadata:
   name: python-example
spec:
   image: <YOUR CUSTOM PYFLINK IMAGE>
   emrReleaseLabel: "emr-6.12.0-flink-latest"
   flinkVersion: v1_16
   flinkConfiguration:
     taskmanager.numberOfTaskSlots: "1"
```

```
serviceAccount: flink
  jobManager:
    highAvailabilityEnabled: false
    replicas: 1
    resource:
      memory: "2048m"
      cpu: 1
  taskManager:
    resource:
      memory: "2048m"
      cpu: 1
  job:
    jarURI: "s3://<S3-BUCKET>/scripts/pyflink.py" # Note, this will trigger the
 artifact download process
    entryClass: "org.apache.flink.client.python.PythonDriver"
    args: ["-pyclientexec", "/usr/local/bin/python3", "-py", "/opt/flink/usrlib/
pyflink.py"]
    parallelism: 1
    upgradeMode: stateless
```

Conectores do S3 para Flink

O Flink vem com dois conectores do S3 (listados abaixo). As seções a seguir debatem sobre o momento de usar cada conector.

Ponto de verificação: conector do S3 para Presto

- Defina o esquema do S3 como s3p://.
- O conector recomendado a ser usado como ponto de verificação para o s3. Para obter mais informações, consulte S3-specific na documentação do Apache Flink.

Exemplo de FlinkDeployment especificação:

```
apiVersion: flink.apache.org/v1beta1
kind: FlinkDeployment
metadata:
   name: basic-example
spec:
   flinkConfiguration:
     taskmanager.numberOfTaskSlots: "2"
     state.checkpoints.dir: s3p://<BUCKET-NAME>/flink-checkpoint/
```

Leitura e gravação no S3: conector Hadoop S3

- Defina o esquema do S3 como s3:// ou (s3a://).
- O conector recomendado para a leitura e a gravação de arquivos do S3 (somente um conector do S3 que implementa a interface Flink Filesystem).
- Por padrão, definimos fs.s3a.aws.credentials.provider no arquivo flink-conf.yaml, que é com.amazonaws.auth.WebIdentityTokenCredentialsProvider. Se você substituir completamente o flink-conf padrão e estiver interagindo com o S3, certifique-se de usar este provedor.

Exemplo de FlinkDeployment especificação

```
apiVersion: flink.apache.org/v1beta1
kind: FlinkDeployment
metadata:
   name: basic-example
spec:
   job:
     jarURI: local:///opt/flink/examples/streaming/WordCount.jar
     args: [ "--input", "s3a://<INPUT BUCKET>/PATH", "--output", "s3a://<OUTPUT BUCKET>/
PATH" ]
   parallelism: 2
   upgradeMode: stateless
```

JobManager do Flink

A alta disponibilidade (HA) para implantações do Flink permite que os trabalhos continuem progredindo mesmo que um erro transitório seja encontrado e você falhe. JobManager Com a HA habilitada, os trabalhos serão reiniciados, mas a partir do último ponto de verificação com êxito. Sem o HA ativado, o Kubernetes reiniciará o seu JobManager, mas seu trabalho começará como um novo trabalho e perderá o progresso. Depois de configurar o HA, podemos pedir ao Kubernetes que armazene os metadados de HA em um armazenamento persistente para referência no caso de uma falha transitória no JobManager e, em seguida, retome nossos trabalhos a partir do último ponto de verificação bem-sucedido.

A HA é habilitada, por padrão, para os trabalhos do Flink (a contagem de réplicas é definida como dois, o que exigirá que você forneça um local de armazenamento do S3 para que os metadados de HA persistam).

Configurações de HA

```
apiVersion: flink.apache.org/v1beta1
kind: FlinkDeployment
metadata:
  name: basic-example
spec:
  flinkConfiguration:
    taskmanager.numberOfTaskSlots: "2"
  executionRoleArn: "<JOB EXECUTION ROLE ARN>"
  emrReleaseLabel: "emr-6.13.0-flink-latest"
  jobManager:
    resource:
      memory: "2048m"
      cpu: 1
    replicas: 2
    highAvailabilityEnabled: true
    storageDir: "s3://<S3 PERSISTENT STORAGE DIR>"
  taskManager:
    resource:
      memory: "2048m"
      cpu: 1
```

A seguir estão as descrições das configurações de HA apresentadas acima no Job Manager (definidas em .spec.jobManager):

- highAvailabilityEnabled (opcional, o padrão é "true"): defina como false se você não desejar que a HA seja habilitada e não quiser usar as configurações de HA fornecidas. Você ainda pode manipular o campo "réplicas" para configurar manualmente a HA.
- replicas(opcional, o padrão é 2): definir esse número como maior que 1 cria outro modo de espera JobManagers e permite uma recuperação mais rápida do seu trabalho. Se você desabilitar a HA, deverá definir a contagem de réplicas como um ou continuará recebendo erros de validação (somente uma réplica tem suporte se a HA não estiver habilitada).
- storageDir (obrigatório): por usar a contagem de réplicas como dois, por padrão, é necessário fornecer um storageDir persistente. No momento, este campo aceita somente caminhos do S3 como local de armazenamento.

Localidade de pod

Se você habilitar o HA, também tentaremos colocar pods na mesma AZ, o que leva a um melhor desempenho (latência de rede reduzida ao ter pods no mesmo). AZs Este é um processo de melhor esforço, ou seja, se você não tiver recursos suficientes na AZ em que a maioria dos seus pods está programada, os pods restantes ainda serão programados, mas poderão acabar em um nó externo a esta AZ.

Determinação da réplica líder

Se o HA estiver habilitado, as réplicas usarão uma concessão para determinar qual delas JMs é a líder e usam um K8s Configmap como armazenamento de dados para armazenar esses metadados. Se quiser determinar o líder, você pode consultar o conteúdo do Configmap e a chave org.apache.flink.k8s.leader.restserver nos dados para encontrar o pod do K8s com o endereço IP. Você também pode usar os comandos bash apresentados a seguir.

```
ip=$(kubectl get configmap -n <NAMESPACE> <JOB-NAME>-cluster-config-map -o json | jq -
r ".data[\"org.apache.flink.k8s.leader.restserver\"]" | awk -F: '{print $2}' | awk -F
 '/' '{print $3}')
kubectl get pods -n NAMESPACE -o json | jq -r ".items[] | select(.status.podIP ==
 \"$ip\") | .metadata.name"
```

Trabalho do Flink: Kubernetes nativo

As versões 6.13.0 e superiores do Amazon EMR oferecem suporte ao Kubernetes nativo para Flink para a execução de aplicações do Flink no modo de alta disponibilidade em um cluster do Amazon EKS.



Note

Você deve ter um bucket do Amazon S3 criado para armazenar os metadados de alta disponibilidade ao enviar o trabalho do Flink. Se não desejar esse atributo, você poderá desativá-lo. Por padrão, ele é habilitado.

Para ativar o recurso de alta disponibilidade do Flink, forneça os parâmetros do Flink descritos a seguir ao executar o comando da CLI run-application. Os parâmetros são definidos abaixo do exemplo.

```
-Dhigh-availability.type=kubernetes \
-Dhigh-availability.storageDir=S3://DOC-EXAMPLE-STORAGE-BUCKET \
```

```
-
Dfs.s3a.aws.credentials.provider="com.amazonaws.auth.WebIdentityTokenCredentialsProvider"
\
-Dkubernetes.jobmanager.replicas=3 \
-Dkubernetes.cluster-id=example-cluster
```

 Dhigh-availability.storageDir: o bucket do Amazon S3 onde você deseja armazenar os metadados de alta disponibilidade para o trabalho.

Dkubernetes.jobmanager.replicas: o número de pods do gerenciador de trabalhos a serem criados como um número inteiro maior que 1.

Dkubernetes.cluster-id: um ID exclusivo que identifica o cluster do Flink.

Otimização dos tempos de reinicialização de trabalhos do Flink para operações de recuperação e ajuste de escala de tarefas com o Amazon EMR no EKS

Quando uma tarefa falha ou ocorre uma operação de ajuste de escala, o Flink tenta reexecutar a tarefa com base no último ponto de verificação concluído. O processo de reinicialização pode levar um minuto ou mais para ser executado, dependendo do tamanho do estado do ponto de verificação e do número de tarefas paralelas. Durante o período de reinicialização, as tarefas de backlog podem se acumular para o trabalho. No entanto, existem algumas maneiras de o Flink otimizar a velocidade de recuperação e reinicialização dos gráficos de execução para melhorar a estabilidade do trabalho.

Esta página descreve algumas maneiras que o Flink no Amazon EMR pode usar para melhorar o tempo de reinício de trabalhos durante as operações de recuperação ou ajuste de escala de tarefas em instâncias spot. As instâncias spot são capacidades computacionais não utilizadas que estão disponíveis com desconto. Elas têm comportamentos únicos, incluindo interrupções ocasionais, por isso é importante entender como o Amazon EMR no EKS lida com elas e como o Amazon EMR no EKS realiza a desativação e a reinicialização do trabalho.

Tópicos

- Recuperação local de tarefas
- Recuperação local de tarefas por meio da montagem de volume do Amazon EBS
- Ponto de verificação incremental genérico baseado em logs
- Recuperação refinada

Mecanismo de reinício combinado no agendador adaptável

Recuperação local de tarefas



Note

A recuperação local de tarefas é compatível com o Flink nas versões 6.14.0 e superiores do Amazon EMR no EKS.

Com os pontos de verificação do Flink, cada tarefa produz um snapshot do próprio estado que o Flink grava em um armazenamento distribuído, como o Amazon S3. Em casos de recuperação, as tarefas restauram o estado no armazenamento distribuído. O armazenamento distribuído oferece tolerância a falhas e pode redistribuir o estado durante o reajuste de escala por ser acessível a todos os nós.

No entanto, um armazenamento distribuído remoto também tem uma desvantagem: todas as tarefas devem ler seu estado de um local remoto na rede. Isso pode resultar em tempos de recuperação longos para estados grandes durante operações de recuperação ou ajuste de escala de tarefas.

Esse problema de tempo de recuperação longo é resolvido pela recuperação local de tarefas. As tarefas gravam seu estado no ponto de verificação em um armazenamento secundário local para a tarefa, como em um disco local. Elas também armazenam seu estado no armazenamento principal, ou no Amazon S3, no nosso caso. Durante a recuperação, o agendador programa as tarefas no mesmo gerenciador de tarefas em que elas foram executadas anteriormente, para que possam se recuperar do armazenamento de estado local em vez de ler do armazenamento de estado remoto. Para obter mais informações, consulte Task-Local Recovery na Documentação do Apache Flink.

Nossos testes de benchmark com exemplos de trabalhos mostraram que o tempo de recuperação foi reduzido de minutos para alguns segundos com a recuperação local da tarefa ativada.

Para habilitar a recuperação local de tarefas, defina as configurações a seguir no seu arquivo flink-conf.yaml. Especifique o valor do intervalo do ponto de verificação em milissegundos.

```
state.backend.local-recovery: true
state.backend: hasmap or rocksdb
state.checkpoints.dir: s3://STORAGE-BUCKET-PATH/checkpoint
execution.checkpointing.interval: 15000
```

Recuperação local de tarefas por meio da montagem de volume do Amazon EBS



Note

A recuperação local de tarefas por meio do Amazon EBS é compatível com o Flink nas versões 6.15.0 e superiores do Amazon EMR no EKS.

Com o Flink para Amazon EMR no EKS, você pode provisionar automaticamente os volumes do Amazon EBS nos pods do TaskManager para a recuperação local de tarefas. A montagem de sobreposição padrão vem com um volume de 10 GB, o que é suficiente para trabalhos com um estado inferior. Trabalhos com estados grandes podem habilitar a opção de montagem automática de volume do EBS. Os pods do TaskManager são criados e montados automaticamente durante a criação do pod e removidos durante a exclusão do pod.

Use as etapas a seguir para habilitar a montagem automática de volume do EBS para o Flink no Amazon EMR no EKS.

Exporte os valores das seguintes variáveis que você usará nas próximas etapas:

```
export AWS_REGION=aa-example-1
export FLINK_EKS_CLUSTER_NAME=my-cluster
export AWS_ACCOUNT_ID=111122223333
```

Crie ou atualize um arquivo kubeconfig YAML para o cluster.

```
aws eks update-kubeconfig --name $FLINK_EKS_CLUSTER_NAME --region $AWS_REGION
```

Crie uma conta de serviço do IAM para o driver de interface de armazenamento de contêiner 3. (CSI) do Amazon EBS no cluster do Amazon EKS.

```
eksctl create iamserviceaccount \
   --name ebs-csi-controller-sa \
   --namespace kube-system \
   --region $AWS_REGION \
   --cluster $FLINK_EKS_CLUSTER_NAME\
   --role-name TLR_${AWS_REGION}_${FLINK_EKS_CLUSTER_NAME} \
   --role-only \
   --attach-policy-arn arn:aws:iam::aws:policy/service-role/
AmazonEBSCSIDriverPolicy \
```

```
--approve
```

4. Crie o driver CSI do Amazon EBS com o seguinte comando:

```
eksctl create addon \
    --name aws-ebs-csi-driver \
    --region $AWS_REGION \
    --cluster $FLINK_EKS_CLUSTER_NAME \
    --service-account-role-arn arn:aws:iam::${AWS_ACCOUNT_ID}:role/TLR_
${AWS_REGION}_${FLINK_EKS_CLUSTER_NAME}
```

5. Crie a classe de armazenamento do Amazon EBS com o seguinte comando:

```
cat # EOF # storage-class.yaml
apiVersion: storage.k8s.io/v1
kind: StorageClass
metadata:
   name: ebs-sc
provisioner: ebs.csi.aws.com
volumeBindingMode: WaitForFirstConsumer
EOF
```

Depois, aplique a classe:

```
kubectl apply -f storage-class.yaml
```

6. O Helm instala o operador do Kubernetes para Flink do Amazon EMR com opções para criar uma conta de serviço. Isso cria o emr-containers-sa-flink para ser usado na implantação do Flink.

```
helm install flink-kubernetes-operator flink-kubernetes-operator/ \
    --set jobServiceAccount.create=true \
    --set rbac.jobRole.create=true \
    --set rbac.jobRoleBinding.create=true
```

7. Para enviar o trabalho do Flink e habilitar o provisionamento automático de volumes do EBS para a recuperação local de tarefas, defina as configurações a seguir no seu arquivo flink-conf.yaml. Ajuste o limite de tamanho para o tamanho do estado do trabalho. Defina serviceAccount como emr-containers-sa-flink. Especifique o valor do intervalo do ponto de verificação em milissegundos. Omita o executionRoleArn.

```
flinkConfiguration:
    task.local-recovery.ebs.enable: true
    kubernetes.taskmanager.local-recovery.persistentVolumeClaim.sizeLimit: 10Gi
    state.checkpoints.dir: s3://BUCKET-PATH/checkpoint
    state.backend.local-recovery: true
    state.backend: hasmap or rocksdb
    state.backend.incremental: "true"
    execution.checkpointing.interval: 15000
  serviceAccount: emr-containers-sa-flink
```

Quando estiver tudo pronto para excluir o plug-in do driver CSI do Amazon EBS, use os seguintes comandos:

```
# Detach Attached Policy
  aws iam detach-role-policy --role-name TLR ${$AWS_REGION} ${FLINK_EKS_CLUSTER_NAME}
 --policy-arn arn:aws:iam::aws:policy/service-role/AmazonEBSCSIDriverPolicy
  # Delete the created Role
  aws iam delete-role --role-name TLR_${$AWS_REGION}_${FLINK_EKS_CLUSTER_NAME}
  # Delete the created service account
  eksctl delete iamserviceaccount --name ebs-csi-controller-sa --namespace kube-system
 --cluster $FLINK_EKS_CLUSTER_NAME --region $AWS_REGION
  # Delete Addon
  eksctl delete addon --name aws-ebs-csi-driver --cluster $FLINK_EKS_CLUSTER_NAME --
region $AWS_REGION
  # Delete the EBS storage class
  kubectl delete -f storage-class.yaml
```

Ponto de verificação incremental genérico baseado em logs



O ponto de verificação incremental genérico baseado em logs é compatível com o Flink para as versões 6.14.0 e superiores do Amazon EMR no EKS.

O recurso de ponto de verificação incremental genérico baseado em logs foi adicionado ao Flink 1.16 para melhorar a velocidade dos pontos de verificação. Um intervalo de ponto de verificação mais rápido costuma resultar em uma redução do trabalho de recuperação porque menos eventos precisam ser reprocessados após a recuperação. Para obter mais informações, consulte Improving

speed and stability of checkpointing with generic log-based incremental checkpoints no blog do Apache Flink.

Com exemplos de trabalhos, nossos testes de comparação mostraram que o tempo do ponto de verificação foi reduzido de minutos para alguns segundos com o ponto de verificação incremental genérico baseado em logs.

Para habilitar os pontos de verificação incrementais genéricos baseados em logs, defina as configurações a seguir no seu arquivo flink-conf.yaml. Especifique o valor do intervalo do ponto de verificação em milissegundos.

```
state.backend.changelog.enabled: true
state.backend.changelog.storage: filesystem
dstl.dfs.base-path: s3://bucket-path/changelog
state.backend.local-recovery: true
state.backend: rocksdb
state.checkpoints.dir: s3://bucket-path/checkpoint
execution.checkpointing.interval: 15000
```

Recuperação refinada



O suporte de recuperação refinada para o agendador padrão é compatível com o Flink nas versões 6.14.0 e superiores do Amazon EMR no EKS. O suporte de recuperação refinada no agendador adaptável está disponível com o Flink nas versões 6.15.0 e superiores do Amazon EMR no EKS.

Quando uma tarefa falha durante a execução, o Flink redefine todo o gráfico de execução e aciona a reexecução completa com base no último ponto de verificação concluído. Isso é mais caro do que apenas reexecutar as tarefas que falharam. A recuperação refinada reinicia somente o componente conectado ao pipeline da tarefa que falhou. No exemplo a seguir, o gráfico do trabalho tem cinco vértices (de A a E). É feito o pipeline de todas as conexões entre os vértices com distribuição pontual, e o parallelism. default do trabalho é definido como 2.

```
A # B # C # D # E
```

Neste exemplo, há um total de dez tarefas em execução. O primeiro pipeline (de a1 a e1) é executado em um TaskManager (TM1) e o segundo pipeline (de a2 a e2) é executado em outro TaskManager (TM2).

```
a1 # b1 # c1 # d1 # e1
a2 # b2 # c2 # d2 # e2
```

Há dois componentes conectados por pipeline: a1 # e1 e a2 # e2. Se TM1 ou TM2 falhar, a falha afetará somente as cinco tarefas no pipeline em que TaskManager estava em execução. A estratégia de reinicialização só inicia o componente do pipeline afetado.

A recuperação refinada funciona somente com trabalhos perfeitamente paralelos do Flink. Não é compatível com operações de keyBy() ou redistribute(). Para obter mais informações, consulte FLIP-1: Fine Grained Recovery from Task Failures no projeto Jira Flink Improvement Proposal.

Para habilitar a recuperação refinada, defina as configurações a seguir no seu arquivo flinkconf.yaml.

```
jobmanager.execution.failover-strategy: region
restart-strategy: exponential-delay or fixed-delay
```

Mecanismo de reinício combinado no agendador adaptável



Note

O mecanismo de reinicialização combinado no agendador adaptável é compatível com o Flink nas versões 6.15.0 e superiores do Amazon EMR no EKS.

O agendador adaptável pode ajustar o paralelismo do trabalho com base nos slots disponíveis. Ele reduz automaticamente o paralelismo se não houver slots suficientes disponíveis para atender ao paralelismo do trabalho configurado. Se novos slots ficarem disponíveis, o trabalho aumentará a escala verticalmente mais uma vez para o paralelismo do trabalho configurado. Um agendador adaptável evita o tempo de inatividade no trabalho quando não há recursos suficientes disponíveis. Esse é o agendador compatível com o escalador automático do Flink. Recomendamos o agendador adaptável com o Flink no Amazon EMR por esses motivos. No entanto, os agendadores adaptáveis podem fazer várias reinicializações em um curto período; uma reinicialização para cada novo recurso adicionado. Isso pode levar a uma queda na performance do trabalho.

Com o Amazon EMR 6.15.0 e versões superiores, o Flink tem um mecanismo de reinício combinado no agendador adaptável que abre uma janela de reinicialização quando o primeiro recurso é adicionado e aguarda pelo intervalo configurado da janela do padrão de um minuto. Ele executa uma única reinicialização quando há recursos suficientes disponíveis para executar o trabalho com o paralelismo configurado ou quando o intervalo expira.

Com os exemplos de trabalhos, nossos testes de comparação mostraram que esse recurso processa 10% a mais dos registros do que o comportamento padrão ao usar o agendador adaptável e o escalador automático do Flink.

Para habilitar o mecanismo de reinício combinado, defina as configurações a seguir no seu arquivo flink-conf.yaml.

```
jobmanager.adaptive-scheduler.combined-restart.enabled: true
jobmanager.adaptive-scheduler.combined-restart.window-interval: 1m
```

Desativação tranquila de instâncias spot com Flink no Amazon EMR no EKS

O Flink com o Amazon EMR no EKS pode melhorar o tempo de reinicialização de trabalhos durante as operações de recuperação e ajuste de escala de tarefas.

Visão geral

As versões 6.15.0 e superiores do Amazon EMR no EKS oferecem suporte à desativação tranquila de gerenciadores de tarefas em instâncias spot no Amazon EMR no EKS com o Flink para Apache. Como parte desse recurso, o Amazon EMR no EKS com o Flink fornece as funcionalidades a seguir.

- Just-in-time ponto de verificação As tarefas de streaming do Flink podem responder à
 interrupção da Instância Spot, realizar o ponto de verificação just-in-time (JIT) das tarefas em
 execução e impedir o agendamento de tarefas adicionais nessas Instâncias Spot. O ponto de
 verificação de JIT é compatível com agendadores padrão e adaptáveis.
- Mecanismo de reinicialização combinado: um mecanismo de reinicialização combinado faz
 a melhor tentativa de reiniciar o trabalho depois que ele atinge o paralelismo de recursos de
 destino ou o final da janela configurada atual. Isso também evita reinicializações consecutivas de
 trabalhos que podem ser causadas por vários encerramentos de instâncias spot. O mecanismo de
 reinicialização combinado está disponível somente com o agendador adaptável.

Desativação tranquila 71

Esses recursos fornecem os seguintes benefícios:

- Você pode aproveitar as instâncias spot para executar gerenciadores de tarefas e reduzir os gastos com clusters.
- O aprimoramento de liveness do gerenciador de tarefas da instância spot resulta em maior resiliência e agendamento de trabalhos mais eficiente.
- Seus trabalhos do Flink terão mais tempo de atividade porque haverá menos reinicializações após o encerramento da instância spot.

Como funciona a desativação normal

Considere o seguinte exemplo: você provisiona um cluster do Amazon EMR no EKS executando o Flink para Apache e especifica nós sob demanda para o gerenciador de trabalhos e nós de instância spot para o gerenciador de tarefas. Dois minutos antes do encerramento, o gerenciador de tarefas recebe um aviso de interrupção.

Nesse cenário, o gerenciador de trabalhos manipularia o sinal de interrupção da instância spot, bloquearia o agendamento de tarefas adicionais na instância spot e iniciaria o ponto de verificação de JIT do trabalho de streaming.

Em seguida, o gerenciador de trabalhos reiniciaria o gráfico do trabalho somente depois que houvesse disponibilidade suficiente de novos recursos para satisfazer o paralelismo atual de trabalhos na janela atual do intervalo de reinicialização. O intervalo da janela de reinicialização é decidido com base na duração da substituição da instância spot, na criação de novos pods do gerenciador de tarefas e no registro no gerenciador de trabalhos.

Pré-requisitos

Para usar a desativação normal, crie e execute um trabalho de streaming em um cluster do Amazon EMR no EKS executando o Apache Flink. Habilite o agendador adaptável e os gerenciadores de tarefas agendados em, pelo menos, uma instância spot, conforme mostrado no exemplo a seguir. Você deve usar nós sob demanda para o gerenciador de trabalhos e pode usar nós sob demanda para os gerenciadores de tarefas, desde que também haja pelo menos uma instância spot.

apiVersion: flink.apache.org/v1beta1

kind: FlinkDeployment

metadata:

name: deployment_name

spec:

Desativação tranquila 72

```
flinkVersion: v1_17
flinkConfiguration:
  taskmanager.numberOfTaskSlots: "2"
  cluster.taskmanager.graceful-decommission.enabled: "true"
  execution.checkpointing.interval: "240s"
  jobmanager.adaptive-scheduler.combined-restart.enabled: "true"
  jobmanager.adaptive-scheduler.combined-restart.window-interval : "1m"
serviceAccount: flink
jobManager:
  resource:
   memory: "2048m"
    cpu: 1
  nodeSelector:
    'eks.amazonaws.com/capacityType': 'ON_DEMAND'
taskManager:
  resource:
   memory: "2048m"
    cpu: 1
  nodeSelector:
    'eks.amazonaws.com/capacityType': 'SPOT'
job:
  jarURI: flink_job_jar_path
```

Configuração

Esta seção aborda a maioria das configurações que você pode especificar para as suas necessidades de desativação.

Chave	Descrição	Valor padrão	Valores aceitos
cluster.t askmanage r.gracefu l-decommi ssion.ena bled	Habilita a desativação tranquila do gerenciador de tarefas.	true	true, false
<pre>jobmanage r.adaptiv e-schedul er.combin</pre>	Habilita o mecanismo de reinicial ização combinado no agendador adaptável.	false	true, false

Desativação tranquila 73

Chave	Descrição	Valor padrão	Valores aceitos
ed-restar t.enabled			
<pre>jobmanage r.adaptiv e-schedul er.combin ed-restar t.window- interval</pre>	O intervalo combinado da janela de reinicialização para efetuar reinicial izações mescladas do trabalho. Um número inteiro sem uma unidade é interpretado como milissegundos.	1m	Exemplos: 30, 60s, 3m, 1h

Uso do Autoscaler para aplicações do Flink

O Autoscaler do operador pode ajudar a aliviar a contrapressão ao coletar métricas de trabalhos do Flink e ajustar automaticamente o paralelismo em nível de vértice do trabalho. Confira o seguinte exemplo de como a configuração pode se parecer:

```
apiVersion: flink.apache.org/v1beta1
kind: FlinkDeployment
metadata:
...
spec:
...
flinkVersion: v1_18
flinkConfiguration:
job.autoscaler.enabled: "true"
job.autoscaler.stabilization.interval: 1m
job.autoscaler.metrics.window: 5m
job.autoscaler.target.utilization: "0.6"
job.autoscaler.target.utilization.boundary: "0.2"
job.autoscaler.restart.time: 2m
job.autoscaler.catch-up.duration: 5m
pipeline.max-parallelism: "720"
...
```

Essa configuração usa valores padrão para a versão mais recente do Amazon EMR. Se você usar outras versões, é possível que tenha valores diferentes.

Uso do Autoscaler 74

Note

A partir do Amazon EMR 7.2.0, você não precisa incluir o prefixo kubernetes.operator na sua configuração. Se você usa a versão 7.1.0 ou inferior, deve usar o prefixo antes de cada configuração. Por exemplo, você deve especificar kubernetes.operator.job.autoscaler.scaling.enabled.

A seguir, são apresentadas as opções de configuração para o Autoscaler.

- job.autoscaler.scaling.enabled: especifica se a execução do ajuste de escala de vértices deve ser habilitada pelo escalador automático. O padrão é true. Se você desabilitar essa configuração, o escalador automático coletará apenas métricas e avaliará o paralelismo sugerido para cada vértice, mas não atualizará os trabalhos.
- job.autoscaler.stabilization.interval: o período de estabilização no qual nenhuma nova escalabilidade será executada. O padrão é de cinco minutos.
- job.autoscaler.metrics.window: o tamanho da janela de agregação de métricas de escalabilidade. Quanto mais ampla for a janela, mais harmoniosa e estável ela será, mas o Autoscaler pode demorar mais para reagir a alterações repentinas de carga. O padrão é de 15 minutos. Recomendamos que você experimente usando um valor entre 3 e 60 minutos.
- job.autoscaler.target.utilization: a utilização desejada para o vértice para fornecer uma performance de trabalho estável e algum buffer para flutuações de carga. O padrão é 0.7 atingir 70% utilization/load para os vértices do trabalho.
- job.autoscaler.target.utilization.boundary: o limite da utilização desejada para o vértice, que serve como buffer extra para evitar uma escalabilidade imediata em flutuações de carga. O padrão é 0.3, o que significa que é permitido um desvio de 30% da utilização desejada antes do acionamento de uma ação de ajuste de escala.
- ob.autoscaler.restart.time: o tempo esperado para reiniciar a aplicação. O padrão é de cinco minutos.
- job.autoscaler.catch-up.duration: o tempo esperado para a recuperação, ou seja, o processamento total de qualquer backlog após a conclusão de uma operação de escalabilidade. O padrão é de cinco minutos. Ao reduzir a duração da recuperação, é necessário que o Autoscaler reserve uma capacidade extra para as ações de escalabilidade.
- pipeline.max-parallelism: o paralelismo máximo que o Autoscaler pode usar. O Autoscaler ignora esse limite se ele for maior que o paralelismo máximo configurado na configuração do

Uso do Autoscaler 75 Flink ou diretamente em cada operador. O padrão é -1. Observe que o Autoscaler calcula o paralelismo como um divisor do número de paralelismo máximo, portanto, é recomendado escolher configurações de paralelismo máximo com muitos divisores em vez de confiar nos padrões fornecidos pelo Flink. Recomendamos usar múltiplos de 60 para esta configuração, por exemplo, 120, 180, 240, 360, 720 etc.

Para obter uma página de referência de configurações mais detalhada, consulte <u>Autoscaler</u> Configuration.

Ajuste automático dos parâmetros do escalador automático

Esta seção descreve o comportamento de ajuste automático para várias versões do Amazon EMR. Ela também detalha as diferentes configurações de ajuste de escala automático.



O Amazon EMR 7.2.0 e versões posteriores usam a configuração de código aberto job.autoscaler.restart.time-tracking.enabled para habilitar a estimativa do tempo de redimensionamento. A estimativa do tempo de redimensionamento tem a mesma funcionalidade do ajuste automático do Amazon EMR, então você não precisa atribuir manualmente valores empíricos ao horário de reinicialização.

Você ainda pode usar o ajuste automático do Amazon EMR se estiver utilizando o Amazon EMR 7.1.0 ou inferior.

7.2.0 and higher

O Amazon EMR 7.2.0 e versões posteriores medem o tempo real de reinicialização necessário para aplicar decisões de ajuste de escala automático. Nas versões 7.1.0 e inferiores, você precisava usar a configuração job.autoscaler.restart.time para configurar manualmente o tempo máximo estimado de reinicialização. Ao usar a configuração job.autoscaler.restart.time-tracking.enabled, você só precisa inserir um horário de reinicialização para o primeiro ajuste de escala. Depois, o operador registra o tempo real de reinicialização e o usará para ajustes de escala subsequentes.

Para habilitar esse rastreamento, use o seguinte comando:

job.autoscaler.restart.time-tracking.enabled: true

A seguir estão as configurações relacionadas para a estimativa do tempo de redimensionamento.

Configuração	Obrigatório	Padrão	Descrição
job.autoscaler.restart.time- tracking.enabled	Não	Falso	Indica se o escalador automático do Flink deve ajustar automaticamente as configurações ao longo do tempo para otimizar as decisões de escalabilidade. Observe que o escalador automático só pode ajustar automaticamente o parâmetro restart.time do escalador automático.
job.autoscaler.restart.time	Não	5 minutos	O tempo de reinicialização esperado que o Amazon EMR no EKS usa até que o operador possa determina r o tempo real de reinicial ização com base nos ajustes de escala anteriores.
job.autoscaler.restart.time- tracking.limit	Não	15 minutos	O tempo máximo de reinicial ização observado quando job.autoscaler.res tart.time-tracking .enabled está definido como true.

Confira este exemplo de especificação de implantação que você pode usar para testar a estimativa do tempo de redimensionamento:

apiVersion: flink.apache.org/v1beta1

kind: FlinkDeployment

metadata:

```
name: autoscaling-example
spec:
 flinkVersion: v1_18
  flinkConfiguration:
    # Autoscaler parameters
    job.autoscaler.enabled: "true"
    job.autoscaler.scaling.enabled: "true"
    job.autoscaler.stabilization.interval: "5s"
    job.autoscaler.metrics.window: "1m"
    job.autoscaler.restart.time-tracking.enabled: "true"
    job.autoscaler.restart.time: "2m"
    job.autoscaler.restart.time-tracking.limit: "10m"
    jobmanager.scheduler: adaptive
    taskmanager.numberOfTaskSlots: "1"
    pipeline.max-parallelism: "12"
  executionRoleArn: < JOB ARN>
  emrReleaseLabel: emr-7.10.0-flink-latest
  jobManager:
    highAvailabilityEnabled: false
    storageDir: s3://<s3_bucket>/flink/autoscaling/ha/
    replicas: 1
    resource:
      memory: "1024m"
      cpu: 0.5
  taskManager:
    resource:
      memory: "1024m"
      cpu: 0.5
  job:
    jarURI: s3://<s3_bucket>/some-job-with-back-pressure
    parallelism: 1
    upgradeMode: stateless
```

Para simular a contrapressão, use a especificação de implantação a seguir.

```
job:
    jarURI: s3://<s3_bucket>/pyflink-script.py
    entryClass: "org.apache.flink.client.python.PythonDriver"
    args: ["-py", "/opt/flink/usrlib/pyflink-script.py"]
```

```
parallelism: 1
upgradeMode: stateless
```

Faça upload do script Python a seguir no bucket do S3.

```
import logging
import sys
import time
import random
from pyflink.datastream import StreamExecutionEnvironment
from pyflink.table import StreamTableEnvironment
TABLE_NAME="orders"
QUERY=f"""
CREATE TABLE {TABLE_NAME} (
  id INT,
  order_time AS CURRENT_TIMESTAMP,
 WATERMARK FOR order_time AS order_time - INTERVAL '5' SECONDS
)
WITH (
  'connector' = 'datagen',
  'rows-per-second'='10',
  'fields.id.kind'='random',
  'fields.id.min'='1',
  'fields.id.max'='100'
);
.....
def create_backpressure(i):
    time.sleep(2)
    return i
def autoscaling_demo():
    env = StreamExecutionEnvironment.get_execution_environment()
    t_env = StreamTableEnvironment.create(env)
    t_env.execute_sql(QUERY)
    res_table = t_env.from_path(TABLE_NAME)
    stream = t_env.to_data_stream(res_table) \
      .shuffle().map(lambda x: create_backpressure(x))\
      .print()
    env.execute("Autoscaling demo")
```

```
if __name__ == '__main__':
    logging.basicConfig(stream=sys.stdout, level=logging.INFO, format="%(message)s")
    autoscaling_demo()
```

Para verificar se a estimativa do tempo de redimensionamento está funcionando, confira se o registro em log em nível de DEBUG do operador do Flink está habilitado. O exemplo abaixo demonstra como atualizar o arquivo values. yaml do chart do Helm. Em seguida, reinstale o chart do Helm atualizado e execute o trabalho do Flink novamente.

```
log4j-operator.properties: |+
    # Flink Operator Logging Overrides
    rootLogger.level = DEBUG
```

Obtenha o nome do pod líder.

```
ip=$(kubectl get configmap -n $NAMESPACE < job-name > -cluster-config-map -o json | jq
-r ".data[\"org.apache.flink.k8s.leader.restserver\"]" | awk -F: '{print $2}' | awk
-F '/' '{print $3}')

kubectl get pods -n $NAMESPACE -o json | jq -r ".items[] | select(.status.podIP ==
\"$ip\") | .metadata.name"
```

Execute o comando a seguir para obter o tempo real de reinicialização usado nas avaliações de métricas.

```
kubectl logs <FLINK-OPERATOR-POD-NAME> -c flink-kubernetes-operator -n <OPERATOR-
NAMESPACE> -f | grep "Restart time used in scaling summary computation"
```

Você deve ver logs semelhantes aos mostrados a seguir. Observe que somente o primeiro ajuste de escala usa job.autoscaler.restart.time. Os ajustes de escala subsequentes usam o tempo de reinicialização observado.

```
2024-05-16 17:17:32,590 o.a.f.a.ScalingExecutor [DEBUG][default/autoscaler-example] Restart time used in scaling summary computation: PT2M 2024-05-16 17:19:03,787 o.a.f.a.ScalingExecutor [DEBUG][default/autoscaler-example] Restart time used in scaling summary computation: PT14S 2024-05-16 17:19:18,976 o.a.f.a.ScalingExecutor [DEBUG][default/autoscaler-example] Restart time used in scaling summary computation: PT14S
```

2024-05-16 17:20:50,283 o.a.f.a.ScalingExecutor [DEBUG][default/autoscaler-example] Restart time used in scaling summary computation: PT14S
2024-05-16 17:22:21,691 o.a.f.a.ScalingExecutor [DEBUG][default/autoscaler-example] Restart time used in scaling summary computation: PT14S

7.0.0 and 7.1.0

O escalador automático integrado do Flink e de código aberto usa várias métricas para tomar as melhores decisões de ajuste de escala. No entanto, os valores padrão que ele usa para seus cálculos devem ser aplicáveis à maioria das workloads e podem não ser ideais para um determinado trabalho. O recurso de ajuste automático adicionado à versão do Amazon EMR no EKS do operador do Flink analisa as tendências históricas observadas em métricas específicas capturadas e, em seguida, tenta calcular o valor mais adequado para o trabalho em questão.

Configuração	Obrigatório	Padrão	Descrição
kubernetes.operator.job.aut oscaler.autotune.enable	Não	Falso	Indica se o escalador automático do Flink deve ajustar automaticamente as configurações ao longo do tempo para otimizar as decisões de escalabilidade do escalador automático. Atualmente, o escalador automático só pode ajustar automáticamente seu parâmetro restart.time .
kubernetes.operator.job.aut oscaler.autotune.metrics.hi story.max.count	Não	3	Indica quantas métricas históricas do Amazon EMR no EKS o escalador automático mantém no mapa de configuração de métricas do Amazon EMR no EKS.
kubernetes.operator.job.aut oscaler.autotune.metrics.re start.count	Não	3	Indica quantas reinicial izações o escalador automático executa antes de

Configuração	Obrigatório	Padrão	Descrição
			começar a calcular o tempo médio de reinicialização de um determinado trabalho.

Para habilitar o ajuste automático, você deve concluir as seguintes etapas:

- Defina kubernetes.operator.job.autoscaler.autotune.enable: como true
- Defina metrics.job.status.enable: como TOTAL_TIME
- Seguir a configuração em <u>Using Autoscaler for Flink applications</u> para habilitar o ajuste de escala automático.

Confira a seguir um exemplo de especificação de implantação que você pode usar para experimentar o ajuste automático.

```
apiVersion: flink.apache.org/v1beta1
kind: FlinkDeployment
metadata:
  name: autoscaling-example
spec:
  flinkVersion: v1_18
  flinkConfiguration:
    # Autotuning parameters
    kubernetes.operator.job.autoscaler.autotune.enable: "true"
    kubernetes.operator.job.autoscaler.autotune.metrics.history.max.count: "2"
    kubernetes.operator.job.autoscaler.autotune.metrics.restart.count: "1"
   metrics.job.status.enable: TOTAL_TIME
    # Autoscaler parameters
    kubernetes.operator.job.autoscaler.enabled: "true"
    kubernetes.operator.job.autoscaler.scaling.enabled: "true"
    kubernetes.operator.job.autoscaler.stabilization.interval: "5s"
    kubernetes.operator.job.autoscaler.metrics.window: "1m"
    jobmanager.scheduler: adaptive
    taskmanager.numberOfTaskSlots: "1"
    state.savepoints.dir: s3://<S3_bucket>/autoscaling/savepoint/
```

```
state.checkpoints.dir: s3://<S3_bucket>/flink/autoscaling/checkpoint/
  pipeline.max-parallelism: "4"
executionRoleArn: <JOB ARN>
emrReleaseLabel: emr-6.14.0-flink-latest
jobManager:
  highAvailabilityEnabled: true
  storageDir: s3://<S3_bucket>/flink/autoscaling/ha/
  replicas: 1
  resource:
   memory: "1024m"
    cpu: 0.5
taskManager:
  resource:
   memory: "1024m"
   cpu: 0.5
job:
  jarURI: s3://<S3_bucket>/some-job-with-back-pressure
  parallelism: 1
  upgradeMode: last-state
```

Para simular a contrapressão, use a especificação de implantação a seguir.

```
job:
   jarURI: s3://<S3_bucket>/pyflink-script.py
   entryClass: "org.apache.flink.client.python.PythonDriver"
   args: ["-py", "/opt/flink/usrlib/pyflink-script.py"]
   parallelism: 1
   upgradeMode: last-state
```

Faça upload do script Python a seguir no bucket do S3.

```
import logging
import sys
import time
import random

from pyflink.datastream import StreamExecutionEnvironment
from pyflink.table import StreamTableEnvironment

TABLE_NAME="orders"
QUERY=f"""
CREATE TABLE {TABLE_NAME} (
```

```
id INT,
  order_time AS CURRENT_TIMESTAMP,
  WATERMARK FOR order_time AS order_time - INTERVAL '5' SECONDS
)
WITH (
  'connector' = 'datagen',
  'rows-per-second'='10',
  'fields.id.kind'='random',
  'fields.id.min'='1',
  'fields.id.max'='100'
);
.....
def create_backpressure(i):
    time.sleep(2)
    return i
def autoscaling_demo():
    env = StreamExecutionEnvironment.get_execution_environment()
    t_env = StreamTableEnvironment.create(env)
    t_env.execute_sql(QUERY)
    res_table = t_env.from_path(TABLE_NAME)
    stream = t_env.to_data_stream(res_table) \
      .shuffle().map(lambda x: create_backpressure(x))\
      .print()
    env.execute("Autoscaling demo")
if __name__ == '__main__':
    logging.basicConfig(stream=sys.stdout, level=logging.INFO, format="%(message)s")
    autoscaling_demo()
```

Para verificar se o escalador automático está funcionando, use os comandos a seguir. Observe que você deve usar as informações do seu próprio pod líder no operador do Flink.

Primeiro, obtenha o nome do seu pod líder.

```
ip=$(kubectl get configmap -n $NAMESPACE <job-name>-cluster-config-map -o json | jq
-r ".data[\"org.apache.flink.k8s.leader.restserver\"]" | awk -F: '{print $2}' | awk
-F '/' '{print $3}')
kubectl get pods -n $NAMESPACE -o json | jq -r ".items[] | select(.status.podIP ==
\"$ip\") | .metadata.name"
```

Depois de ter o nome do pod líder, será possível executar o comando a seguir.

```
kubectl logs -n $NAMESPACE -c flink-kubernetes-operator --follow <YOUR-FLINK-
OPERATOR-POD-NAME> | grep -E 'EmrEks|autotun|calculating|restart|autoscaler'
```

Você deve ver logs semelhantes aos mostrados a seguir.

```
[m[33m2023-09-13 20:10:35,941[m [36mc.a.c.f.k.o.a.EmrEksMetricsAutotuner[m
    [36m[DEBUG][flink/autoscaling-example] Using the latest
    Emr Eks Metric for calculating restart.time for autotuning:
    EmrEksMetrics(restartMetric=RestartMetric(restartingTime=65, numRestarts=1))

[m[33m2023-09-13 20:10:35,941[m [36mc.a.c.f.k.o.a.EmrEksMetricsAutotuner[m
    [32m[INFO ][flink/autoscaling-example] Calculated average restart.time metric via autotuning to be: PT0.065S
```

Manutenção e solução de problemas para trabalhos do Flink no Amazon EMR no EKS

As seções a seguir descrevem como manter seus trabalhos de execução prolongada do Flink e fornecem orientações sobre como solucionar alguns problemas comuns em trabalhos do Flink.

Manutenção de aplicações do Flink

Tópicos

Modos de atualização

As aplicações do Flink costumam ser projetadas para execuções por longos períodos, como semanas, meses ou até anos. Como acontece com todos os serviços de execução prolongada, as aplicações de streaming do Flink precisam de manutenção. Isso inclui correções de erros, melhorias e migração para um cluster do Flink de uma versão mais recente.

Quando o arquivo spec muda para os recursos FlinkDeployment e FlinkSessionJob, você precisa atualizar a aplicação em execução. Para isso, o operador interrompe o trabalho em execução (a menos que já esteja suspenso) e o reimplanta com o arquivo spec mais recente e, em aplicações com estado, o estado da execução anterior.

Os usuários controlam como gerenciar o estado quando aplicações com estado param e são restauradas com a configuração upgradeMode do JobSpec.

Modos de atualização

Introdução opcional

Sem estado

Atualizações de aplicações sem estado começando do estado vazio.

Último estado

Atualizações rápidas em qualquer estado da aplicação (mesmo em trabalhos com falha) não exigem um trabalho íntegro, pois sempre usam o ponto de verificação bem-sucedido mais recente. A recuperação manual pode ser necessária se os metadados de HA forem perdidos. Para limitar o tempo em que o trabalho pode efetuar fallback ao escolher o ponto de verificação mais recente, você pode configurar kubernetes.operator.job.upgrade.last-state.max.allowed.checkpoint.age. Se o ponto de verificação for mais antigo que o valor configurado, um ponto de salvamento será usado para trabalhos íntegros. Isso não é compatível com o modo Sessão.

Ponto de salvamento

Use o ponto de salvamento para atualização, fornecendo máxima segurança e possibilidade de servir como ponto. backup/fork O ponto de salvamento será criado durante o processo de atualização. Observe que o trabalho do Flink precisa estar em execução para permitir que o ponto de salvamento seja criado. Se o trabalho não estiver íntegro, o último ponto de verificação será usado (a menos que kubernetes.operator.job.upgrade). last-state-fallback.enabled está definido como false). Se o ponto de verificação mais recente não estiver disponível, a atualização do trabalho falhará.

Solução de problemas

Esta seção descreve como solucionar problemas com o Amazon EMR no EKS. Para obter informações sobre como solucionar problemas gerais com o Amazon EMR, consulte <u>Troubleshoot a cluster</u> no Guia de gerenciamento do Amazon EMR.

- Solução de problemas de trabalhos que usam PersistentVolumeClaims (PVC)
- Solução de problemas de escalabilidade automática vertical do Amazon EMR no EKS

Solução de problemas do operador do Spark do Amazon EMR no EKS

Solução de problemas do Apache Flink para Amazon EMR no EKS

Mapeamento de recursos não encontrado ao instalar o chart do Helm

Você pode encontrar a mensagem de erro apresentada a seguir ao instalar o chart do Helm.

```
Error: INSTALLATION FAILED: pulling from host 1234567890.dkr.ecr.us-
west-2.amazonaws.com failed with status code [manifests 6.13.0]: 403 Forbidden Error:
INSTALLATION FAILED: unable to build kubernetes objects from release manifest:
[resource mapping not found for name: "flink-operator-serving-cert" namespace: "<the
namespace to install your operator>" from "": no matches for kind "Certificate" in
version "cert-manager.io/v1"

ensure CRDs are installed first, resource mapping not found for name: "flink-operator-
selfsigned-issuer" namespace: "<the namespace to install your operator>" " from "": no
matches for kind "Issuer" in version "cert-manager.io/v1"

ensure CRDs are installed first].
```

Para resolver esse erro, instale o cert-manager para habilitar a adição do componente webhook. Você deve instalar o cert-manager em cada cluster do Amazon EKS que usa.

```
kubectl apply -f https://github.com/cert-manager/cert-manager/releases/download/v1.12.0
```

AWS service (Servico da AWS) erro de acesso negado

Se você ver um erroaccess denied, confirme se o perfil do IAM para operatorExecutionRoleArn no arquivo values.yaml do chart do Helm tem as permissões corretas. Além disso, certifique-se de que o perfil do IAM em executionRoleArn na sua especificação de FlinkDeployment tenha as permissões corretas.

FlinkDeployment está preso

Se seu FlinkDeployment ficar em um estado preso, use as etapas a seguir para forçar a exclusão da implantação:

Edite a execução da implantação.

```
kubectl edit -n Flink Namespace flinkdeployments/App Name
```

2. Remova este finalizador.

```
finalizers:
   - flinkdeployments.flink.apache.org/finalizer
```

3. Exclua a implantação.

```
kubectl delete -n Flink Namespace flinkdeployments/App Name
```

AWSBadRequestException problema s3a ao executar um aplicativo Flink em um opt-in Região da AWS

Se você executar um aplicativo Flink em um opt-in Região da AWS, poderá ver os seguintes erros:

```
Caused by: org.apache.hadoop.fs.s3a.AWSBadRequestException: getFileStatus on s3://flink.txt: com.amazonaws.services.s3.model.AmazonS3Exception: Bad Request (Service: Amazon S3; Status Code: 400; Error Code: 400 Bad Request; Request ID: ABCDEFGHIJKL; S3 Extended Request ID:

ABCDEFGHIJKLMNOP=; Proxy: null), S3 Extended Request ID: ABCDEFGHIJKLMNOP=:400 Bad Request: Bad Request
(Service: Amazon S3; Status Code: 400; Error Code: 400 Bad Request; Request ID: ABCDEFGHIJKL; S3 Extended Request ID: ABCDEFGHIJKLMNOP=; Proxy: null)
```

```
Caused by: org.apache.hadoop.fs.s3a.AWSBadRequestException: getS3Region on flink-application: software.amazon.awssdk.services.s3.model.S3Exception: null (Service: S3, Status Code: 400, Request ID: ABCDEFGHIJKLMNOP, Extended Request ID: ABCDEFGHIJKLMNOPQRST==):null: null (Service: S3, Status Code: 400, Request ID: ABCDEFGHIJKLMNOP, Extended Request ID: AH142uDNaTUFOus/5IIVNvSakBcMjMCH7dd37ky0vE6jhABCDEFGHIJKLMNOPQRST==)
```

Para corrigir esses erros, use a configuração a seguir no arquivo FlinkDeployment de definição.

```
spec:
  flinkConfiguration:
    taskmanager.numberOfTaskSlots: "2"
    fs.s3a.endpoint.region: OPT_IN_AWS_REGION_NAME
```

Também recomendamos que você use o provedor de SDKv2 credenciais:

```
fs.s3a.aws.credentials.provider:
software.amazon.awssdk.auth.credentials.WebIdentityTokenFileCredentialsProvider
```

Se você quiser usar o provedor de SDKv1 credenciais, certifique-se de que seu SDK seja compatível com sua região de inscrição. Para obter mais informações, consulte o <u>aws-sdk-java GitHub</u> repositório.

Se você obtiver o erro S3 AWSBadRequestException ao executar instruções SQL do Flink em uma região de ingresso, certifique-se de definir a configuração fs.s3a.endpoint.region:

OPT_IN_AWS_REGION_NAME na especificação de configuração do Flink.

S3A AWSBad RequestException ao executar um trabalho de sessão do Flink em regiões CN

Nas versões 6.15.0 à 7.2.0 do Amazon EMR, você pode encontrar as mensagens de erro a seguir ao executar um trabalho de sessão do Flink nas regiões da China. Estes incluem China (Pequim) e China (Ningxia):

```
Error:
 {"type":"org.apache.flink.kubernetes.operator.exception.ReconciliationException","message":"or
                    getFileStatus on s3://ABCDPath:
software.amazon.awssdk.services.s3.model.S3Exception: null (Service: S3, Status Code:
400, Request ID: ABCDEFGH, Extended Request ID:
                    ABCDEFGH:null: null (Service: S3, Status Code: 400, Request ID:
ABCDEFGH, Extended Request ID: ABCDEFGH", "additionalMetadata":{}, "throwableList":
 [{"type":"org.apache.hadoop.fs.s3a.AWSBadRequestException","message":"getFileStatus on
s3://ABCDPath: software.amazon.awssdk.services.s3.model.S3Exception:
                    null (Service: S3, Status Code: 400, Request ID: ABCDEFGH, Extended
Request ID: ABCDEFGH:null: null (Service: S3, Status Code: 400, Request ID: ABCDEFGH,
                    Extended Request ID: ABCDEFGH", "additionalMetadata":{}},
{"type": "software.amazon.awssdk.services.s3.model.S3Exception", "message": "null
 (Service: S3, Status Code: 400,
                    Request ID: ABCDEFGH, Extended Request ID:
ABCDEFGH", "additionalMetadata":{}}]}
```

Esse problema é conhecido. A equipe está trabalhando na correção dos operadores do Flink em todas essas versões de lançamento. No entanto, antes de concluirmos o patch, para corrigir esse erro, você precisa baixar o chart do Helm do operador do Flink, descompactá-lo (extrair o arquivo compactado) e fazer alterações na configuração do chart do Helm.

As etapas especificas são as seguintes:

1. Altere, especificamente, os diretórios da sua pasta local do chart do Helm e execute a linha de comando a seguir para extrair o chart do Helm e descompactá-lo (extraí-lo).

```
helm pull oci://public.ecr.aws/emr-on-eks/flink-kubernetes-operator \
--version $VERSION \
--namespace $NAMESPACE
```

```
tar -zxvf flink-kubernetes-operator-$VERSION.tgz
```

- 2. Vá até a pasta do chart do Helm e encontre o arquivo templates/flink-operator.yaml.
- 3. Localize flink-operator-config ConfigMap e adicione a seguinte fs.s3a.endpoint.region configuração noflink-conf.yaml. Por exemplo:

```
{{- if .Values.defaultConfiguration.create }}
apiVersion: v1
kind: ConfigMap
metadata:
   name: flink-operator-config
   namespace: {{ .Release.Namespace }}
   labels:
     {{- include "flink-operator.labels" . | nindent 4 }}
data:
   flink-conf.yaml: |+
fs.s3a.endpoint.region: {{ .Values.emrContainers.awsRegion }}
```

4. Instale o chart do Helm local e execute o trabalho.

Versões compatíveis para Amazon EMR no EKS com Apache Flink

O Apache Flink está disponível com as versões do Amazon EMR no EKS apresentadas a seguir. Para obter informações sobre todas as versões disponíveis, consulte <u>Versões do Amazon EMR no EKS</u>.

Rótulo da versão	Java	Flink	Operador do Flink
emr-7.2.0-flink-latest	17	1.18.1	-
emr-7.2.0-flink-k8s-operator-latest	11	-	1.8.0

Versões compatíveis 90

Rótulo da versão	Java	Flink	Operador do Flink
emr-7.1.0-flink-latest	17	1.18.1	-
emr-7.1.0-flink-k8s-operator-latest	11	-	1.6.1
emr-7.0.0-flink-latest	11	1.18.0	-
emr-7.0.0-flink-k8s-operator-latest	11	-	1.6.1
emr-6.15.0-flink-latest	11	1.17.1	-
emr-6.15.0-flink-k8s-operator-latest	11	-	1.6.0
emr-6.14.0-flink-latest	11	1.17.1	-
emr-6.14.0-flink-k8s-operator-latest	11	-	1.6.0
emr-6.13.0-flink-latest	11	1.17.0	-
emr-6.13.0-flink-k8s-operator-latest	11	-	1.5.0

Versões compatíveis 91

Execução de trabalhos do Spark com o Amazon EMR no **EKS**

A execução de um trabalho é uma unidade de trabalho, como um Spark jar, PySpark script ou consulta SparkSQL, que você envia para o Amazon EMR no EKS. Este tópico fornece uma visão geral do gerenciamento de execuções de trabalhos usando o AWS CLI, da visualização de execuções de trabalhos usando o console do Amazon EMR e da solução de problemas comuns de execução de trabalhos.

Observe que você não pode executar trabalhos do IPv6 Spark no Amazon EMR no EKS



Note

Antes de enviar uma execução de trabalho com o Amazon EMR no EKS, você deve concluir as etapas em Configuração do Amazon EMR no EKS.

Tópicos

- Execução de trabalhos do Spark com StartJobRun
- Execução de trabalhos do Spark com o operador do Spark
- Execução de trabalhos do Spark com spark-submit
- Uso do Apache Livy com o Amazon EMR no EKS
- Gerenciamento de execuções de trabalhos do Amazon EMR no EKS
- Uso de modelos de trabalho
- Uso de modelos de pod
- Uso de políticas de repetição de trabalho
- Uso da alternância de log de eventos do Spark
- Uso da alternância de log do contêiner do Spark
- Uso da escalabilidade automática vertical com trabalhos do Spark no Amazon EMR

Execução de trabalhos do Spark com StartJobRun

Esta seção inclui etapas detalhadas de configuração para preparar seu ambiente para executar trabalhos do Spark e, em seguida, fornece step-by-step instruções para enviar um trabalho executado com parâmetros especificados.

Tópicos

- Configuração do Amazon EMR no EKS
- Envio de uma execução de trabalho com StartJobRun
- Uso da classificação de envio de trabalho
- Usando a classificação padrão do contêiner Amazon EMR

Configuração do Amazon EMR no EKS

Conclua as tarefas apresentadas a seguir para se preparar para usar o Amazon EMR no EKS. Se você já se inscreveu na Amazon Web Services (AWS) e usa o Amazon EKS, está com quase tudo pronto para usar o Amazon EMR no EKS. Ignore qualquer tarefa que já tenha concluído.



Você também pode acompanhar o <u>Amazon EMR on EKS Workshop</u> para configurar todos os recursos obrigatórios para executar trabalhos do Spark no Amazon EMR no EKS. O workshop também fornece automação usando CloudFormation modelos para criar os recursos necessários para você começar. Para outros modelos e melhores práticas, consulte nosso Guia de práticas recomendadas de contêineres do EMR em. GitHub

- 1. Instale ou atualize para a versão mais recente do AWS CLI
- 2. Configuração do kubectl e eksctl
- 3. Conceitos básicos do Amazon EKS: eksctl
- 4. Habilite o acesso ao cluster para o Amazon EMR no EKS
- 5. Ativar funções do IAM para o cluster EKS
- 6. Concessão de acesso ao Amazon EMR no EKS para os usuários
- 7. Registro do cluster do Amazon EKS com o Amazon EMR

StartJobRun 93

Habilitação do acesso ao cluster para o Amazon EMR no EKS

As seções a seguir mostram algumas maneiras de habilitar o acesso ao cluster. A primeira é usando o gerenciamento de acesso ao cluster (CAM) do Amazon EKS e a segunda mostra como executar etapas manuais para habilitar o acesso ao cluster.

Como habilitar o acesso ao cluster usando a entrada de acesso do EKS (recomendado)



Note

O aws-auth ConfigMap está obsoleto. O método recomendado para gerenciar o acesso ao Kubernetes APIs é o Access Entries.

O Amazon EMR é integrado ao gerenciamento de acesso ao cluster (CAM) do Amazon EKS, para que você possa automatizar a configuração das políticas AuthN e AuthZ necessárias para executar trabalhos Spark do Amazon EMR em namespaces de clusters do Amazon EKS. Quando você cria um cluster virtual usando um namespace de cluster do Amazon EKS, o Amazon EMR configura automaticamente todas as permissões necessárias, para que você não precise adicionar nenhuma etapa extra aos fluxos de trabalho atuais.



Note

A integração do Amazon EMR com o CAM do Amazon EKS é compatível somente com novos clusters virtuais do Amazon EMR no EKS. Não é possível migrar clusters virtuais existentes para usar essa integração.

Pré-requisitos

- Verifique se você está executando a versão 2.15.3 ou superior do AWS CLI
- O cluster do Amazon EKS deve estar na versão 1.23 ou posterior.

Configuração

Para configurar a integração entre o Amazon EMR e as operações de AccessEntry API do Amazon EKS, certifique-se de ter concluído os seguintes itens:

 Verifique se authenticationMode do cluster do Amazon EKS está configurado como API_AND_CONFIG_MAP.

```
aws eks describe-cluster --name <eks-cluster-name>
```

Se ainda não estiver, defina authenticationMode como API_AND_CONFIG_MAP.

```
aws eks update-cluster-config
    --name <eks-cluster-name>
    --access-config authenticationMode=API_AND_CONFIG_MAP
```

Para obter mais informações sobre modos de autenticação, consulte Cluster authentication modes.

Certifique-se de que o <u>perfil do IAM</u> que você está usando para executar
 CreateVirtualCluster e as operações da API DeleteVirtualCluster também tenha as seguintes permissões:

```
"Effect": "Allow",
  "Action": [
    "eks:CreateAccessEntry"
  ],
  "Resource":
 "arn:<aws_PARTITION>:eks:<aws_REGION>:<aws_ACCOUNT_ID>:cluster/<EKS_CLUSTER_NAME>"
},
  "Effect": "Allow",
  "Action": [
    "eks:DescribeAccessEntry",
    "eks:DeleteAccessEntry",
    "eks:ListAssociatedAccessPolicies",
    "eks:AssociateAccessPolicy",
    "eks:DisassociateAccessPolicy"
  ],
  "Resource": "arn:<a href="AWS_PARTITION">:eks:<AWS_REGION>:<AWS_ACCOUNT_ID>:access-entry/"
<EKS_CLUSTER_NAME>/role/<AWS_ACCOUNT_ID>/AWSServiceRoleForAmazonEMRContainers/*"
}
```

Conceitos e terminologia

A seguir está uma lista de terminologias e conceitos relacionados ao CAM do Amazon EKS.

- Cluster virtual (VC): representação lógica do namespace criado no Amazon EKS. É um link 1:1 para um namespace de cluster do Amazon EKS. Você pode usá-lo para executar workloads do Amazon EMR em um cluster do Amazon EKS dentro do namespace especificado.
- Namespace: mecanismo para isolar grupos de recursos em um único cluster do EKS.
- Política de acesso: permissões que concedem acesso e ações a um perfil do IAM em um cluster do EKS.
- Entrada de acesso: uma entrada criada com um ARN de perfil. Você pode vincular a entrada de acesso a uma política de acesso para atribuir permissões específicas no cluster do Amazon EKS.
- Cluster virtual integrado de entrada de acesso do EKS: o cluster virtual criado usando operações de API de entrada de acesso do Amazon EKS.

Como habilitar o acesso ao cluster usando aws-auth

Você deve permitir o acesso do Amazon EMR no EKS a um namespace específico em seu cluster ao executar as seguintes ações: criar um perfil do Kubernetes, associar o perfil a um usuário do Kubernetes e mapear o usuário do Kubernetes com o perfil vinculado ao serviço AWSServiceRoleForAmazonEMRContainers. Essas ações são automatizadas em eksctl quando o comando de mapeamento de identidade do IAM é usado com emr-containers como nome do serviço. Você pode executar essas operações facilmente ao usar o comando apresentado a seguir.

```
eksctl create iamidentitymapping \
    --cluster my_eks_cluster \
    --namespace kubernetes_namespace \
    --service-name "emr-containers"
```

my eks clusterSubstitua pelo nome do seu cluster Amazon EKS e kubernetes namespace substitua pelo namespace Kubernetes criado para executar cargas de trabalho do Amazon EMR.

Important

Você deve fazer o download do eksctl mais recente usando a etapa anterior Configuração do kubectl e eksctl para usar esta funcionalidade.

Etapas manuais para habilitar o acesso ao cluster para o Amazon EMR no EKS

Você também pode usar as etapas manuais apresentadas a seguir para habilitar o acesso ao cluster para o Amazon EMR no EKS.

1. Crie um perfil do Kubernetes em um namespace específico

Amazon EKS 1.22 - 1.29

Com o Amazon EKS, versões 1.22 a 1.29, execute o comando a seguir para criar um perfil do Kubernetes em um namespace específico. Esse perfil concede as permissões de RBAC obrigatórias ao Amazon EMR no EKS.

```
namespace=my-namespace
cat - >>EOF | kubectl apply -f - >>namespace "${namespace}"
apiVersion: rbac.authorization.k8s.io/v1
kind: Role
metadata:
  name: emr-containers
  namespace: ${namespace}
rules:
  - apiGroups: [""]
    resources: ["namespaces"]
    verbs: ["get"]
  - apiGroups: [""]
    resources: ["serviceaccounts", "services", "configmaps", "events", "pods",
 "pods/log"]
    verbs: ["get", "list", "watch", "describe", "create", "edit", "delete",
 "deletecollection", "annotate", "patch", "label"]
  - apiGroups: [""]
    resources: ["secrets"]
    verbs: ["create", "patch", "delete", "watch"]
  - apiGroups: ["apps"]
    resources: ["statefulsets", "deployments"]
    verbs: ["get", "list", "watch", "describe", "create", "edit", "delete",
 "annotate", "patch", "label"]
  - apiGroups: ["batch"]
    resources: ["jobs"]
    verbs: ["get", "list", "watch", "describe", "create", "edit", "delete",
 "annotate", "patch", "label"]
  - apiGroups: ["extensions", "networking.k8s.io"]
    resources: ["ingresses"]
```

```
verbs: ["get", "list", "watch", "describe", "create", "edit", "delete",
"annotate", "patch", "label"]
  - apiGroups: ["rbac.authorization.k8s.io"]
    resources: ["roles", "rolebindings"]
    verbs: ["get", "list", "watch", "describe", "create", "edit", "delete",
"deletecollection", "annotate", "patch", "label"]
    - apiGroups: [""]
    resources: ["persistentvolumeclaims"]
    verbs: ["get", "list", "watch", "describe", "create", "edit", "delete",
"deletecollection", "annotate", "patch", "label"]
EOF
```

Amazon EKS 1.21 and below

Com o Amazon EKS, versões 1.21 e anteriores, execute o comando a seguir para criar um perfil do Kubernetes em um namespace específico. Esse perfil concede as permissões de RBAC obrigatórias ao Amazon EMR no EKS.

```
namespace=my-namespace
cat - >>EOF | kubectl apply -f - >>namespace "${namespace}"
apiVersion: rbac.authorization.k8s.io/v1
kind: Role
metadata:
  name: emr-containers
  namespace: ${namespace}
rules:
  - apiGroups: [""]
    resources: ["namespaces"]
    verbs: ["get"]
  - apiGroups: [""]
    resources: ["serviceaccounts", "services", "configmaps", "events", "pods",
 "pods/log"]
    verbs: ["get", "list", "watch", "describe", "create", "edit", "delete",
 "deletecollection", "annotate", "patch", "label"]
  - apiGroups: [""]
    resources: ["secrets"]
    verbs: ["create", "patch", "delete", "watch"]
  - apiGroups: ["apps"]
    resources: ["statefulsets", "deployments"]
    verbs: ["get", "list", "watch", "describe", "create", "edit", "delete",
 "annotate", "patch", "label"]
  - apiGroups: ["batch"]
```

```
resources: ["jobs"]
   verbs: ["get", "list", "watch", "describe", "create", "edit", "delete",
 "annotate", "patch", "label"]
  - apiGroups: ["extensions"]
   resources: ["ingresses"]
   verbs: ["get", "list", "watch", "describe", "create", "edit", "delete",
 "annotate", "patch", "label"]
  - apiGroups: ["rbac.authorization.k8s.io"]
   resources: ["roles", "rolebindings"]
   verbs: ["get", "list", "watch", "describe", "create", "edit", "delete",
 "deletecollection", "annotate", "patch", "label"]
 - apiGroups: [""]
   resources: ["persistentvolumeclaims"]
   verbs: ["get", "list", "watch", "describe", "create", "edit", "delete",
 "deletecollection", "annotate", "patch", "label"]
EOF
```

2. Crie uma associação de perfil do Kubernetes com escopo definido para o namespace

Execute o comando apresentado a seguir para criar uma associação de perfil do Kubernetes no namespace especificado. Essa associação de perfil concede as permissões definidas no perfil criado na etapa anterior a um usuário chamado emr-containers. Esse usuário identifica perfis vinculados ao serviço para o Amazon EMR no EKS e, dessa forma, permite que o Amazon EMR no EKS execute ações conforme definido pelo perfil criado.

```
namespace=my-namespace

cat - <<EOF | kubectl apply -f - --namespace "${namespace}"
apiVersion: rbac.authorization.k8s.io/v1
kind: RoleBinding
metadata:
   name: emr-containers
   namespace: ${namespace}
subjects:
   - kind: User
   name: emr-containers
   apiGroup: rbac.authorization.k8s.io
roleRef:
   kind: Role
   name: emr-containers
   apiGroup: rbac.authorization.k8s.io</pre>
```

EOF

Atualize mapa de configuração aws-auth do Kubernetes

Você pode usar uma das opções apresentadas a seguir para mapear o perfil vinculado ao serviço do Amazon EMR no EKS com o usuário emr-containers que foi associado ao perfil do Kubernetes na etapa anterior.

Opção 1: usar o eksctl

Execute o comando eksctl apresentado a seguir para mapear o perfil vinculado ao serviço do Amazon EMR no EKS com o usuário emr-containers.

```
eksctl create iamidentitymapping \
    --cluster my-cluster-name \
    --arn "arn:aws:iam::my-account-id:role/AWSServiceRoleForAmazonEMRContainers" \
    --username emr-containers
```

Opção 2: sem usar o eksctl

1. Execute o comando apresentado a seguir para abrir o mapa de configuração aws-auth no editor de texto.

```
kubectl edit -n kube-system configmap/aws-auth
```



Note

Se você receber um erro informandoError from server (NotFound): configmaps "aws-auth" not found, consulte as etapas em Adicionar funções de usuário no Guia do usuário do Amazon EKS para aplicar o estoque ConfigMap.

2. Adicione detalhes do perfil vinculado ao serviço do Amazon EMR no EKS à seção mapRoles do ConfigMap, em data. Adicione essa seção se ela ainda não existir no arquivo. A seção mapRoles atualizada em dados deve ser semelhante ao exemplo a seguir.

```
apiVersion: v1
data:
 mapRoles: |
    - rolearn: arn:aws:iam::<your-account-id>:role/
AWSServiceRoleForAmazonEMRContainers
```

```
username: emr-containers
- ... <other previously existing role entries, if there's any>.
```

3. Salve o arquivo e saia do seu editor de texto.

Ativar funções do IAM para o cluster EKS

Os tópicos a seguir detalham as opções para habilitar as funções do IAM.

Tópicos

- Opção 1: ativar o EKS Pod Identity no EKS Cluster
- Opção 2: habilitar funções do IAM para contas de serviço (IRSA) no cluster EKS

Opção 1: ativar o EKS Pod Identity no EKS Cluster

As associações do Amazon EKS Pod Identity oferecem a capacidade de gerenciar credenciais para seus aplicativos, da mesma forma que os perfis de EC2 instância da Amazon fornecem credenciais às instâncias da Amazon EC2. O Amazon EKS Pod Identity fornece credenciais para suas workloads com uma API EKS Auth adicional e um pod de agente executado em cada nó.

O Amazon EMR no EKS começa a oferecer suporte à identidade de pod do EKS desde o lançamento do emr-7.3.0 para o modelo de envio. StartJobRun

Para obter mais informações sobre identidades de pod do EKS, consulte Entenda como funciona o EKS Pod Identity.

Por que usar o EKS Pod Identities?

Como parte da configuração do EMR, o Job Execution Role precisa estabelecer limites de confiança entre uma função do IAM e as contas de serviço em um namespace específico (de clusters virtuais do EMR). Com o IRSA, isso foi conseguido atualizando a política de confiança do EMR Job Execution Role. No entanto, devido ao limite rígido de 4096 caracteres no tamanho da política de confiança do IAM, havia uma restrição para compartilhar uma única função do IAM de execução de trabalhos em um máximo de doze (12) clusters EKS.

Com o suporte do EMR para identidades de pods, o limite de confiança entre as funções do IAM e as contas de serviço agora está sendo gerenciado pela equipe do EKS por meio da associação da identidade do pod do EKS. APIs



Note

O limite de segurança da identidade do pod EKS ainda está no nível da conta de serviço, não no nível do pod.

Considerações sobre a identidade do pod

Para obter informações sobre as limitações do Pod Identity, consulte Considerações sobre o EKS Pod Identity.

Prepare a identidade do EKS Pod no EKS Cluster

Verifique se a permissão necessária existe em NodeInstanceRole

A função de nó NodeInstanceRole precisa de uma permissão para que o agente realize a AssumeRoleForPodIdentity ação na API EKS Auth. Você pode adicionar o seguinte à Amazon EKSWorker NodePolicy, que é definida no Guia do usuário do Amazon EKS, ou usar uma política personalizada.

Se seu cluster EKS foi criado com a versão eksctl superior a 0.181.0, a Amazon, incluindo a AssumeRoleForPodIdentity permissão necessária EKSWorkerNodePolicy, será anexada automaticamente à função de nó. Se a permissão não estiver presente, adicione manualmente a seguinte permissão à Amazon, EKSWorker NodePolicy que permite assumir uma função para a identidade do pod. Essa permissão é necessária pelo agente de identidade de pods do EKS para recuperar as credenciais dos pods.

JSON

```
"Version": "2012-10-17",
"Statement": [
  {
    "Effect": "Allow",
    "Action": [
      "eks-auth:AssumeRoleForPodIdentity"
    ],
    "Resource": [
      ·· * ··
    ],
```

```
"Sid": "AllowEKSAUTHAssumeroleforpodidentity"
     }
]
```

Crie o complemento do agente de identidade EKS pod

Use o comando a seguir para criar o complemento EKS Pod Identity Agent com a versão mais recente:

```
aws eks create-addon --cluster-name cluster-name --addon-name eks-pod-identity-agent kubectl get pods -n kube-system | grep 'eks-pod-identity-agent'
```

Use as etapas a seguir para criar o complemento EKS Pod Identity Agent a partir do console Amazon EKS:

- 1. Abra o console do Amazon EKS: console do Amazon EKS.
- No painel de navegação esquerdo, selecione Clusters e depois o nome do cluster para o qual você deseja configurar o complemento do EKS Pod Identity Agent.
- 3. Escolha a guia Add-ons (Complementos).
- 4. Escolha Obter mais complementos.
- 5. Selecione a caixa no canto superior direito da caixa do complemento do EKS Pod Identity Agent e escolha Editar.
- 6. Na página Definir configurações de complementos selecionados, selecione qualquer versão na lista suspensa Versão.
- 7. (Opcional) Expanda Configurações opcionais para inserir configurações adicionais. Por exemplo, é possível indicar um local alternativo para a imagem do contêiner e ImagePullSecrets. O esquema JSON com chaves aceitas é mostrado no Esquema de configuração do complemento.

Insira as chaves e os valores de configuração em Valores de configuração.

- 8. Escolha Próximo.
- 9. Confirme se os pods do agente estão sendo executados no seu cluster por meio da CLI.

kubectl get pods -n kube-system | grep 'eks-pod-identity-agent'

Um exemplo de saída é o seguinte:

	NAME	READY	STATUS	RESTARTS	AGE
	eks-pod-identity-agent-gmqp7	1/1	Running	1 (24h ago)	24h
	eks-pod-identity-agent-prnsh	1/1	Running	1 (24h ago)	24h
- (

Isso configura um novo DaemonSet no kube-system namespace. O Amazon EKS Pod Identity Agent, executado em cada nó do EKS, usa a AssumeRoleForPodIdentityação para recuperar credenciais temporárias da API EKS Auth. Essas credenciais são então disponibilizadas para AWS SDKs que você execute dentro de seus contêineres.

Para obter mais informações, verifique o pré-requisito no documento público: Configurar o Amazon EKS Pod Identity Agent.

Criar uma função de execução de trabalhos

Crie ou atualize a função de execução de tarefas que permite o EKS Pod Identity

Para executar cargas de trabalho com o Amazon EMR no EKS, você precisa criar uma função do IAM. Referimo-nos a esse perfil como perfil de execução de trabalho nesta documentação. Para obter mais informações sobre como criar a função do IAM, consulte Como criar funções do IAM no Guia do usuário.

Além disso, você deve criar uma política do IAM que especifique as permissões necessárias para a função de execução do trabalho e, em seguida, anexe essa política à função para habilitar o EKS Pod Identity.

Por exemplo, você tem a seguinte função de execução de tarefas. Para obter mais informações, consulte Criar uma função de execução de tarefas.

arn:aws:iam::111122223333:role/PodIdentityJobExecutionRole

Important

O Amazon EMR no EKS cria automaticamente contas de serviço do Kubernetes, com base no nome da função de execução do trabalho. Certifique-se de que o nome da função não seja muito longo, pois seu trabalho pode falhar se a combinação de cluster_namepod_name, e service_account_name exceder o limite de duração.

Configuração da função de execução do trabalho — Certifique-se de que a função de execução do trabalho seja criada com a permissão de confiança abaixo para o EKS Pod Identity. Para atualizar uma função de execução de trabalho existente, configure-a para confiar no seguinte diretor de serviço do EKS como uma permissão adicional na política de confiança. Essa permissão de confiança pode coexistir com as políticas de confiança existentes da IRSA.

```
cat >trust-relationship.json <<EOF</pre>
{
    "Version": "2012-10-17",
    "Statement": [
        {
             "Sid": "AllowEksAuthToAssumeRoleForPodIdentity",
             "Effect": "Allow",
             "Principal": {
                 "Service": "pods.eks.amazonaws.com"
            },
             "Action": [
                 "sts:AssumeRole",
                 "sts:TagSession"
            ]
        }
    ]
}
EOF
```

Permissão do usuário: os usuários precisam da iam: PassRole permissão para executar chamadas de StartJobRun API ou enviar trabalhos. Essa permissão permite que os usuários passem a função de execução do trabalho para o EMR no EKS. Os administradores de trabalhos devem ter a permissão por padrão.

Abaixo está a permissão necessária para um usuário:

```
"Effect": "Allow",
"Action": "iam:PassRole",
"Resource": "arn:aws:iam::111122223333:role/PodIdentityJobExecutionRole",
"Condition": {
    "StringEquals": {
        "iam:PassedToService": "pods.eks.amazonaws.com"
    }
}
```

}

Para restringir ainda mais o acesso do usuário a clusters EKS específicos, adicione o filtro de AssociatedResourceArn atributos à política do IAM. Isso limita a suposição de funções aos clusters EKS autorizados, fortalecendo seus controles de segurança em nível de recursos.

Configurar associações de identidade de pod EKS

Pré-requisito

Certifique-se de que a identidade do IAM que cria a associação de identidade do pod, como um usuário administrador do EKS, tenha a permissão eks:CreatePodIdentityAssociation iam:PassRole e.

JSON

```
"Version": "2012-10-17",
"Statement": [
 {
    "Effect": "Allow",
    "Action": [
      "eks:CreatePodIdentityAssociation"
   ],
    "Resource": [
      "arn:aws:eks:*:*:cluster/*"
   ],
    "Sid": "AllowEKSCreatepodidentityassociation"
 },
    "Effect": "Allow",
    "Action": [
      "iam:PassRole"
    ],
```

Crie associações para a função e a conta de serviço do EMR

Create EMR role associations through the AWS CLI

Quando você envia um trabalho para um namespace do Kubernetes, um administrador deve criar associações entre a função de execução do trabalho e a identidade da conta de serviço gerenciado do EMR. Observe que a conta de serviço gerenciado do EMR é criada automaticamente no envio do trabalho, com escopo definido para o namespace no qual o trabalho é enviado.

Com a versão 2.24.0 AWS CLI (acima), execute o comando a seguir para criar associações de função com a identidade do pod.

Execute o comando a seguir para criar associações de função com a identidade do pod:

Nota:

- Cada cluster pode ter um limite de 1.000 associações. Cada função de execução do trabalho
 o mapeamento de namespace exigirá 3 associações para os pods de remetente, driver e executor do trabalho.
- Você só pode associar funções que estejam na mesma AWS conta do cluster. É possível delegar o acesso de outra conta ao perfil nessa conta configurada por você para o EKS Pod

Identities usar. Para ver um tutorial sobre delegação de acessoAssumeRole, consulte o <u>tutorial</u> do IAM: Delegar acesso entre AWS contas usando funções do IAM.

Create EMR role associations through Amazon EKS

O EMR cria uma conta de serviço com determinado padrão de nomenclatura quando um trabalho é enviado. Para fazer associações manuais ou integrar esse fluxo de trabalho ao AWS SDK, siga estas etapas:

Nome da conta de serviço do Construct:

```
emr-containers-sa-spark-%(SPARK_ROLE)s-%(AWS_ACCOUNT_ID)s-
%(BASE36_ENCODED_ROLE_NAME)s
```

Os exemplos abaixo criam uma associação de função para um exemplo de função de execução de Job JobExecutionRoleIRSAv2.

Exemplos de associações de funções:

```
RoleName: JobExecutionRoleIRSAv2
Base36EncodingOfRoleName: 2eum5fah1jc1kwyjc19ikdhdkdegh1n26vbe
```

Exemplo de comando CLI:

```
# setup for the client service account (used by job runner pod)
# emr-containers-sa-spark-client-111122223333-2eum5fah1jc1kwyjc19ikdhdkdegh1n26vbe
aws eks create-pod-identity-association --cluster-name mycluster
    --role-arn arn:aws:iam::111122223333:role/JobExecutionRoleIRSAv2
    --namespace mynamespace --service-account emr-containers-sa-spark-
client-111122223333-2eum5fah1jc1kwyjc19ikdhdkdegh1n26vbe

# driver service account
# emr-containers-sa-spark-driver-111122223333-2eum5fah1jc1kwyjc19ikdhdkdegh1n26vbe

aws eks create-pod-identity-association --cluster-name mycluster
    --role-arn arn:aws:iam::111122223333:role/JobExecutionRoleIRSAv2
    --namespace mynamespace --service-account emr-containers-sa-spark-
driver-111122223333-2eum5fah1jc1kwyjc19ikdhdkdegh1n26vbe

# executor service account
```

```
# emr-containers-sa-spark-executor-111122223333-2eum5fah1jc1kwyjc19ikdhdkdegh1n26vbe
aws eks create-pod-identity-association --cluster-name mycluster
   --role-arn arn:aws:iam::111122223333:role/JobExecutionRoleIRSAv2
   --namespace mynamespace --service-account emr-containers-sa-spark-
executor-111122223333-2eum5fah1jc1kwyjc19ikdhdkdegh1n26vbe
```

Depois de concluir todas as etapas necessárias para a identidade do pod EKS, você pode pular as seguintes etapas para a configuração do IRSA:

- Ativar funções do IAM para contas de serviço (IRSA) no cluster EKS
- Crie uma função de execução de tarefas
- Atualize a política de confiança da função de execução do trabalho

Você pode pular diretamente para a seguinte etapa: <u>Conceder aos usuários acesso ao Amazon EMR</u> no EKS

Excluir associações de função

Sempre que você excluir um cluster virtual ou uma função de execução de tarefas e não quiser mais dar acesso ao EMR às suas contas de serviço, exclua as associações da função. Isso ocorre porque o EKS permite associações com recursos inexistentes (namespace e conta de serviço). O Amazon EMR no EKS recomenda excluir as associações se o namespace for excluído ou se a função não estiver mais em uso, para liberar espaço para outras associações.



As associações persistentes podem potencialmente afetar sua capacidade de escalar se você não as excluir, pois o EKS tem limitações no número de associações que você pode criar (limite flexível: 1000 associações por cluster). Você pode listar as associações de identidade do pod em um determinado namespace para verificar se há alguma associação remanescente que precise ser limpa:

aws eks list-pod-identity-associations --cluster-name mycluster --namespace mynamespace

Com a AWS CLI (versão 2.24.0 ou superior), execute o seguinte comando emr-containers para excluir as associações de função do EMR:

Migre automaticamente o IRSA existente para o Pod Identity

Você pode usar a ferramenta eksctl para migrar as funções existentes do IAM para contas de serviço (IRSA) para associações de identidade de pod:

```
eksctl utils migrate-to-pod-identity \
    --cluster mycluster \
    --remove-oidc-provider-trust-relationship \
    --approve
```

A execução do comando sem o --approve sinalizador produzirá apenas um plano que reflete as etapas da migração, e nenhuma migração real ocorrerá.

Solução de problemas

Meu trabalho falhou com NoClassDefinitionFound ou com ClassNotFound exceção do provedor de credenciais, ou falhou em obter o provedor de credenciais.

O EKS Pod Identity usa o Container Credentials Provider para recuperar as credenciais necessárias. Se você especificou um provedor de credenciais personalizado, verifique se ele está funcionando corretamente. Como alternativa, verifique se você está usando uma versão correta do AWS SDK compatível com o EKS Pod Identity. Para obter mais informações, consulte Comece a usar o Amazon EKS.

O trabalho falhou com o erro "Falha na recuperação de credenciais devido ao limite de tamanho [x]" mostrado no eks-pod-identity-agent registro.

O EMR no EKS cria contas de serviço do Kubernetes com base no nome da função de execução do trabalho. Se o nome da função for muito longo, o EKS Auth não conseguirá recuperar as credenciais porque a combinação decluster_name,pod_name, e service_account_name excede o limite de comprimento. Identifique qual componente está ocupando mais espaço e ajuste o tamanho de acordo.

O trabalho falhou com o erro "Failed to Retrieve Credentials xxx" exibido no eks-pod-identity registro.

Uma possível causa desse problema pode ser que o cluster EKS esteja configurado em sub-redes privadas sem ser configurado corretamente PrivateLink para o cluster. Verifique se seu cluster está em uma rede privada e configure AWS PrivateLink para resolver o problema. Para obter instruções detalhadas, consulte Comece a usar o Amazon EKS.

Opção 2: habilitar funções do IAM para contas de serviço (IRSA) no cluster EKS

O recurso de perfis do IAM para contas de serviço está disponível nas versões 1.14 e posteriores do Amazon EKS e para clusters do EKS atualizados para as versões 1.13 ou posteriores ou após 3 de setembro de 2019. Para usar esse recurso, é possível atualizar os clusters do EKS existentes para a versão 1.14 ou posterior. Para obter mais informações, consulte Atualizar uma versão do Kubernetes do cluster do Amazon EKS.

Se o seu cluster oferecer suporte para perfis do IAM para contas de serviço, ele terá um URL do emissor OpenID Connect associado a ele. Você pode visualizar essa URL no console do Amazon EKS ou usar o AWS CLI comando a seguir para recuperá-la.



Important

Você deve usar a versão mais recente do AWS CLI para receber a saída adequada desse comando.

aws eks describe-cluster --name cluster_name --query "cluster.identity.oidc.issuer" -output text

A saída esperada é semelhante à apresentada a seguir.

```
https://oidc.eks.<region-code>.amazonaws.com/id/EXAMPLED539D4633E53DE1B716D3041E
```

Para usar perfis do IAM para contas de serviço em seu cluster, é necessário criar um provedor de identidades OIDC usando o eksctl ou o AWS Management Console.

Para criar um provedor de identidade OIDC do IAM para o cluster com o eksct1

Verifique a versão do eksctl com o comando a seguir. Este procedimento pressupõe que você instalou o eksctl e que a versão do eksctl seja 0.32.0 ou posterior.

```
eksctl version
```

Para obter mais informações sobre como instalar ou atualizar o eksctl, consulte <u>Instalar ou atualizar o</u> eksctl.

Crie o provedor de identidade OIDC para o cluster com o seguinte comando. Substitua *cluster_name* pelos seus próprios valores.

```
eksctl utils associate-iam-oidc-provider --cluster cluster_name --approve
```

Para criar um provedor de identidade IAM OIDC para seu cluster com o AWS Management Console

Recupere a URL do emissor do OIDC na descrição do console Amazon EKS do seu cluster ou use o comando a seguir. AWS CLI

Use o comando apresentado a seguir para recuperar o URL do emissor OIDC da AWS CLI.

```
aws eks describe-cluster --name <cluster_name> --query "cluster.identity.oidc.issuer"
    --output text
```

Use as etapas apresentadas a seguir para recuperar o URL do emissor OIDC do console do Amazon EKS.

- 1. Abra o console do IAM em https://console.aws.amazon.com/iam/.
- 2. No painel de navegação, escolha Provedores de identidade e, em seguida, selecione Criar provedor.
 - 1. Para Tipo de provedor, escolha Escolher um tipo de provedor e escolha OpenID Connect.
 - 2. Em Provider URL (URL do provedor), cole o URL emissor OIDC do cluster.
 - 3. Para Público, digite sts.amazonaws.com e escolha Próxima etapa.
- 3. Verifique se as informações do provedor estão corretas e escolha Create (Criar) para criar seu provedor de identidade.

Criação de um perfil de execução de trabalho

Para executar workloads no Amazon EMR no EKS, você precisa criar um perfil do IAM. Referimo-nos a esse perfil como perfil de execução de trabalho nesta documentação. Para obter mais informações sobre como criar perfis do IAM, consulte Criação de perfis do IAM no Guia do usuário do IAM.

Você também deve criar uma política do IAM que especifique as permissões para o perfil de execução de trabalho e, em seguida, anexar a política do IAM ao perfil de execução de trabalho.

A política a seguir para a função de execução do trabalho permite acesso a destinos de recursos, Amazon S3 e. CloudWatch Essas permissões são necessárias para monitorar trabalhos e acessar logs. Para seguir o mesmo processo usando o AWS CLI:

Crie a função do IAM para execução do trabalho: vamos criar o papel que o EMR usará para a execução do trabalho. Essa é a função que os trabalhos do EMR assumirão quando forem executados no EKS.

Em seguida, precisamos anexar as políticas de IAM necessárias à função para que ela possa gravar registros no s3 e no cloudwatch.

```
"s3:ListBucket"
             ],
             "Resource": "arn:aws:s3:::amzn-s3-demo-bucket"
         },
         {
             "Effect": "Allow",
             "Action": [
                 "logs:PutLogEvents",
                 "logs:CreateLogStream",
               "logs:DescribeLogGroups",
                 "logs:DescribeLogStreams"
             ],
             "Resource": [
                 "arn:aws:logs:*:*:*"
             ]
         }
     ]
}
EoF
aws iam put-role-policy --role-name EMRContainers-JobExecutionRole --policy-name
EMR-Containers-Job-Execution --policy-document file://~/environment/EMRContainers-
JobExecutionRole.json
```

Note

O acesso deve ter um escopo adequado e não ser concedido a todos os objetos do S3 no perfil de execução de trabalho.

JSON

```
"Resource": [
        "arn:aws:s3:::amzn-s3-demo-bucket"
      ],
      "Sid": "AllowS3Putobject"
    },
      "Effect": "Allow",
      "Action": [
        "logs:PutLogEvents",
        "logs:CreateLogStream",
        "logs:DescribeLogGroups",
        "logs:DescribeLogStreams"
      ],
      "Resource": [
        "arn:aws:logs:*:*:*"
      ],
      "Sid": "AllowLOGSPutlogevents"
    }
  ]
}
```

Para obter mais informações, consulte <u>Usando funções de execução de trabalho</u>, <u>Configurar uma execução de trabalho para usar registros do S3</u> e <u>Configurar uma execução de trabalho para usar CloudWatch registros</u>.

Atualização da política de confiança do perfil de execução de trabalho

Ao usar perfis do IAM para contas de serviço (IRSA) com a finalidade de executar trabalhos em um namespace do Kubernetes, um administrador deve criar uma relação de confiança entre o perfil de execução de trabalho e a identidade da conta de serviço gerenciado do EMR. A relação de confiança pode ser criada ao atualizar a política de confiança do perfil de execução de trabalho. Observe que a conta de serviço gerenciado do EMR é criada automaticamente no envio do trabalho, com escopo definido para o namespace no qual o trabalho é enviado.

Execute o comando apresentado a seguir para atualizar a política de confiança.

```
aws emr-containers update-role-trust-policy \
     --cluster-name cluster \
     --namespace namespace \
     --role-name iam_role_name_for_job_execution
```

Para obter mais informações, consulte Uso de perfis de execução de trabalho com o Amazon EMR no EKS.

♠ Important

O operador que executa o comando acima deve ter estas permissões: eks:DescribeCluster, iam:GetRole e iam:UpdateAssumeRolePolicy.

Concessão de acesso ao Amazon EMR no EKS para os usuários

Para as ações executadas no Amazon EMR no EKS, uma permissão do IAM correspondente para essa ação é necessária. Você deve criar uma política do IAM que permita executar ações do Amazon EMR no EKS e anexá-la ao perfil ou ao usuário do IAM que você usa.

Este tópico fornece etapas para a criação de uma nova política e para o anexo dela a um usuário. Ele também abrange as permissões básicas necessárias para configurar o ambiente do Amazon EMR no EKS. Recomendamos redefinir as permissões para recursos específicos sempre que possível com base nas suas necessidades de negócios.

Criação de uma nova política do IAM e anexo dela a um usuário no console do IAM

Crie uma nova política do IAM

- Faça login no AWS Management Console e abra o console do IAM em https:// 1. console.aws.amazon.com/iam/.
- 2. No painel de navegação esquerdo do console do IAM, escolha Políticas.
- Na página Policies (Políticas), escolha Create Policy (Criar política). 3.
- Na janela Criar política, navegue até a guia Editar JSON. Crie um documento de política 4. com uma ou mais instruções JSON, conforme mostrado nos exemplos que seguem este procedimento. Em seguida, escolha Analisar política.
- Na tela Review Policy (Revisar política), insira o Policy Name (Nome da política), por exemplo, AmazonEMROnEKSPolicy. Insira uma descrição opcional e, em seguida, escolha Criar política.

Anexe a política a um usuário ou a um perfil

Faça login no AWS Management Console e abra o console do IAM em https:// 1. console.aws.amazon.com/iam/

- 2. No painel de navegação, selecione Políticas.
- 3. Na lista de políticas, marque a caixa de seleção ao lado da política criada na seção anterior. Você pode usar o menu Filtro e a caixa de pesquisa para filtrar a lista de políticas.
- 4. Selecione Ações da política e escolha Anexar.
- 5. Escolha o usuário ou o perfil ao qual a política será anexada. Você pode usar o menu Filter (Filtro) e a caixa de pesquisa para filtrar a lista de entidades principais. Após escolher o usuário ou o perfil ao qual a política será anexada, selecione Anexar política.

Permissões para o gerenciamento de clusters virtuais

Para gerenciar clusters virtuais em sua AWS conta, crie uma política do IAM com as seguintes permissões. Essas permissões permitem que você crie, liste, descreva e exclua clusters virtuais em sua AWS conta.

JSON

```
"Version": "2012-10-17",
"Statement": [
  {
    "Effect": "Allow",
    "Action": [
      "iam:CreateServiceLinkedRole"
    ],
    "Resource": [
      11 * 11
   ],
    "Condition": {
      "StringLike": {
        "iam:AWSServiceName": "emr-containers.amazonaws.com"
      }
    },
    "Sid": "AllowIAMCreateservicelinkedrole"
 },
 {
    "Effect": "Allow",
    "Action": [
      "emr-containers:CreateVirtualCluster",
      "emr-containers:ListVirtualClusters",
      "emr-containers:DescribeVirtualCluster",
```

```
"emr-containers:DeleteVirtualCluster"
],
    "Resource": [
        "*"
      ],
      "Sid": "AllowEMRCONTAINERSCreatevirtualcluster"
    }
]
```

O Amazon EMR é integrado ao gerenciamento de acesso ao cluster (CAM) do Amazon EKS, para que você possa automatizar a configuração das políticas AuthN e AuthZ necessárias para executar trabalhos Spark do Amazon EMR em namespaces de clusters do Amazon EKS. Para fazer isso, você deve ter as seguintes permissões:

```
"Effect": "Allow",
  "Action": [
    "eks:CreateAccessEntry"
  ],
  "Resource":
 "arn:<AWS_PARTITION>:eks:<AWS_REGION>:<AWS_ACCOUNT_ID>:cluster/<EKS_CLUSTER_NAME>"
},
{
  "Effect": "Allow",
  "Action": [
    "eks:DescribeAccessEntry",
    "eks:DeleteAccessEntry",
    "eks:ListAssociatedAccessPolicies",
    "eks:AssociateAccessPolicy",
    "eks:DisassociateAccessPolicy"
  ],
  "Resource": "arn:<ahvs_partition>:eks:<ahvs_region>:<ahvs_account_id>:access-
entry/<EKS_CLUSTER_NAME>/role/<AWS_ACCOUNT_ID>/AWSServiceRoleForAmazonEMRContainers/*"
}
```

Para obter mais informações, consulte Automate enabling cluster access for Amazon EMR on EKS.

Quando a CreateVirtualCluster operação é invocada pela primeira vez a partir de uma AWS conta, você também precisa das CreateServiceLinkedRole permissões para criar a função

vinculada ao serviço para o Amazon EMR no EKS. Para obter mais informações, consulte <u>Uso de</u> perfis vinculados ao serviço para o Amazon EMR no EKS.

Permissões para o envio de trabalhos

Para enviar trabalhos nos clusters virtuais da sua AWS conta, crie uma política do IAM com as seguintes permissões. Essas permissões permitem iniciar, listar, descrever e cancelar execuções de trabalhos para todos os clusters virtuais em sua conta. Você deve considerar adicionar permissões para listar ou descrever clusters virtuais, o que permite verificar o estado do cluster virtual antes de enviar trabalhos.

JSON

```
"Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Action": [
        "emr-containers:StartJobRun",
        "emr-containers:ListJobRuns",
        "emr-containers:DescribeJobRun",
        "emr-containers:CancelJobRun"
      ],
      "Resource": [
        11 * 11
      ],
      "Sid": "AllowEMRCONTAINERSStartjobrun"
  ]
}
```

Permissões para a depuração e o monitoramento

Para obter acesso aos registros enviados para o Amazon S3 e CloudWatch, ou para visualizar os registros de eventos do aplicativo no console do Amazon EMR, crie uma política do IAM com as seguintes permissões. Recomendamos redefinir as permissões para recursos específicos sempre que possível com base nas suas necessidades de negócios.

▲ Important

Se você não criou um bucket do Amazon S3, será necessário adicionar a permissão s3:CreateBucket à instrução de política. Se você não criou um grupo de logs, será necessário adicionar logs: CreateLogGroup à instrução de política.

JSON

```
{
 "Version": "2012-10-17",
 "Statement": [
      "Effect": "Allow",
      "Action": [
        "emr-containers:DescribeJobRun",
        "elasticmapreduce:CreatePersistentAppUI",
        "elasticmapreduce:DescribePersistentAppUI",
        "elasticmapreduce:GetPersistentAppUIPresignedURL"
      ],
      "Resource": [
        11 * 11
      "Sid": "AllowEMRCONTAINERSDescribejobrun"
    },
      "Effect": "Allow",
      "Action": [
        "s3:GetObject",
        "s3:ListBucket"
      ],
      "Resource": [
        11 * 11
      "Sid": "AllowS3Getobject"
   },
      "Effect": "Allow",
      "Action": [
        "logs:Get*",
        "logs:DescribeLogGroups",
        "logs:DescribeLogStreams"
```

```
],
    "Resource": [
         "*"
        ],
        "Sid": "AllowLOGSGet"
      }
]
```

Para obter mais informações sobre como configurar uma execução de trabalho para enviar registros para o Amazon S3 CloudWatch, consulte Configurar uma execução de trabalho para usar registros do S3 e Configurar uma execução de trabalho para usar registros. CloudWatch

Registro do cluster do Amazon EKS com o Amazon EMR

Registrar o cluster é a etapa final obrigatória para configurar o Amazon EMR no EKS para executar workloads.

Use o comando apresentado a seguir para criar um cluster virtual com um nome de sua escolha para o cluster e para o namespace do Amazon EKS configurados nas etapas anteriores.

Note

Cada cluster virtual deve ter um nome exclusivo em todos os clusters do EKS. Se dois clusters virtuais tiverem o mesmo nome, o processo de implantação falhará mesmo que os dois clusters virtuais pertençam a clusters do EKS diferentes.

```
aws emr-containers create-virtual-cluster \
--name virtual_cluster_name \
--container-provider '{
    "id": "cluster_name",
    "type": "EKS",
    "info": {
        "eksInfo": {
            "namespace_name"
        }
    }
}'
```

Como alternativa, você pode criar um arquivo JSON que inclua os parâmetros obrigatórios para o cluster virtual e, em seguida, executar o comando create-virtual-cluster com o caminho para o arquivo JSON. Para obter mais informações, consulte Gerenciamento de clusters virtuais.



Note

Para validar a criação com êxito de um cluster virtual, visualize o status dos clusters virtuais ao usar a operação list-virtual-clusters ou acessar a página Clusters virtuais no console do Amazon EMR.

Envio de uma execução de trabalho com StartJobRun

Enviar uma execução de trabalho com um arquivo JSON com parâmetros especificados

Crie um arquivo start-job-run-request. json e especifique os parâmetros obrigatórios para a execução de trabalho, como demonstrado pelo arquivo JSON de exemplo a seguir. Para obter mais informações sobre os parâmetros, consulte Opções para a configuração de uma execução de trabalho.

```
{
  "name": "myjob",
  "virtualClusterId": "123456",
  "executionRoleArn": "iam_role_name_for_job_execution",
  "releaseLabel": "emr-6.2.0-latest",
  "jobDriver": {
    "sparkSubmitJobDriver": {
      "entryPoint": "entryPoint_location",
      "entryPointArguments": ["argument1", "argument2", ...],
       "sparkSubmitParameters": "--class <main_class> --conf
 spark.executor.instances=2 --conf spark.executor.memory=2G --conf
 spark.executor.cores=2 --conf spark.driver.cores=1"
    }
 },
  "configurationOverrides": {
    "applicationConfiguration": [
        "classification": "spark-defaults",
        "properties": {
          "spark.driver.memory":"2G"
```

```
}
],
"monitoringConfiguration": {
    "persistentAppUI": "ENABLED",
    "cloudWatchMonitoringConfiguration": {
        "logGroupName": "my_log_group",
        "logStreamNamePrefix": "log_stream_prefix"
    },
    "s3MonitoringConfiguration": {
        "logUri": "s3://my_s3_log_location"
    }
}
```

2. Use o comando start-job-run com um caminho para o arquivo start-job-runrequest.json armazenado localmente.

```
aws emr-containers start-job-run \
--cli-input-json file://./start-job-run-request.json
```

Iniciar uma execução de trabalho usando o comando start-job-run

1. Forneça todos os parâmetros especificados no comando StartJobRun, como demonstrado pelo exemplo a seguir.

```
aws emr-containers start-job-run \
--virtual-cluster-id 123456 \
--name myjob \
--execution-role-arn execution-role-arn \
--release-label emr-6.2.0-latest \
--job-driver '{"sparkSubmitJobDriver": {"entryPoint": "entryPoint_location",
    "entryPointArguments": ["argument1", "argument2", ...], "sparkSubmitParameters":
    "--class <main_class> --conf spark.executor.instances=2 --conf
    spark.executor.memory=26 --conf spark.executor.cores=2 --conf
    spark.driver.cores=1"}}' \
--configuration-overrides '{"applicationConfiguration": [{"classification":
    "spark-defaults", "properties": {"spark.driver.memory": "2G"}}],
    "monitoringConfiguration": {"cloudWatchMonitoringConfiguration":
    {"log_group_name", "logStreamNamePrefix": "log_stream_prefix"},
```

```
"persistentAppUI":"ENABLED", "s3MonitoringConfiguration": {"logUri":
"s3://my_s3_log_location" }}}'
```

2. Para o Spark SQL, forneça todos os parâmetros especificados no comando StartJobRun, como demonstrado pelo exemplo a seguir.

Uso da classificação de envio de trabalho

Visão geral

A solicitação StartJobRun do Amazon EMR no EKS cria um pod de envio de trabalho (também conhecido como pod job-runner) para gerar o driver do Spark. Você pode usar a emr-job-submitter classificação para configurar seletores de nós para o pod remetente do trabalho, bem como definir a imagem, a CPU e a memória para o contêiner de registro do pod remetente do trabalho.

As seguintes configurações estão disponíveis na emr-job-submitter classificação:

jobsubmitter.node.selector.[labelKey]

É adicionada ao seletor de nó do pod de envio de trabalho, com a chave <code>labelKey</code> e o valor como o valor de configuração para a configuração. Por exemplo, você pode definir <code>jobsubmitter.node.selector.identifier</code> como myIdentifier e o pod de envio de trabalho terá um seletor de nó com um valor de identificador de chave myIdentifier. Isso pode

ser usado para especificar em quais nós o pod remetente do trabalho pode ser colocado. Para adicionar diversas chaves seletoras de nós, defina diversas configurações com esse prefixo.

jobsubmitter.logging.image

Define uma imagem personalizada a ser usada para o contêiner de registro no pod remetente do trabalho.

jobsubmitter.logging.request.cores

Define um valor personalizado para o número CPUs, em unidades de CPU, do contêiner de registro no pod remetente do trabalho. Por padrão, isso é definido como 100m.

jobsubmitter.logging.request.memory

Define um valor personalizado para a quantidade de memória, em bytes, para o contêiner de registro no pod remetente do trabalho. Por padrão, isso é definido como 200Mi. Um mebibyte é uma unidade de medida semelhante a um megabyte.

Recomendamos colocar pods de remetentes de trabalhos em instâncias sob demanda. A colocação de pods do remetente do trabalho em instâncias spot pode resultar em uma falha no trabalho se a instância em que o pod do remetente do trabalho é executado estiver sujeita a uma interrupção da instância spot. Você também pode <u>estabelecer o pod de envio de trabalho em uma única zona de</u> disponibilidade ou usar quaisquer rótulos do Kubernetes aplicados aos nós.

Exemplos de classificação de envio de trabalho

Nesta seção

- Solicitação StartJobRun com nó sob demanda estabelecido para o pod de envio de trabalho
- Solicitação StartJobRun com nó em uma única zona de disponibilidade estabelecido para o pod de envio de trabalho
- <u>StartJobRunsolicitação com posicionamento do tipo de EC2 instância Single-AZ e Amazon para o</u> pod remetente do trabalho
- StartJobRunsolicitação com imagem de contêiner de registro personalizada, CPU e memória

Solicitação **StartJobRun** com nó sob demanda estabelecido para o pod de envio de trabalho

```
cat >spark-python-in-s3-nodeselector-job-submitter.json << EOF
{
    "name": "spark-python-in-s3-nodeselector",</pre>
```

```
"virtualClusterId": "virtual-cluster-id",
  "executionRoleArn": "execution-role-arn",
  "releaseLabel": "emr-6.11.0-latest",
  "jobDriver": {
    "sparkSubmitJobDriver": {
      "entryPoint": "s3://S3-prefix/trip-count.py",
       "sparkSubmitParameters": "--conf spark.driver.cores=5 --conf
 spark.executor.memory=20G --conf spark.driver.memory=15G --conf
 spark.executor.cores=6"
    }
  },
  "configurationOverrides": {
    "applicationConfiguration": [
        "classification": "spark-defaults",
        "properties": {
          "spark.dynamicAllocation.enabled":"false"
      },
        "classification": "emr-job-submitter",
        "properties": {
            "jobsubmitter.node.selector.eks.amazonaws.com/capacityType": "ON_DEMAND"
        }
      }
    ],
    "monitoringConfiguration": {
      "cloudWatchMonitoringConfiguration": {
        "logGroupName": "/emr-containers/jobs",
        "logStreamNamePrefix": "demo"
      },
      "s3MonitoringConfiguration": {
        "logUri": "s3://joblogs"
      }
    }
  }
}
EOF
aws emr-containers start-job-run --cli-input-json file:///spark-python-in-s3-
nodeselector-job-submitter.json
```

Solicitação **StartJobRun** com nó em uma única zona de disponibilidade estabelecido para o pod de envio de trabalho

```
cat >spark-python-in-s3-nodeselector-job-submitter-az.json << EOF</pre>
  "name": "spark-python-in-s3-nodeselector",
  "virtualClusterId": "virtual-cluster-id",
  "executionRoleArn": "execution-role-arn",
  "releaseLabel": "emr-6.11.0-latest",
  "jobDriver": {
    "sparkSubmitJobDriver": {
      "entryPoint": "s3://S3-prefix/trip-count.py",
       "sparkSubmitParameters": "--conf spark.driver.cores=5 --conf
 spark.executor.memory=20G --conf spark.driver.memory=15G --conf
 spark.executor.cores=6"
    }
  },
  "configurationOverrides": {
    "applicationConfiguration": [
        "classification": "spark-defaults",
        "properties": {
          "spark.dynamicAllocation.enabled":"false"
         }
      },
        "classification": "emr-job-submitter",
        "properties": {
            "jobsubmitter.node.selector.topology.kubernetes.io/zone": "Availability
 Zone"
        }
      }
    ],
    "monitoringConfiguration": {
      "cloudWatchMonitoringConfiguration": {
        "logGroupName": "/emr-containers/jobs",
        "logStreamNamePrefix": "demo"
      },
      "s3MonitoringConfiguration": {
        "logUri": "s3://joblogs"
      }
    }
```

```
}
EOF
aws emr-containers start-job-run --cli-input-json file:///spark-python-in-s3-
nodeselector-job-submitter-az.json
```

StartJobRunsolicitação com posicionamento do tipo de EC2 instância Single-AZ e Amazon para o pod remetente do trabalho

```
{
  "name": "spark-python-in-s3-nodeselector",
  "virtualClusterId": "virtual-cluster-id",
  "executionRoleArn": "execution-role-arn",
  "releaseLabel": "emr-6.11.0-latest",
  "jobDriver": {
    "sparkSubmitJobDriver": {
      "entryPoint": "s3://S3-prefix/trip-count.py",
       "sparkSubmitParameters": "--conf spark.driver.cores=5 --conf
 spark.kubernetes.pyspark.pythonVersion=3 --conf spark.executor.memory=20G
 --conf spark.driver.memory=15G --conf spark.executor.cores=6 --conf
 spark.sql.shuffle.partitions=1000"
    }
  },
  "configurationOverrides": {
    "applicationConfiguration": [
      {
        "classification": "spark-defaults",
        "properties": {
          "spark.dynamicAllocation.enabled": "false",
         }
      },
        "classification": "emr-job-submitter",
        "properties": {
            "jobsubmitter.node.selector.topology.kubernetes.io/zone": "Availability
 Zone",
            "jobsubmitter.node.selector.node.kubernetes.io/instance-type":"<u>m5.4xlarge</u>"
        }
    ],
    "monitoringConfiguration": {
      "cloudWatchMonitoringConfiguration": {
        "logGroupName": "/emr-containers/jobs",
        "logStreamNamePrefix": "demo"
```

```
},
    "s3MonitoringConfiguration": {
        "logUri": "s3://joblogs"
     }
    }
}
```

StartJobRunsolicitação com imagem de contêiner de registro personalizada, CPU e memória

```
{
  "name": "spark-python",
  "virtualClusterId": "virtual-cluster-id",
  "executionRoleArn": "execution-role-arn",
  "releaseLabel": "emr-6.11.0-latest",
  "jobDriver": {
    "sparkSubmitJobDriver": {
      "entryPoint": "s3://S3-prefix/trip-count.py"
    }
  },
  "configurationOverrides": {
    "applicationConfiguration": [
      {
        "classification": "emr-job-submitter",
        "properties": {
            "jobsubmitter.logging.image": "YOUR_ECR_IMAGE_URL",
            "jobsubmitter.logging.request.memory": "200Mi",
            "jobsubmitter.logging.request.cores": "0.5"
        }
      }
    ],
    "monitoringConfiguration": {
      "cloudWatchMonitoringConfiguration": {
        "logGroupName": "/emr-containers/jobs",
        "logStreamNamePrefix": "demo"
      },
      "s3MonitoringConfiguration": {
        "logUri": "s3://joblogs"
      }
    }
  }
}
```

Usando a classificação padrão do contêiner Amazon EMR

Visão geral

As seguintes configurações estão disponíveis na emr-containers-defaults classificação:

job-start-timeout

Por padrão, um trabalho atingirá o tempo limite se não puder ser iniciado e esperar no SUBMITTED estado por 15 minutos. Essa configuração altera o número de segundos de espera antes que o trabalho expire.

logging.image

Define uma imagem personalizada a ser usada para o contêiner de registro nos pods do driver e do executor.

logging.request.cores

Define um valor personalizado para o número CPUs, em unidades de CPU, do contêiner de registro nos pods do driver e do executor. Por padrão, isso não está definido.

logging.request.memory

Define um valor personalizado para a quantidade de memória, em bytes, para o contêiner de registro nos pods do driver e do executor. Por padrão, isso é definido como 512Mi. Um mebibyte é uma unidade de medida semelhante a um megabyte.

Exemplos de classificação de envio de trabalho

Nesta seção

- StartJobRunsolicitação com tempo limite de trabalho personalizado
- StartJobRunsolicitação com imagem de contêiner de registro personalizada, CPU e memória

StartJobRunsolicitação com tempo limite de trabalho personalizado

```
"name": "spark-python",
  "virtualClusterId": "virtual-cluster-id",
  "executionRoleArn": "execution-role-arn",
  "releaseLabel": "emr-6.11.0-latest",
```

```
"jobDriver": {
    "sparkSubmitJobDriver": {
      "entryPoint": "s3://S3-prefix/trip-count.py"
    }
  },
  "configurationOverrides": {
    "applicationConfiguration": [
      {
        "classification": "emr-containers-defaults",
        "properties": {
            "job-start-timeout": "1800"
        }
      }
    ],
    "monitoringConfiguration": {
      "cloudWatchMonitoringConfiguration": {
        "logGroupName": "/emr-containers/jobs",
        "logStreamNamePrefix": "demo"
      },
      "s3MonitoringConfiguration": {
        "logUri": "s3://joblogs"
    }
  }
}
```

StartJobRunsolicitação com imagem de contêiner de registro personalizada, CPU e memória

```
"logging.image": "YOUR_ECR_IMAGE_URL",
            "logging.request.memory": "200Mi",
            "logging.request.cores": "0.5"
        }
      }
    ],
    "monitoringConfiguration": {
      "cloudWatchMonitoringConfiguration": {
        "logGroupName": "/emr-containers/jobs",
        "logStreamNamePrefix": "demo"
      },
      "s3MonitoringConfiguration": {
        "logUri": "s3://joblogs"
    }
  }
}
```

Execução de trabalhos do Spark com o operador do Spark

As versões 6.10.0 e superiores do Amazon EMR oferecem suporte ao operador do Kubernetes para Apache Spark, ou ao operador do Spark, como um modelo de envio de trabalho para o Amazon EMR no EKS. Com o operador do Spark, você pode implantar e gerenciar aplicações do Spark com o runtime da versão do Amazon EMR em seus próprios clusters do Amazon EKS. Após implantar o operador do Spark no cluster do Amazon EKS, você poderá enviar aplicações do Spark diretamente ao operador. O operador gerencia o ciclo de vida das aplicações do Spark.

Note

O Amazon EMR calcula preços no Amazon EKS com base no consumo de memória e vCPU. Esse cálculo se aplica aos pods de drivers e de executores. Esse cálculo começa quando você baixa a imagem da aplicação do Amazon EMR até que o pod do Amazon EKS seja encerrado e arredondado para o segundo mais próximo.

Tópicos

- Configuração do operador do Spark para o Amazon EMR no EKS
- Conceitos básicos do operador do Spark para o Amazon EMR no EKS
- Uso do ajuste de escala automático vertical com o operador do Spark para Amazon EMR no EKS

Operador do Spark 132

- Desinstalação do operador do Spark para o Amazon EMR no EKS
- Usando a configuração de monitoramento para monitorar o operador do Spark Kubernetes e as tarefas do Spark
- Segurança e o operador do Spark com o Amazon EMR no EKS

Configuração do operador do Spark para o Amazon EMR no EKS

Conclua as tarefas apresentadas a seguir para se preparar antes de instalar o operador do Spark no Amazon EKS. Se você já se inscreveu na Amazon Web Services (AWS) e usou o Amazon EKS, está com quase tudo pronto para usar o Amazon EMR no EKS. Conclua as tarefas apresentadas a seguir para se preparar para usar o operador do Spark no Amazon EKS. Se você já completou algum dos pré-requisitos, pode ignorá-los e passar para os próximos.

- Instale ou atualize para a versão mais recente do AWS CLI Se você já instalou o AWS CLI, confirme se você tem a versão mais recente.
- <u>Instale o kubectl e o eksctl</u> : o eksctl é uma ferramenta de linha de comando que você usa para se comunicar com o Amazon EKS.
- <u>Instale o Helm</u>: o gerenciador de pacotes Helm para o Kubernetes ajuda a instalar e gerenciar aplicações em seu cluster do Kubernetes.
- Comece a usar o Amazon EKS (eksctl): siga as etapas para criar um cluster do Kubernetes com nós no Amazon EKS.
- <u>Selecione um URI de imagem base do Amazon EMR</u> (versão 6.10.0 ou superiores): o operador do Spark é compatível com as versões 6.10.0 e superiores do Amazon EMR.

Conceitos básicos do operador do Spark para o Amazon EMR no EKS

Este tópico ajuda você a começar a usar o operador do Spark no Amazon EKS ao implantar uma aplicação do Spark e uma aplicação programada do Spark.

Instalação do operador do Spark

Use as etapas apresentadas a seguir para instalar o operador do Kubernetes para Apache Spark.

 Caso ainda não tenha feito, conclua as etapas em <u>Configuração do operador do Spark para o</u> Amazon EMR no EKS.

Configuração 133

2. Autentique seu cliente Helm no registro do Amazon ECR. No comando a seguir, substitua region-id os valores pelos seus preferidos Região da AWS e pelo ECR-registry-account valor correspondente para a Região na Contas de registro do Amazon ECR por região página.

```
aws ecr get-login-password \
--region region-id | helm registry login \
--username AWS \
--password-stdin ECR-registry-account.dkr.ecr.region-id.amazonaws.com
```

3. Instale o operador do Spark com o comando apresentado a seguir.

Para o parâmetro --version do chart do Helm, use o rótulo de versão do Amazon EMR com o prefixo emr- e o sufixo de data removidos. Por exemplo, com a versão emr-6.12.0-java17-latest, especifique 6.12.0-java17. O exemplo com o comando apresentado a seguir usa a versão emr-7.9.0-latest, portanto, especifica 7.9.0 para o chart do Helm --version.

```
helm install spark-operator-demo \
    oci://895885662937.dkr.ecr.region-id.amazonaws.com/spark-operator \
    --set emrContainers.awsRegion=region-id \
    --version 7.9.0 \
    --namespace spark-operator \
    --create-namespace
```

Por padrão, o comando cria uma conta de serviço emr-containers-sa-spark-operator para o operador do Spark. Para usar uma conta de serviço diferente, forneça o argumento serviceAccounts.sparkoperator.name. Por exemplo:

```
--set serviceAccounts.sparkoperator.name my-service-account-for-spark-operator
```

Para <u>usar o escalonamento automático vertical com o operador Spark</u>, adicione a seguinte linha ao comando de instalação para permitir webhooks para o operador:

```
--set webhook.enable=true
```

4. Verifique se você instalou o chart do Helm com o comando helm list:

```
helm list --namespace spark-operator -o yaml
```

O comando helm list deve retornar as informações da versão do chart do Helm recémimplantado:

```
app_version: v1beta2-1.3.8-3.1.1
chart: spark-operator-7.9.0
name: spark-operator-demo
namespace: spark-operator
revision: "1"
status: deployed
updated: 2023-03-14 18:20:02.721638196 +0000 UTC
```

5. Conclua a instalação com quaisquer opções adicionais necessárias. Para obter mais informações, consulte a spark-on-k8s-operatordocumentação em. GitHub

Execução de uma aplicação do Spark

O operador do Spark é compatível com o Amazon EMR 6.10.0 ou com versões superiores. Ao instalar o operador do Spark, ele cria a conta de serviço emr-containers-sa-spark para executar aplicações do Spark por padrão. Use as etapas apresentadas a seguir para executar uma aplicação do Spark com o operador do Spark no Amazon EMR no EKS 6.10.0 ou em versões superiores.

- Antes de executar uma aplicação do Spark com o operador do Spark, conclua as etapas em <u>Configuração do operador do Spark para o Amazon EMR no EKS</u> e <u>Instalação do operador do</u> Spark.
- Crie um arquivo de definição SparkApplication spark-pi.yaml com o seguinte conteúdo de exemplo:

```
apiVersion: "sparkoperator.k8s.io/v1beta2"
kind: SparkApplication
metadata:
   name: spark-pi
   namespace: spark-operator
spec:
   type: Scala
   mode: cluster
   image: "895885662937.dkr.ecr.us-west-2.amazonaws.com/spark/emr-6.10.0:latest"
   imagePullPolicy: Always
   mainClass: org.apache.spark.examples.SparkPi
```

```
mainApplicationFile: "local:///usr/lib/spark/examples/jars/spark-examples.jar"
sparkVersion: "3.3.1"
restartPolicy:
  type: Never
volumes:
  - name: "test-volume"
    hostPath:
      path: "/tmp"
      type: Directory
driver:
  cores: 1
  coreLimit: "1200m"
 memory: "512m"
 labels:
    version: 3.3.1
  serviceAccount: emr-containers-sa-spark
  volumeMounts:
    - name: "test-volume"
      mountPath: "/tmp"
executor:
  cores: 1
  instances: 1
 memory: "512m"
  labels:
    version: 3.3.1
  volumeMounts:
    - name: "test-volume"
      mountPath: "/tmp"
```

3. Agora, envie a aplicação do Spark com o comando apresentado a seguir. Isso também criará um objeto SparkApplication chamado spark-pi:

```
kubectl apply -f spark-pi.yaml
```

4. Verifique os eventos do objeto SparkApplication com o seguinte comando:

```
kubectl describe sparkapplication spark-pi --namespace spark-operator
```

Para obter mais informações sobre como enviar inscrições para o Spark por meio do operador do Spark, consulte <u>Usando um SparkApplication na spark-on-k8s-operator documentação</u> sobre. GitHub

Uso do Amazon S3 para armazenamento

Para usar o Amazon S3 como opção de armazenamento de arquivos, adicione as configurações a seguir ao arquivo YAML.

```
hadoopConf:
# EMRFS filesystem
  fs.s3.customAWSCredentialsProvider:
 com.amazonaws.auth.WebIdentityTokenCredentialsProvider
  fs.s3.impl: com.amazon.ws.emr.hadoop.fs.EmrFileSystem
  fs.AbstractFileSystem.s3.impl: org.apache.hadoop.fs.s3.EMRFSDelegate
  fs.s3.buffer.dir: /mnt/s3
  fs.s3.getObject.initialSocketTimeoutMilliseconds: "2000"
  mapreduce.fileoutputcommitter.algorithm.version.emr_internal_use_only.EmrFileSystem:
 "2"
  mapreduce.fileoutputcommitter.cleanup-
failures.ignored.emr_internal_use_only.EmrFileSystem: "true"
sparkConf:
 # Required for EMR Runtime
 spark.driver.extraClassPath: /usr/lib/hadoop-lzo/lib/*:/usr/lib/hadoop/hadoop-
aws.jar:/usr/share/aws/aws-java-sdk/*:/usr/share/aws/emr/emrfs/conf:/usr/share/aws/
emr/emrfs/lib/*:/usr/share/aws/emr/emrfs/auxlib/*:/usr/share/aws/emr/security/conf:/
usr/share/aws/emr/security/lib/*:/usr/share/aws/hmclient/lib/aws-glue-datacatalog-
spark-client.jar:/usr/share/java/Hive-JSON-Serde/hive-openx-serde.jar:/usr/share/aws/
sagemaker-spark-sdk/lib/sagemaker-spark-sdk.jar:/home/hadoop/extrajars/*
 spark.driver.extraLibraryPath: /usr/lib/hadoop/lib/native:/usr/lib/hadoop-lzo/lib/
native:/docker/usr/lib/hadoop/lib/native:/docker/usr/lib/hadoop-lzo/lib/native
 spark.executor.extraClassPath: /usr/lib/hadoop-lzo/lib/*:/usr/lib/hadoop/hadoop-
aws.jar:/usr/share/aws/aws-java-sdk/*:/usr/share/aws/emr/emrfs/conf:/usr/share/aws/
emr/emrfs/lib/*:/usr/share/aws/emr/emrfs/auxlib/*:/usr/share/aws/emr/security/conf:/
usr/share/aws/emr/security/lib/*:/usr/share/aws/hmclient/lib/aws-glue-datacatalog-
spark-client.jar:/usr/share/java/Hive-JSON-Serde/hive-openx-serde.jar:/usr/share/aws/
sagemaker-spark-sdk/lib/sagemaker-spark-sdk.jar:/home/hadoop/extrajars/*
 spark.executor.extraLibraryPath: /usr/lib/hadoop/lib/native:/usr/lib/hadoop-lzo/lib/
native:/docker/usr/lib/hadoop/lib/native:/docker/usr/lib/hadoop-lzo/lib/native
```

Se você usa o Amazon EMR nas versões 7.2.0 e superiores, as configurações são incluídas por padrão. Nesse caso, você pode definir o caminho do arquivo como s3://

s3://

sample of the configurações são incluídas por padrão. Nesse caso, você pode definir o caminho do arquivo como s3://

sample of the configurações são incluídas por padrão. Nesse caso, você pode definir o caminho do arquivo como s3://
sample of the configurações são incluídas por padrão. Nesse caso, você pode definir o caminho do arquivo como s3://
sample of the configurações são incluídas por padrão. Nesse caso, você pode definir o caminho do arquivo como s3://
sample of the configurações são incluídas por padrão. Nesse caso, você pode definir o caminho do arquivo como s3://
sample of the configurações são incluídas por padrão. Nesse caso, você pode definir o caminho do arquivo como s3://
sample of the configurações são incluídas por padrão. Nesse caso, você pode definir o caminho do arquivo como sample of the configurações padrão incluídas por padrão. Nesse caso, você pode definir o caminho do arquivo como sample of the configurações padrão incluídas padrão

Em seguida, envie a aplicação do Spark normalmente.

Uso do ajuste de escala automático vertical com o operador do Spark para Amazon EMR no EKS

A partir do Amazon EMR 7.0, você pode usar o ajuste de escala automático vertical do Amazon EMR no EKS para simplificar o gerenciamento de recursos. Ela ajusta automaticamente os recursos de memória e de CPU para se adaptar às necessidades da workload fornecida para aplicações do Spark no Amazon EMR. Para obter mais informações, consulte <u>Uso da escalabilidade automática vertical com trabalhos do Spark no Amazon EMR</u>.

Esta seção descreve como configurar o operador Spark para usar o ajuste de escala automático vertical.

Pré-requisitos

Antes de configurar o monitoramento, certifique-se de concluir as seguintes tarefas de configuração:

- Siga as etapas em Configuração do operador do Spark para o Amazon EMR no EKS.
- (opcional) Se você instalou anteriormente uma versão mais antiga do operador Spark, exclua o SparkApplication/ScheduledSparkApplication CRD.

```
kubectl delete crd sparkApplication
kubectl delete crd scheduledSparkApplication
```

• Siga as etapas em <u>Instalação do operador do Spark</u>. Na etapa 3, adicione a seguinte linha ao comando de instalação para permitir webhooks para o operador:

```
--set webhook.enable=true
```

- Siga as etapas em Configuração da escalabilidade automática vertical para o Amazon EMR no EKS.
- Dê acesso aos arquivos na localização do Amazon S3:
 - 1. Anote sua conta de serviço do driver e do operador com o JobExecutionRole que tem as permissões do S3.

```
kubectl annotate serviceaccount -n spark-operator emr-containers-sa-spark
  eks.amazonaws.com/role-arn=JobExecutionRole
kubectl annotate serviceaccount -n spark-operator emr-containers-sa-spark-
operator eks.amazonaws.com/role-arn=JobExecutionRole
```

2. Atualize a política de confiança do perfil de execução de trabalho nesse namespace.

```
aws emr-containers update-role-trust-policy \
--cluster-name cluster \
--namespace ${Namespace}\
--role-name iam_role_name_for_job_execution
```

3. Edite a política de confiança do perfil do IAM do perfil de execução de trabalhos e atualize a serviceaccount de emr-containers-sa-spark-*-*-xxxx para emr-containers-sa-*.

```
{
    "Effect": "Allow",
    "Principal": {
        "Federated": "OIDC-provider"
},
    "Action": "sts:AssumeRoleWithWebIdentity",
    "Condition": {
        "StringLike": {
            "OIDC": "system:serviceaccount:${Namespace}:emr-containers-sa-*"
        }
    }
}
```

4. Se estiver usando o Amazon S3 como armazenamento de arquivos, adicione os padrões a seguir ao arquivo YAML.

```
hadoopConf:

# EMRFS filesystem
fs.s3.customAWSCredentialsProvider:
com.amazonaws.auth.WebIdentityTokenCredentialsProvider
fs.s3.impl: com.amazon.ws.emr.hadoop.fs.EmrFileSystem
fs.AbstractFileSystem.s3.impl: org.apache.hadoop.fs.s3.EMRFSDelegate
fs.s3.buffer.dir: /mnt/s3
fs.s3.getObject.initialSocketTimeoutMilliseconds: "2000"

mapreduce.fileoutputcommitter.algorithm.version.emr_internal_use_only.EmrFileSystem:
"2"
mapreduce.fileoutputcommitter.cleanup-
failures.ignored.emr_internal_use_only.EmrFileSystem: "true"
sparkConf:
# Required for EMR Runtime
```

```
spark.driver.extraClassPath: /usr/lib/hadoop-lzo/lib/*:/usr/lib/hadoop/hadoop-
aws.jar:/usr/share/aws/aws-java-sdk/*:/usr/share/aws/emr/emrfs/conf:/usr/share/
aws/emr/emrfs/lib/*:/usr/share/aws/emr/emrfs/auxlib/*:/usr/share/aws/emr/
security/conf:/usr/share/aws/emr/security/lib/*:/usr/share/aws/hmclient/lib/aws-
glue-datacatalog-spark-client.jar:/usr/share/java/Hive-JSON-Serde/hive-openx-
serde.jar:/usr/share/aws/sagemaker-spark-sdk/lib/sagemaker-spark-sdk.jar:/home/
hadoop/extrajars/*
spark.driver.extraLibraryPath: /usr/lib/hadoop/lib/native:/usr/lib/hadoop-lzo/
lib/native:/docker/usr/lib/hadoop/lib/native:/docker/usr/lib/hadoop-lzo/lib/
native
spark.executor.extraClassPath: /usr/lib/hadoop-lzo/lib/*:/usr/lib/hadoop/
hadoop-aws.jar:/usr/share/aws/aws-java-sdk/*:/usr/share/aws/emr/emrfs/conf:/usr/
share/aws/emr/emrfs/lib/*:/usr/share/aws/emr/emrfs/auxlib/*:/usr/share/aws/emr/
security/conf:/usr/share/aws/emr/security/lib/*:/usr/share/aws/hmclient/lib/aws-
glue-datacatalog-spark-client.jar:/usr/share/java/Hive-JSON-Serde/hive-openx-
serde.jar:/usr/share/aws/sagemaker-spark-sdk/lib/sagemaker-spark-sdk.jar:/home/
hadoop/extrajars/*
spark.executor.extraLibraryPath: /usr/lib/hadoop/lib/native:/usr/lib/hadoop-
lzo/lib/native:/docker/usr/lib/hadoop/lib/native:/docker/usr/lib/hadoop-lzo/lib/
native
```

Execute um trabalho com ajuste de escala automático vertical no operador do Spark

Antes de executar uma aplicação do Spark com o operador do Spark, conclua as etapas em <u>Prérequisitos</u>.

Para usar o ajuste de escala automático vertical com o operador do Spark, adicione a seguinte configuração ao driver da especificação da aplicação do Spark para ativar o ajuste de escala automático vertical:

```
dynamicSizing:
  mode: Off
  signature: "my-signature"
```

Essa configuração permite o ajuste de escala automático vertical e é uma configuração de assinatura obrigatória que permite escolher uma assinatura para o trabalho.

Para obter mais informações sobre as configurações e os valores dos parâmetros, consulte Configuring vertical autoscaling for Amazon EMR on EKS. Por padrão, seu trabalho é enviado com a opção somente monitoramento no modo Desativado para a escalabilidade automática vertical. Esse

estado de monitoramento permite calcular e visualizar recomendações de recursos sem realizar a escalabilidade automática. Para obter mais informações, consulte Vertical autoscaling modes.

O exemplo a seguir mostra um arquivo de definição SparkApplication chamado sparkpi.yaml com as configurações necessárias para usar ajuste de escala automático vertical.

```
apiVersion: "sparkoperator.k8s.io/v1beta2"
kind: SparkApplication
metadata:
  name: spark-pi
  namespace: spark-operator
spec:
  type: Scala
  mode: cluster
  image: "895885662937.dkr.ecr.us-west-2.amazonaws.com/spark/emr-7.9.0:latest"
  imagePullPolicy: Always
  mainClass: org.apache.spark.examples.SparkPi
  mainApplicationFile: "local:///usr/lib/spark/examples/jars/spark-examples.jar"
  sparkVersion: "3.4.1"
  dynamicSizing:
    mode: Off
    signature: "my-signature"
  restartPolicy:
    type: Never
  volumes:
    - name: "test-volume"
      hostPath:
        path: "/tmp"
        type: Directory
  driver:
    cores: 1
    coreLimit: "1200m"
    memory: "512m"
    labels:
      version: 3.4.1
    serviceAccount: emr-containers-sa-spark
    volumeMounts:
      - name: "test-volume"
        mountPath: "/tmp"
  executor:
    cores: 1
    instances: 1
    memory: "512m"
```

```
labels:
```

version: 3.4.1
volumeMounts:

- name: "test-volume"
mountPath: "/tmp"

Agora, envie a aplicação do Spark com o comando apresentado a seguir. Isso também criará um objeto SparkApplication chamado spark-pi:

```
kubectl apply -f spark-pi.yaml
```

Para obter mais informações sobre como enviar inscrições para o Spark por meio do operador do Spark, consulte <u>Usando um SparkApplication na spark-on-k8s-operator documentação</u> sobre. GitHub

Verificação de funcionalidade da escalabilidade automática vertical

Para verificar se a escalabilidade automática vertical está funcionando corretamente para o trabalho enviado, use kubectl para obter o recurso personalizado verticalpodautoscaler e visualizar suas recomendações de escalabilidade.

```
kubectl get verticalpodautoscalers --all-namespaces \
-l=emr-containers.amazonaws.com/dynamic.sizing.signature=my-signature
```

A saída desta consulta deve ser semelhante à seguinte:

```
NAMESPACE NAME
CPU MEM PROVIDED AGE
spark-operator ds-p73j6mkosvc4xeb3gr7x4xol2bfcw5evqimzqojrlysvj3giozuq-vpa Off
580026651 True 15m
```

Se a saída não for semelhante ou tiver um código de erro, consulte <u>Solução de problemas de</u> <u>escalabilidade automática vertical do Amazon EMR no EKS</u> para obter passos que ajudam a resolver o problema.

Para remover pods e aplicações, execute o seguinte comando:

```
kubectl delete sparkapplication spark-pi
```

Desinstalação do operador do Spark para o Amazon EMR no EKS

Use as etapas apresentadas a seguir para desinstalar o operador do Spark.

 Exclua o operador do Spark usando o namespace correto. Para este exemplo, o namespace é spark-operator-demo.

```
helm uninstall spark-operator-demo -n spark-operator
```

2. Exclua a conta de serviço do operador do Spark:

```
kubectl delete sa emr-containers-sa-spark-operator -n spark-operator
```

Exclua o operador do Spark CustomResourceDefinitions (CRDs):

```
kubectl delete crd sparkapplications.sparkoperator.k8s.io
kubectl delete crd scheduledsparkapplications.sparkoperator.k8s.io
```

Usando a configuração de monitoramento para monitorar o operador do Spark Kubernetes e as tarefas do Spark

A configuração de monitoramento permite que você configure facilmente o arquivamento de registros do seu aplicativo Spark e dos registros do operador no Amazon S3 ou no. Amazon CloudWatch Você pode escolher um ou ambos. Isso adiciona um sidecar do agente de log aos seus pods de operador, driver e executor do Spark e, posteriormente, encaminha os registros desses componentes para seus coletores configurados.

Pré-requisitos

Antes de configurar o monitoramento, certifique-se de concluir as seguintes tarefas de configuração:

 (Opcional) Se você instalou anteriormente uma versão mais antiga do operador Spark, exclua o SparkApplicationScheduledSparkApplication/CRD.

```
kubectl delete crd scheduledsparkapplications.sparkoperator.k8s.io
kubectl delete crd sparkapplications.sparkoperator.k8s.io
```

2. Crie uma função de operator/job execução no IAM se você ainda não tiver uma.

Desinstalar 143

3. Execute o comando a seguir para atualizar a política de confiança da função de operator/job execução que você acabou de criar:

```
aws emr-containers update-role-trust-policy \
--cluster-name cluster \
--namespace namespace \
--role-name iam_role_name_for_operator/job_execution_role
```

4. Edite a política de confiança da função do IAM da sua função de operator/job execução da seguinte forma:

5. Crie uma política de MonitoringConfiguration no IAM com as seguintes permissões:

JSON

```
"arn:aws:logs:*:*:log-group:log_group_name:*"
      ],
      "Sid": "AllowLOGSDescribelogstreams"
    },
    {
      "Effect": "Allow",
      "Action": [
        "logs:DescribeLogGroups"
      ],
      "Resource": [
      ],
      "Sid": "AllowLOGSDescribeloggroups"
    },
      "Effect": "Allow",
      "Action": [
        "s3:PutObject",
        "s3:GetObject",
        "s3:ListBucket"
      ],
      "Resource": [
        "arn:aws:s3:::bucket_name",
        "arn:aws:s3:::bucket_name/*"
      ],
      "Sid": "AllowS3Putobject"
  ]
}
```

6. Anexe a política acima à sua função operator/job de execução.

Registros do operador do Spark

Você pode definir a configuração de monitoramento da seguinte forma ao fazerhelm install:

```
helm install spark-operator spark-operator \
--namespace namespace \
--set emrContainers.awsRegion=aws_region \
--set emrContainers.monitoringConfiguration.image=log_agent_image_url \
--set
emrContainers.monitoringConfiguration.s3MonitoringConfiguration.logUri=S3_bucket_uri \
```

```
--set
emrContainers.monitoringConfiguration.cloudWatchMonitoringConfiguration.logGroupName=log_group
\
--set
emrContainers.monitoringConfiguration.cloudWatchMonitoringConfiguration.logStreamNamePrefix=log
\
--set emrContainers.monitoringConfiguration.sideCarResources.limits.cpuLimit=500m \
--set emrContainers.monitoringConfiguration.sideCarResources.limits.memoryLimit=512Mi \
--set
emrContainers.monitoringConfiguration.containerLogRotationConfiguration.rotationSize=2GB
\
--set
emrContainers.monitoringConfiguration.containerLogRotationConfiguration.maxFilesToKeep=10
\
--set webhook.enable=true \
--set emrContainers.operatorExecutionRoleArn=operator_execution_role_arn
```

Configuração de monitoramento

A seguir estão as opções de configuração disponíveis em MonitoringConfiguration.

- Imagem (opcional) URL da imagem do agente de log. Será buscado emrReleaseLabel se não for fornecido.
- s3 MonitoringConfiguration Defina essa opção para arquivar no Amazon S3.
 - logURI (obrigatório) O caminho do bucket do Amazon S3 em que você deseja armazenar seus registros.
 - A seguir estão exemplos de formatos para os caminhos de bucket do Amazon S3, após o upload dos registros. O primeiro exemplo mostra que não há rotação de log ativada.

```
s3://${logUri}/${POD NAME}/operator/stdout.gz
s3://${logUri}/${POD NAME}/operator/stderr.gz
```

A rotação de registros está ativada por padrão. Você pode ver tanto um arquivo rotacionado, com um índice incrementado, quanto um arquivo atual, que é o mesmo da amostra anterior.

```
s3://${logUri}/${POD NAME}/operator/stdout_YYYYMMDD_index.gz
s3://${logUri}/${POD NAME}/operator/stderr_YYYYMMDD_index.gz
```

 cloudWatchMonitoringConfiguração — A chave de configuração para a qual configurar o encaminhamento. Amazon CloudWatch

- logGroupName(obrigatório) Nome do grupo de Amazon CloudWatch registros para o qual você deseja enviar registros. O grupo é criado automaticamente, se não existir.
- logStreamNamePrefixo (opcional) Nome do fluxo de registros para o qual você deseja enviar os registros. O valor padrão é uma string vazia. O formato em Amazon CloudWatch é o seguinte:

```
${logStreamNamePrefix}/${POD NAME}/STDOUT or STDERR
```

- sideCarResources(opcional) A chave de configuração para definir limites de recursos no contêiner auxiliar Fluentd lançado.
 - MemoryLimit (opcional) O limite de memória. Ajuste de acordo com suas necessidades. O padrão é de 512Mi.
 - CPULimit (opcional) O limite da CPU. Ajuste de acordo com suas necessidades. O padrão é 500m.
- containerLogRotationConfiguração (opcional) controla o comportamento de rotação do log do contêiner. Ele é habilitado por padrão.
 - RotationSize (obrigatório) Especifica o tamanho do arquivo para a rotação do log. O intervalo de valores possíveis é de 2 KB a 2 GB. A parcela numérica da unidade do parâmetro rotationSize é transferida como um número inteiro. Como não há suporte para valores decimais, você pode especificar um tamanho de rotação de 1,5 GB, por exemplo, com o valor 1.500 MB. O padrão é 2 GB.
 - maxFilesToManter (obrigatório) Especifica o número máximo de arquivos a serem retidos no contêiner após a rotação. O valor mínimo é de 1 e o valor máximo é de 50. O padrão é 10.

Depois de configurar o MonitoringConfiguration, você deve ser capaz de verificar os registros do pod do operador do Spark em um bucket do Amazon S3 ou em ambos. Amazon CloudWatch Para um bucket do Amazon S3, você precisa esperar 2 minutos para que o primeiro arquivo de log seja liberado.

Para encontrar os logins Amazon CloudWatch, você pode navegar até o seguinte: CloudWatch> Grupos de registros > > Log group namePod name/operator/stderr

Ou você pode navegar até: CloudWatch> Grupos de registros > > Log group namePod name/ operator/stdout

Registros do aplicativo Spark

Você pode definir essa configuração da maneira apresentada a seguir.

```
apiVersion: "sparkoperator.k8s.io/v1beta2"
kind: SparkApplication
metadata:
  name: spark-pi
  namespace: namespace
spec:
  type: Scala
  mode: cluster
  imagePullPolicy: Always
  mainClass: org.apache.spark.examples.SparkPi
  mainApplicationFile: "local:///usr/lib/spark/examples/jars/spark-examples.jar"
  sparkVersion: "3.3.1"
  emrReleaseLabel: emr_release_label
  executionRoleArn: job_execution_role_arn
  restartPolicy:
    type: Never
  volumes:
    - name: "test-volume"
      hostPath:
        path: "/tmp"
        type: Directory
  driver:
    cores: 1
    coreLimit: "1200m"
    memory: "512m"
    labels:
      version: 3.3.1
    volumeMounts:
      - name: "test-volume"
        mountPath: "/tmp"
  executor:
    cores: 1
    instances: 1
    memory: "512m"
    labels:
      version: 3.3.1
    volumeMounts:
      - name: "test-volume"
        mountPath: "/tmp"
  monitoringConfiguration:
```

```
image: "log_agent_image"
s3MonitoringConfiguration:
    logUri: "S3_bucket_uri"
cloudWatchMonitoringConfiguration:
    logGroupName: "log_group_name"
    logStreamNamePrefix: "log_stream_prefix"
sideCarResources:
    limits:
        cpuLimit: "500m"
        memoryLimit: "250Mi"
containerLogRotationConfiguration:
        rotationSize: "2GB"
        maxFilesToKeep: "10"
```

A seguir estão as opções de configuração disponíveis em MonitoringConfiguration.

- Imagem (opcional) URL da imagem do agente de log. Será buscado emrReleaseLabel se não for fornecido.
- s3 MonitoringConfiguration Defina essa opção para arquivar no Amazon S3.
 - logURI (obrigatório) O caminho do bucket do Amazon S3 em que você deseja armazenar seus registros. O primeiro exemplo mostra que não há rotação de registros ativada:

```
s3://${logUri}/${APPLICATION NAME}-${APPLICATION UID}/${POD NAME}/stdout.gz
s3://${logUri}/${APPLICATION NAME}-${APPLICATION UID}/${POD NAME}/stderr.gz
```

A rotação de registros é ativada por padrão. Você pode usar um arquivo girado (com índice incrementado) e um arquivo atual (um sem o carimbo de data).

```
s3://${logUri}/${APPLICATION NAME}-${APPLICATION UID}/${POD NAME}/
stdout_YYYYMMDD_index.gz
s3://${logUri}/${APPLICATION NAME}-${APPLICATION UID}/${POD NAME}/
stderr_YYYYMMDD_index.gz
```

- cloudWatchMonitoringConfiguração A chave de configuração para a qual configurar o encaminhamento. Amazon CloudWatch
 - logGroupName(obrigatório) O nome do grupo de logs do Cloudwatch para o qual você deseja enviar registros. O grupo é criado automaticamente se não existir.
 - logStreamNamePrefixo (opcional) O nome do fluxo de registros para o qual você deseja enviar os registros. O valor padrão é uma string vazia. O formato em CloudWatch é o seguinte:

\$\logStreamNamePrefix\rangle\\$\langle APPLICATION NAME\rangle -\\$\langle APPLICATION UID\rangle\\$\logStreamNamePrefix\rangle\\$\langle APPLICATION NAME\rangle -\\$\langle APPLICATION UID\rangle\\$\logPOD NAME\rangle\stderr

- sideCarResources(opcional) A chave de configuração para definir limites de recursos no contêiner auxiliar Fluentd lançado.
 - MemoryLimit (opcional) O limite de memória. Ajuste de acordo com suas necessidades. O padrão é 250Mi.
 - CPULimit O limite da CPU. Ajuste de acordo com suas necessidades. O padrão é 500m.
- containerLogRotationConfiguração (opcional) controla o comportamento de rotação do log do contêiner. Ele é habilitado por padrão.
 - RotationSize (obrigatório) Especifica o tamanho do arquivo para a rotação do log. O intervalo de valores possíveis é de 2 KB a 2 GB. A parcela numérica da unidade do parâmetro rotationSize é transferida como um número inteiro. Como não há suporte para valores decimais, você pode especificar um tamanho de rotação de 1,5 GB, por exemplo, com o valor 1.500 MB. O padrão é 2 GB.
 - maxFilesToManter (obrigatório) Especifica o número máximo de arquivos a serem retidos no contêiner após a rotação. O valor mínimo é 1. O valor máximo é 50. O padrão é 10.

Depois de configurar o MonitoringConfiguration, você deve ser capaz de verificar os registros do driver e do executor do aplicativo Spark em um bucket do Amazon S3 ou em ambos. CloudWatch Para um bucket do Amazon S3, você precisa esperar 2 minutos para que o primeiro arquivo de log seja liberado. Por exemplo, no Amazon S3, o caminho do bucket aparece da seguinte forma:

Amazon S3 > Buckets > > > **Bucket name** > stderr.gz **Spark application name - UUID Pod Name**

Ou:

Amazon S3 > Buckets > > > **Bucket name** > stdout.gz **Spark application name - UUID Pod Name**

Em CloudWatch, o caminho aparece da seguinte forma:

CloudWatch> Grupos de registros > **Log group name**>Spark application name - UUID/Pod name/stderr

Ou:

CloudWatch> Grupos de registros > **Log group name**>Spark application name - UUID/Pod name/stdout

Segurança e o operador do Spark com o Amazon EMR no EKS

Há algumas maneiras de configurar as permissões de acesso ao cluster ao usar o operador do Spark. A primeira é usar o controle de acesso baseado em perfil. O controle de acesso por perfil (RBAC) restringe o acesso com base no perfil de uma pessoa dentro de uma organização. Ele se tornou-se a principal forma de lidar com o acesso. O segundo método de acesso é assumir uma AWS Identity and Access Management função, que fornece acesso a recursos por meio de permissões específicas atribuídas.

Tópicos

- Configuração de permissões de acesso ao cluster com o controle de acesso por perfil (RBAC)
- Configuração de permissões de acesso ao cluster com perfis do IAM para contas de serviço (IRSA)

Configuração de permissões de acesso ao cluster com o controle de acesso por perfil (RBAC)

Para implantar o operador do Spark, o Amazon EMR no EKS cria dois perfis e duas contas de serviço para o operador do Spark e para as aplicações do Spark.

Tópicos

- Perfil e conta de serviço do operador
- Perfil e conta de serviço do Spark

Perfil e conta de serviço do operador

O Amazon EMR no EKS cria o perfil e a conta de serviço do operador para gerenciar SparkApplications para trabalhos do Spark e para outros recursos, como serviços.

O nome padrão para esta conta de serviço é emr-containers-sa-spark-operator.

As seguintes regras se aplicam para este perfil de serviço:

rules:

apiGroups:

```
_ ""
 resources:
  - pods
  verbs:
  _ "*"
- apiGroups:
  _ ""
  resources:
  - services
  - configmaps
  - secrets
  verbs:
  - create
 - get
  - delete
  - update
- apiGroups:
  - extensions
  - networking.k8s.io
  resources:
  - ingresses
  verbs:
 - create
  - get
  - delete
- apiGroups:
  _ ""
  resources:
  - nodes
  verbs:
  - get
- apiGroups:
  _ ""
  resources:
  - events
 verbs:
  - create
  - update
  - patch
- apiGroups:
  _ ""
  resources:
  - resourcequotas
  verbs:
```

```
- get
 - list
 - watch
- apiGroups:
 - apiextensions.k8s.io
 resources:
 - customresourcedefinitions
 verbs:
 - create
 - get
 - update
 - delete
- apiGroups:
 - admissionregistration.k8s.io
 resources:
 - mutatingwebhookconfigurations
 - validatingwebhookconfigurations
 verbs:
 - create
 - get
 - update
  - delete
- apiGroups:
 - sparkoperator.k8s.io
 resources:
 - sparkapplications
 - sparkapplications/status
 - scheduledsparkapplications
 - scheduledsparkapplications/status
 verbs:
 _ "*"
 {{- if .Values.batchScheduler.enable }}
 # required for the `volcano` batch scheduler
- apiGroups:
 - scheduling.incubator.k8s.io
 - scheduling.sigs.dev
 - scheduling.volcano.sh
 resources:
 - podgroups
 verbs:
 _ "*"
 {{- end }}
 {{ if .Values.webhook.enable }}
 apiGroups:
```

```
- batch
resources:
- jobs
verbs:
- delete
{{- end }}
```

Perfil e conta de serviço do Spark

Um pod do driver do Spark precisa de uma conta de serviço do Kubernetes no mesmo namespace que o pod. Esta conta de serviço precisa de permissões para criar, obter, listar, aplicar patches e excluir pods de executores e para criar um serviço descentralizado do Kubernetes para o driver. O driver apresentará falhas e será encerrado sem a conta de serviço, a menos que a conta de serviço padrão no namespace do pod tenha as permissões obrigatórias.

O nome padrão para esta conta de serviço é emr-containers-sa-spark.

As seguintes regras se aplicam para este perfil de serviço:

```
rules:
- apiGroups:
  _ ""
  resources:
  - pods
  verbs:
  _ "*"
- apiGroups:
  _ ""
  resources:
  - services
  verbs:
  _ "*"
- apiGroups:
  _ ""
  resources:
  - configmaps
  verbs:
  _ "*"
- apiGroups:
  resources:
  - persistentvolumeclaims
  verbs:
```

_ "*"

Configuração de permissões de acesso ao cluster com perfis do IAM para contas de serviço (IRSA)

Esta seção usa um exemplo para demonstrar como configurar uma conta de serviço do Kubernetes para assumir uma função. AWS Identity and Access Management Os pods que usam a conta de serviço podem então acessar qualquer AWS serviço que a função tenha permissão para acessar.

O exemplo apresentado a seguir executa uma aplicação do Spark para contar as palavras de um arquivo no Amazon S3. Para fazer isso, você pode configurar perfis do IAM para contas de serviço (IRSA) com a finalidade de autenticar e autorizar as contas de serviço do Kubernetes.



Note

Este exemplo usa o namespace "spark-operator" para o operador do Spark e para o namespace no qual a aplicação do Spark é enviada.

Pré-requisitos

Antes de testar o exemplo apresentado nesta página, complete os seguintes pré-requisitos:

- Prepare-se para usar o operador do Spark.
- Instalação do operador do Spark.
- Crie um bucket do Amazon S3.
- Salve seu poema favorito em um arquivo de texto chamado poem.txt e faça o upload do arquivo em seu bucket do S3. A aplicação do Spark criada nesta página realizará a leitura do conteúdo do arquivo de texto. Para obter mais informações sobre como fazer o upload de arquivos para o S3, consulte Fazer upload de um objeto para o seu bucket no Guia do usuário do Amazon Simple Storage Service.

Configuração de uma conta de serviço do Kubernetes para assumir um perfil do IAM

Use as etapas a seguir para configurar uma conta de serviço do Kubernetes para assumir um papel do IAM que os pods possam usar para acessar AWS serviços que o papel tem permissões para acessar.

1. Depois de concluir o<u>Pré-requisitos</u>, use o AWS Command Line Interface para criar um example-policy.json arquivo que permita acesso somente para leitura ao arquivo que você enviou para o Amazon S3:

```
cat >example-policy.json <<EOF
{
    "Version": "2012-10-17",
    "Statement": 「
        {
            "Effect": "Allow",
            "Action": [
                 "s3:GetObject",
                 "s3:ListBucket"
            ],
            "Resource": [
                "arn:aws:s3:::my-pod-bucket",
                "arn:aws:s3:::my-pod-bucket/*"
            ]
        }
    ]
}
EOF
```

2. Em seguida, crie uma política do IAM example-policy:

```
aws iam create-policy --policy-name example-policy --policy-document file://
example-policy.json
```

3. Depois disso, crie um perfil do IAM example-role e associe-o a uma conta de serviço do Kubernetes para o driver do Spark:

```
eksctl create iamserviceaccount --name driver-account-sa --namespace spark-operator
\
--cluster my-cluster --role-name "example-role" \
--attach-policy-arn arn:aws:iam::111122223333:policy/example-policy --approve
```

4. Crie um arquivo em YAML com as associações de perfil do cluster que são obrigatórias para a conta de serviço do driver do Spark:

```
cat >spark-rbac.yaml <<EOF
apiVersion: v1
kind: ServiceAccount</pre>
```

```
metadata:
   name: driver-account-sa
---
apiVersion: rbac.authorization.k8s.io/v1
kind: ClusterRoleBinding
metadata:
   name: spark-role
roleRef:
   apiGroup: rbac.authorization.k8s.io
   kind: ClusterRole
   name: edit
subjects:
   - kind: ServiceAccount
   name: driver-account-sa
   namespace: spark-operator
EOF
```

5. Aplique as configurações de associação de perfil do cluster:

```
kubectl apply -f spark-rbac.yaml
```

O comando kubectl deve confirmar a criação com êxito da conta:

```
serviceaccount/driver-account-sa created
clusterrolebinding.rbac.authorization.k8s.io/spark-role configured
```

Execução de uma aplicação do operador do Spark

Após <u>configurar a conta de serviço do Kubernetes</u>, você poderá executar uma aplicação do Spark que conta o número de palavras no arquivo de texto carregado como parte dos Pré-requisitos.

 Crie um novo arquivo word-count.yaml, com uma definição SparkApplication para sua aplicação de contagem de palavras.

```
cat >word-count.yaml <<EOF
apiVersion: "sparkoperator.k8s.io/v1beta2"
kind: SparkApplication
metadata:
   name: word-count
   namespace: spark-operator
spec:</pre>
```

```
type: Java
 mode: cluster
  image: "895885662937.dkr.ecr.us-west-2.amazonaws.com/spark/emr-6.10.0:latest"
  imagePullPolicy: Always
 mainClass: org.apache.spark.examples.JavaWordCount
 mainApplicationFile: local:///usr/lib/spark/examples/jars/spark-examples.jar
  arguments:
    - s3://my-pod-bucket/poem.txt
  hadoopConf:
  # EMRFS filesystem
   fs.s3.customAWSCredentialsProvider:
 com.amazonaws.auth.WebIdentityTokenCredentialsProvider
   fs.s3.impl: com.amazon.ws.emr.hadoop.fs.EmrFileSystem
   fs.AbstractFileSystem.s3.impl: org.apache.hadoop.fs.s3.EMRFSDelegate
   fs.s3.buffer.dir: /mnt/s3
   fs.s3.getObject.initialSocketTimeoutMilliseconds: "2000"
mapreduce.fileoutputcommitter.algorithm.version.emr_internal_use_only.EmrFileSystem:
 "2"
   mapreduce.fileoutputcommitter.cleanup-
failures.ignored.emr_internal_use_only.EmrFileSystem: "true"
  sparkConf:
    # Required for EMR Runtime
    spark.driver.extraClassPath: /usr/lib/hadoop-lzo/lib/*:/usr/lib/hadoop/
hadoop-aws.jar:/usr/share/aws/aws-java-sdk/*:/usr/share/aws/emr/emrfs/conf:/usr/
share/aws/emr/emrfs/lib/*:/usr/share/aws/emr/emrfs/auxlib/*:/usr/share/aws/emr/
security/conf:/usr/share/aws/emr/security/lib/*:/usr/share/aws/hmclient/lib/aws-
glue-datacatalog-spark-client.jar:/usr/share/java/Hive-JSON-Serde/hive-openx-
serde.jar:/usr/share/aws/sagemaker-spark-sdk/lib/sagemaker-spark-sdk.jar:/home/
hadoop/extrajars/*
    spark.driver.extraLibraryPath: /usr/lib/hadoop/lib/native:/usr/lib/hadoop-lzo/
lib/native:/docker/usr/lib/hadoop/lib/native:/docker/usr/lib/hadoop-lzo/lib/native
    spark.executor.extraClassPath: /usr/lib/hadoop-lzo/lib/*:/usr/lib/hadoop/
hadoop-aws.jar:/usr/share/aws/aws-java-sdk/*:/usr/share/aws/emr/emrfs/conf:/usr/
share/aws/emr/emrfs/lib/*:/usr/share/aws/emr/emrfs/auxlib/*:/usr/share/aws/emr/
security/conf:/usr/share/aws/emr/security/lib/*:/usr/share/aws/hmclient/lib/aws-
glue-datacatalog-spark-client.jar:/usr/share/java/Hive-JSON-Serde/hive-openx-
serde.jar:/usr/share/aws/sagemaker-spark-sdk/lib/sagemaker-spark-sdk.jar:/home/
hadoop/extrajars/*
    spark.executor.extraLibraryPath: /usr/lib/hadoop/lib/native:/usr/lib/hadoop-
lzo/lib/native:/docker/usr/lib/hadoop/lib/native:/docker/usr/lib/hadoop-lzo/lib/
native
  sparkVersion: "3.3.1"
 restartPolicy:
```

```
type: Never
driver:
    cores: 1
    coreLimit: "1200m"
    memory: "512m"
    labels:
       version: 3.3.1
    serviceAccount: my-spark-driver-sa
executor:
    cores: 1
    instances: 1
    memory: "512m"
    labels:
       version: 3.3.1
EOF
```

2. Envie a aplicação do Spark.

```
kubectl apply -f word-count.yaml
```

O comando kubectl deve retornar a confirmação de que você criou com êxito um objeto SparkApplication chamado word-count.

```
sparkapplication.sparkoperator.k8s.io/word-count configured
```

3. Para verificar os eventos do objeto SparkApplication, execute o seguinte comando:

```
kubectl describe sparkapplication word-count -n spark-operator
```

O comando kubecti deve retornar a descrição da SparkApplication com os eventos:

```
Events:
 Type
          Reason
                                                Age
                                                                       From
   Message
 Normal
          SparkApplicationSpecUpdateProcessed 3m2s (x2 over 17h)
                                                                       spark-
operator Successfully processed spec update for SparkApplication word-count
 Warning SparkApplicationPendingRerun
                                                3m2s (x2 over 17h)
                                                                       spark-
operator SparkApplication word-count is pending rerun
 Normal
          SparkApplicationSubmitted
                                                2m58s (x2 over 17h)
                                                                       spark-
operator SparkApplication word-count was submitted successfully
```

Normal	SparkDriverRunning	2m56s (x2 over 17h)	spark-
operator	Driver word-count-driver is running		
Normal	SparkExecutorPending	2m50s	spark-
operator	Executor [javawordcount-fdd1698807392c66-exec-1] is pending		
Normal	SparkExecutorRunning	2m48s	spark-
operator	Executor [javawordcount-fdd1698807392c66-exec-1] is running		
Normal	SparkDriverCompleted	2m31s (x2 over 17h)	spark-
operator	Driver word-count-driver completed		
Normal	SparkApplicationCompleted	2m31s (x2 over 17h)	spark-
operator	SparkApplication word-count completed		
Normal	SparkExecutorCompleted	2m31s (x2 over 2m31s)	spark-
operator	Executor [javawordcount-fdd1698807392c66-exec-1] completed		

Agora, a aplicação está realizando a contagem das palavras em seu arquivo do S3. Para localizar a contagem de palavras, consulte os arquivos de log do seu driver:

```
kubectl logs pod/word-count-driver -n spark-operator
```

O comando kubectl deve retornar o conteúdo do arquivo de log com os resultados da sua aplicação de contagem de palavras.

Para obter mais informações sobre como enviar aplicativos ao Spark por meio do operador Spark, consulte <u>Usando um SparkApplication</u> na documentação do operador Kubernetes para Apache Spark (operador 8s) em. spark-on-k GitHub

Execução de trabalhos do Spark com spark-submit

As versões 6.10.0 e superiores do Amazon EMR oferecem suporte ao spark-submit como uma ferramenta de linha de comando que você pode usar para enviar e executar aplicações do Spark para um cluster do Amazon EMR no EKS.



O Amazon EMR calcula preços no Amazon EKS com base no consumo de memória e vCPU. Esse cálculo se aplica aos pods de drivers e de executores. Esse cálculo começa quando

spark-submit 160

você baixa a imagem da aplicação do Amazon EMR até que o pod do Amazon EKS seja encerrado e arredondado para o segundo mais próximo.

Tópicos

Amazon EMR

- Configuração do spark-submit para o Amazon EMR no EKS
- Conceitos básicos do spark-submit para o Amazon EMR no EKS
- · Verificação dos requisitos de segurança da conta de serviço do driver do Spark para spark-submit

Configuração do spark-submit para o Amazon EMR no EKS

Conclua as tarefas apresentadas a seguir para se preparar antes de executar uma aplicação com spark-submit no Amazon EMR no EKS. Se você já se inscreveu na Amazon Web Services (AWS) e usou o Amazon EKS, está com quase tudo pronto para usar o Amazon EMR no EKS. Se você já completou algum dos pré-requisitos, pode ignorá-los e passar para os próximos.

- <u>Instale ou atualize para a versão mais recente do AWS CLI</u> Se você já instalou o AWS CLI, confirme se você tem a versão mais recente.
- <u>Instale o kubectl e o eksctl</u>: o eksctl é uma ferramenta de linha de comando que você usa para se comunicar com o Amazon EKS.
- Comece a usar o Amazon EKS (eksctl): siga as etapas para criar um cluster do Kubernetes com nós no Amazon EKS.
- <u>Selecione um URI de imagem base do Amazon EMR</u> (versão 6.10.0 ou superiores): o comando spark-submit é compatível com as versões 6.10.0 e superiores do Amazon EMR.
- Confirme se a conta de serviço do driver tem permissões apropriadas para criar e monitorar os pods de executores. Para obter mais informações, consulte <u>Verificação dos requisitos de</u> segurança da conta de serviço do driver do Spark para spark-submit.
- Configure seu perfil de credenciais locais da AWS.
- No console do Amazon EKS, escolha seu cluster do EKS e encontre o endpoint do cluster do EKS, localizado em Visão geral, Detalhes e, em seguida, Endpoint do servidor de API.

Configuração 161

Conceitos básicos do spark-submit para o Amazon EMR no EKS

O Amazon EMR 6.10.0 e versões superiores oferecem suporte ao spark-submit para a execução de aplicações do Spark em um cluster do Amazon EKS. A seção a seguir mostra como enviar um comando para uma aplicação do Spark.

Execução de uma aplicação do Spark

Para executar a aplicação do Spark, siga estas etapas:

- 1. Antes de executar uma aplicação do Spark com o comando spark-submit, conclua as etapas em Configuração do spark-submit para o Amazon EMR no EKS.
- 2. Execute um contêiner com uma imagem base do Amazon EMR no EKS. Consulte <u>How to select</u> a base image URI para obter mais informações.

```
kubectl run -it containerName --image=EMRonEKSImage --command -n namespace /bin/
bash
```

3. Defina os valores para as seguintes variáveis de ambiente:

```
export SPARK_HOME=spark-home
export MASTER_URL=k8s://Amazon EKS-cluster-endpoint
```

4. Agora, envie a aplicação do Spark com o seguinte comando:

```
$SPARK_HOME/bin/spark-submit \
--class org.apache.spark.examples.SparkPi \
--master $MASTER_URL \
--conf spark.kubernetes.container.image=895885662937.dkr.ecr.us-
west-2.amazonaws.com/spark/emr-6.10.0:latest \
--conf spark.kubernetes.authenticate.driver.serviceAccountName=spark \
--deploy-mode cluster \
--conf spark.kubernetes.namespace=spark-operator \
local:///usr/lib/spark/examples/jars/spark-examples.jar 20
```

Para obter mais informações sobre o envio de aplicações para o Spark, consulte <u>Submitting</u> <u>applications</u> na documentação do Apache Spark.



M Important

O spark-submit oferece suporte somente para o modo de cluster como o mecanismo de envio.

Verificação dos requisitos de segurança da conta de serviço do driver do Spark para spark-submit

O pod do driver do Spark usa uma conta de serviço do Kubernetes para acessar o servidor da API do Kubernetes para criar e monitorar pods de executores. A conta de serviço do driver deve ter as permissões apropriadas para listar, criar, editar, aplicar patches e excluir pods em seu cluster. É possível verificar se consegue listar esses recursos ao executar o seguinte comando:

```
kubectl auth can-i list|create|edit|delete|patch pods
```

Verifique se você tem as permissões necessárias executando cada comando.

```
kubectl auth can-i list pods
kubectl auth can-i create pods
kubectl auth can-i edit pods
kubectl auth can-i delete pods
kubectl auth can-i patch pods
```

As seguintes regras se aplicam para este perfil de serviço:

```
rules:
- apiGroups:
  _ ""
  resources:
  - pods
  verbs:
  _ "*"
- apiGroups:
  _ ""
  resources:
  - services
  verbs:
  _ "*"
 apiGroups:
```

```
resources:
- configmaps
verbs:
- "*"
- apiGroups:
- ""
resources:
- persistentvolumeclaims
verbs:
- "*"
```

Configuração de perfis do IAM para contas de serviço (IRSA) do spark-submit

As seções a seguir explicam como configurar perfis do IAM para contas de serviço (IRSA) com a finalidade de autenticar e autorizar as contas de serviço do Kubernetes para que você possa executar aplicações do Spark armazenadas no Amazon S3.

Pré-requisitos

Antes de testar qualquer um dos exemplos desta documentação, certifique-se de ter cumprido os seguintes pré-requisitos:

- · Concluir a configuração do spark-submit
- Criar um bucket do S3 e fazer upload do JAR da aplicação do Spark

Configuração de uma conta de serviço do Kubernetes para assumir um perfil do IAM

As etapas a seguir abordam como configurar uma conta de serviço do Kubernetes para assumir um perfil do AWS Identity and Access Management (IAM). Depois de configurar os pods para usar a conta de serviço, eles poderão acessar qualquer um AWS service (Serviço da AWS) que a função tenha permissão para acessar.

1. Crie um arquivo de política para permitir acesso somente leitura ao objeto do Amazon S3 do qual você fez upload:

2. Crie a política do IAM.

```
aws iam create-policy --policy-name my-policy --policy-document file://my-policy.json
```

3. Crie um perfil do IAM e associe-o a uma conta de serviço do Kubernetes para o driver do Spark

```
eksctl create iamserviceaccount --name my-spark-driver-sa --namespace spark-
operator \
--cluster my-cluster --role-name "my-role" \
--attach-policy-arn arn:aws:iam::111122223333:policy/my-policy --approve
```

4. Crie um arquivo YAML com as <u>permissões</u> obrigatórias para a conta de serviço do driver do Spark:

```
cat >spark-rbac.yaml <<EOF
apiVersion: rbac.authorization.k8s.io/v1
kind: Role
metadata:
   namespace: default
   name: emr-containers-role-spark
rules:
   - apiGroups:
   - ""
   resources:
   - pods
   verbs:
   - "*"
   - apiGroups:</pre>
```

```
_ ""
  resources:
  - services
  verbs:
  _ "*"
- apiGroups:
  _ ""
  resources:
  - configmaps
  verbs:
  _ "*"
- apiGroups:
  _ ""
 resources:
  - persistentvolumeclaims
  verbs:
  _ "*"
apiVersion: rbac.authorization.k8s.io/v1
kind: RoleBinding
metadata:
  name: spark-role-binding
  namespace: default
roleRef:
  apiGroup: rbac.authorization.k8s.io
  kind: Role
  name: emr-containers-role-spark
subjects:
- kind: ServiceAccount
  name: emr-containers-sa-spark
  namespace: default
EOF
```

5. Aplique as configurações de associação de perfil do cluster.

```
kubectl apply -f spark-rbac.yaml
```

6. O comando kubect1 deve retornar a confirmação da conta criada.

```
serviceaccount/emr-containers-sa-spark created clusterrolebinding.rbac.authorization.k8s.io/emr-containers-role-spark configured
```

Execução da aplicação do Spark

O Amazon EMR 6.10.0 e versões superiores oferecem suporte ao spark-submit para a execução de aplicações do Spark em um cluster do Amazon EKS. Para executar a aplicação do Spark, siga estas etapas:

- 1. Certifique-se de ter concluído as etapas em Setting up spark-submit for Amazon EMR on EKS.
- 2. Defina os valores para as seguintes variáveis de ambiente:

```
export SPARK_HOME=spark-home
export MASTER_URL=k8s://Amazon EKS-cluster-endpoint
```

3. Agora, envie a aplicação do Spark com o seguinte comando:

```
$SPARK_HOME/bin/spark-submit \
 --class org.apache.spark.examples.SparkPi \
 --master $MASTER_URL \
 --conf spark.kubernetes.container.image=895885662937.dkr.ecr.us-
west-2.amazonaws.com/spark/emr-6.15.0:latest \
 --conf spark.kubernetes.authenticate.driver.serviceAccountName=emr-containers-sa-
spark \
 --deploy-mode cluster \
 --conf spark.kubernetes.namespace=default \
 --conf "spark.driver.extraClassPath=/usr/lib/hadoop-lzo/lib/*:/usr/lib/hadoop/
hadoop-aws.jar:/usr/share/aws/aws-java-sdk/*:/usr/share/aws/emr/emrfs/conf:/usr/
share/aws/emr/emrfs/lib/*:/usr/share/aws/emr/emrfs/auxlib/*:/usr/share/aws/emr/
security/conf:/usr/share/aws/emr/security/lib/*:/usr/share/aws/hmclient/lib/aws-
glue-datacatalog-spark-client.jar:/usr/share/java/Hive-JSON-Serde/hive-openx-
serde.jar:/usr/share/aws/sagemaker-spark-sdk/lib/sagemaker-spark-sdk.jar:/home/
hadoop/extrajars/*" \
 --conf "spark.driver.extraLibraryPath=/usr/lib/hadoop/lib/native:/usr/lib/hadoop-
lzo/lib/native:/docker/usr/lib/hadoop/lib/native:/docker/usr/lib/hadoop-lzo/lib/
native" \
--conf "spark.executor.extraClassPath=/usr/lib/hadoop-lzo/lib/*:/usr/lib/hadoop/
hadoop-aws.jar:/usr/share/aws/aws-java-sdk/*:/usr/share/aws/emr/emrfs/conf:/usr/
share/aws/emr/emrfs/lib/*:/usr/share/aws/emr/emrfs/auxlib/*:/usr/share/aws/emr/
security/conf:/usr/share/aws/emr/security/lib/*:/usr/share/aws/hmclient/lib/aws-
glue-datacatalog-spark-client.jar:/usr/share/java/Hive-JSON-Serde/hive-openx-
serde.jar:/usr/share/aws/sagemaker-spark-sdk/lib/sagemaker-spark-sdk.jar:/home/
hadoop/extrajars/*" \
 --conf "spark.executor.extraLibraryPath=/usr/lib/hadoop/lib/native:/usr/lib/
hadoop-lzo/lib/native:/docker/usr/lib/hadoop/lib/native:/docker/usr/lib/hadoop-lzo/
lib/native" \
```

```
--conf
spark.hadoop.fs.s3.customAWSCredentialsProvider=com.amazonaws.auth.WebIdentityTokenCredent
\
--conf spark.hadoop.fs.s3.impl=com.amazon.ws.emr.hadoop.fs.EmrFileSystem \
--conf
spark.hadoop.fs.AbstractFileSystem.s3.impl=org.apache.hadoop.fs.s3.EMRFSDelegate \
--conf spark.hadoop.fs.s3.buffer.dir=/mnt/s3 \
--conf spark.hadoop.fs.s3.getObject.initialSocketTimeoutMilliseconds="2000" \
--conf
spark.hadoop.mapreduce.fileoutputcommitter.algorithm.version.emr_internal_use_only.EmrFileSystem="true" \
--conf spark.hadoop.mapreduce.fileoutputcommitter.cleanup-
failures.ignored.emr_internal_use_only.EmrFileSystem="true" \
s3://my-pod-bucket/spark-examples.jar 20
```

4. Depois que o driver do Spark terminar o trabalho do Spark, você verá uma linha de log no final do envio indicando que o trabalho do Spark foi concluído.

```
23/11/24 17:02:14 INFO LoggingPodStatusWatcherImpl: Application org.apache.spark.examples.SparkPi with submission ID default:org-apache-spark-examples-sparkpi-4980808c03ff3115-driver finished 23/11/24 17:02:14 INFO ShutdownHookManager: Shutdown hook called
```

Limpeza

Quando terminar de executar as aplicações, você poderá realizar a limpeza com o comando a seguir.

```
kubectl delete -f spark-rbac.yaml
```

Uso do Apache Livy com o Amazon EMR no EKS

Com as versões 7.1.0 e posteriores do Amazon EMR, é possível usar o Apache Livy para enviar trabalhos no Amazon EMR no EKS. Usando o Apache Livy, você pode configurar seu próprio endpoint REST do Apache Livy e usá-lo para implantar e gerenciar aplicações do Spark em clusters do Amazon EKS. Após instalar o Livy no cluster do Amazon EKS, você poderá usar o endpoint do Livy para enviar aplicações do Spark ao servidor do Livy. O servidor gerencia o ciclo de vida das aplicações do Spark.

Apache Livy 168



Note

O Amazon EMR calcula preços no Amazon EKS com base no consumo de memória e vCPU. Esse cálculo se aplica aos pods de drivers e de executores. Esse cálculo começa quando você baixa a imagem da aplicação do Amazon EMR até que o pod do Amazon EKS seja encerrado e arredondado para o segundo mais próximo.

Tópicos

- Configuração do Apache Livy para o Amazon EMR no EKS
- Conceitos básicos do Apache Livy para Amazon EMR no EKS
- Execução de uma aplicação do Spark com o Apache Livy para o Amazon EMR no EKS
- Desinstalação do Apache Livy com o Amazon EMR no EKS
- Segurança para o Apache Livy com o Amazon EMR no EKS
- Propriedades de instalação do Apache Livy em versões do Amazon EMR no EKS
- Solução de erros comuns de formato de variável de ambiente

Configuração do Apache Livy para o Amazon EMR no EKS

Antes de instalar o Apache Livy no cluster do Amazon EKS, é necessário instalar e configurar um conjunto de ferramentas de pré-requisitos. Isso inclui o AWS CLI, que é uma ferramenta de linha de comando fundamental para trabalhar com AWS recursos, ferramentas de linha de comando para trabalhar com o Amazon EKS e um controlador usado nesse caso de uso para disponibilizar seu aplicativo de cluster na Internet e rotear o tráfego de rede.

- Instale ou atualize para a versão mais recente do AWS CLI Se você já instalou o AWS CLI, confirme se você tem a versão mais recente.
- Instale o kubectl e o eksctl : o eksctl é uma ferramenta de linha de comando que você usa para se comunicar com o Amazon EKS.
- Instale o Helm: o gerenciador de pacotes Helm para o Kubernetes ajuda a instalar e gerenciar aplicações em seu cluster do Kubernetes.
- Comece a usar o Amazon EKS (eksctl): siga as etapas para criar um cluster do Kubernetes com nós no Amazon EKS.
- Selecione um rótulo de versão do Amazon EMR: o Apache Livy é compatível com as versões 7.1.0 e superiores do Amazon EMR.

Configuração 169 Instale o controlador ALB — o controlador ALB gerencia o AWS Elastic Load Balancing para clusters Kubernetes. Ele cria um AWS Network Load Balancer (NLB) quando você cria um Kubernetes Ingress enquanto configura o Apache Livy.

Conceitos básicos do Apache Livy para Amazon EMR no EKS

Conclua as etapas a seguir para instalar o Apache Livy. Eles incluem a configuração do gerenciador de pacotes, a criação de um namespace para executar workloads do Spark, a instalação do Livy, a configuração do balanceamento de carga e as etapas de verificação. Você precisa concluir essas etapas para executar um trabalho em lotes com o Spark.

- 1. Se ainda não tiver feito isso, configure o Apache Livy para o Amazon EMR no EKS.
- Autentique seu cliente Helm no registro do Amazon ECR. Você pode encontrar o ECRregistry-account valor correspondente para suas contas Região da AWS de registro do Amazon ECR por região.

```
aws ecr get-login-password \--region < AWS_REGION> | helm registry login \
--username AWS \
--password-stdin <eccount>.dkr.ecr.<region-id>.amazonaws.com
```

- A configuração do Livy cria uma conta de serviço para o servidor do Livy e outra conta para a aplicação do Spark. Para configurar o IRSA nas contas de serviço, consulte Setting up access permissions with IAM roles for service accounts (IRSA).
- 4. Crie um namespace para executar as workloads do Spark.

```
kubectl create ns <spark-ns>
```

Use o comando a seguir para instalar o Livy. 5.

Esse endpoint do Livy só está disponível internamente para a VPC no cluster do EKS. Para habilitar o acesso além da VPC, defina —-set loadbalancer.internal=false no comando de instalação do Helm.



Note

Por padrão, o SSL não está habilitado nesse endpoint do Livy, e o endpoint só é visível dentro da VPC do cluster do EKS. Se você definir loadbalancer.internal=false e ssl.enabled=false, estará expondo um endpoint inseguro fora da sua VPC. Para

configurar um endpoint seguro do Livy, consulte Configuring a secure Apache Livy endpoint with TLS/SSL.

```
helm install livy-demo \
   oci://895885662937.dkr.ecr.region-id.amazonaws.com/livy \
   --version 7.10.0 \
   --namespace livy-ns \
   --set image=ECR-registry-account.dkr.ecr.region-id.amazonaws.com/livy/
emr-7.10.0:latest \
   --set sparkNamespace=<spark-ns> \
   --create-namespace
```

Você verá a saída a seguir.

```
NAME: livy-demo
LAST DEPLOYED: Mon Mar 18 09:23:23 2024
NAMESPACE: livy-ns
STATUS: deployed
REVISION: 1
TEST SUITE: None
NOTES:
The Livy server has been installed.
Check installation status:
1. Check Livy Server pod is running
  kubectl --namespace livy-ns get pods -l "app.kubernetes.io/instance=livy-demo"
2. Verify created NLB is in Active state and it's target groups are healthy (if
loadbalancer.enabled is true)
Access LIVY APIs:
    # Ensure your NLB is active and healthy
    # Get the Livy endpoint using command:
    LIVY_ENDPOINT=$(kubectl get svc -n livy-ns -l app.kubernetes.io/
instance=livy-demo,emr-containers.amazonaws.com/type=loadbalancer -o
 jsonpath='{.items[0].status.loadBalancer.ingress[0].hostname}' | awk '{printf
 "%s:8998\n", $0}')
    # Access Livy APIs using http://$LIVY_ENDPOINT or https://$LIVY_ENDPOINT (if
 SSL is enabled)
    # Note: While uninstalling Livy, makes sure the ingress and NLB are deleted
 after running the helm command to avoid dangling resources
```

Os nomes padrão da conta de serviço para o servidor do Livy e a sessão do Spark são emr-containers-sa-livy e emr-containers-sa-spark-livy. Para usar nomes personalizados, utilize os parâmetros serviceAccounts.name e sparkServiceAccount.name.

```
--set serviceAccounts.name=my-service-account-for-livy
--set sparkServiceAccount.name=my-service-account-for-spark
```

Verifique se você instalou o chart do Helm.

```
helm list -n livy-ns -o yaml
```

O comando helm list deve retornar as informações sobre o novo chart do Helm.

```
app_version: 0.7.1-incubating
chart: livy-emr-7.10.0
name: livy-demo
namespace: livy-ns
revision: "1"
status: deployed
updated: 2024-02-08 22:39:53.539243 -0800 PST
```

Verifique se o Network Load Balancer está ativo.

```
LIVY_NAMESPACE=<livy-ns>
LIVY_APP_NAME=<livy-app-name>
AWS_REGION>

# Get the NLB Endpoint URL
NLB_ENDPOINT=$(kubectl --namespace $LIVY_NAMESPACE get svc -1 "app.kubernetes.io/
instance=$LIVY_APP_NAME,emr-containers.amazonaws.com/type=loadbalancer" -o
jsonpath='{.items[0].status.loadBalancer.ingress[0].hostname}')

# Get all the load balancers in the account's region
ELB_LIST=$(aws elbv2 describe-load-balancers --region $AWS_REGION)

# Get the status of the NLB that matching the endpoint from the Kubernetes service
NLB_STATUS=$(echo $ELB_LIST | grep -A 8 "\"DNSName\": \"$NLB_ENDPOINT\"" | awk '/
Code/{print $2}/}' | tr -d '"},\n')
echo $NLB_STATUS
```

8. Agora, verifique se o grupo de destino no Network Load Balancer está íntegro.

```
LIVY_NAMESPACE=Vy-ns>
LIVY_APP_NAME=<livy-app-name>
AWS_REGION=<AWS_REGION>
# Get the NLB endpoint
NLB_ENDPOINT=$(kubectl --namespace $LIVY_NAMESPACE get svc -1 "app.kubernetes.io/
instance=$LIVY_APP_NAME,emr-containers.amazonaws.com/type=loadbalancer" -o
jsonpath='{.items[0].status.loadBalancer.ingress[0].hostname}')
# Get all the load balancers in the account's region
ELB_LIST=$(aws elbv2 describe-load-balancers --region $AWS_REGION)
# Get the NLB ARN from the NLB endpoint
NLB_ARN=$(echo $ELB_LIST | grep -B 1 "\"DNSName\": \"$NLB_ENDPOINT\"" | awk
 '/"LoadBalancerArn":/,/"/'| awk '/:/{print $2}' | tr -d \",)
# Get the target group from the NLB. Livy setup only deploys 1 target group
TARGET_GROUP_ARN=$(aws elbv2 describe-target-groups --load-balancer-arn $NLB_ARN
--region $AWS_REGION | awk '/"TargetGroupArn":/,/"/' awk '/:/{print $2}' | tr -d
\",)
# Get health of target group
aws elbv2 describe-target-health --target-group-arn $TARGET_GROUP_ARN
```

Este é um exemplo de saída que mostra o status do grupo de destino:

}

Quando o status do NLB se tornar active e o grupo de destino healthy, você poderá continuar. Pode demorar alguns minutos.

9. Recupere o endpoint do Livy da instalação do Helm. Se o endpoint do Livy é seguro ou não, depende de você ter habilitado o SSL.

```
LIVY_NAMESPACE=<livy-ns>
LIVY_APP_NAME=livy-app-name
LIVY_ENDPOINT=$(kubectl get svc -n livy-ns -l app.kubernetes.io/
instance=livy-app-name, emr-containers.amazonaws.com/type=loadbalancer -o
jsonpath='{.items[0].status.loadBalancer.ingress[0].hostname}' | awk '{printf
"%s:8998\n", $0}')
echo "$LIVY_ENDPOINT"
```

10. Recuperação da conta de serviço do Spark na instalação do Helm

```
SPARK_NAMESPACE=spark-ns
LIVY_APP_NAME=<livy-app-name>
SPARK_SERVICE_ACCOUNT=$(kubectl --namespace $SPARK_NAMESPACE
  get sa -l "app.kubernetes.io/instance=$LIVY_APP_NAME" -o
    jsonpath='{.items[0].metadata.name}')
echo "$SPARK_SERVICE_ACCOUNT"
```

Você deverá ver algo semelhante ao resultado a seguir:

```
emr-containers-sa-spark-livy
```

- 11. Se você configurar internalALB=true para habilitar o acesso de fora da sua VPC, crie uma EC2 instância da Amazon e certifique-se de que o Network Load Balancer permita tráfego de rede proveniente da instância. EC2 Você deve fazer isso para que a instância tenha acesso ao endpoint do Livy. Para obter mais informações sobre como expor com segurança o endpoint fora da VPC, consulte Setting up with a secure Apache Livy endpoint with TLS/SSL.
- 12. A instalação do Livy cria a conta de serviço emr-containers-sa-spark para executar as aplicações do Spark. Se seu aplicativo Spark usa AWS recursos como o S3 ou chama operações de AWS API ou CLI, você deve vincular uma função do IAM com as permissões necessárias à sua conta de serviço do Spark. Para obter mais informações, consulte Setting up access permissions with IAM roles for service accounts (IRSA).

O Apache Livy oferece suporte a configurações adicionais que você pode usar ao instalar o Livy. Para obter mais informações, consulte Propriedades de instalação do Apache Livy no Amazon EMR nas versões do EKS.

Execução de uma aplicação do Spark com o Apache Livy para o Amazon EMR no EKS

Antes de executar uma aplicação do Spark com o Apache Livy, certifique-se de ter concluído as etapas em Configuração do Apache Livy para o Amazon EMR no EKS e Getting started with Apache Livy for Amazon EMR on EKS.

Você pode usar o Apache Livy para executar dois tipos de aplicações:

- Sessões em lote: um tipo de workload do Livy para enviar trabalhos em lote do Spark.
- Sessões interativas: um tipo de workload do Livy que fornece uma interface programática e visual para executar consultas do Spark.

Note

Pods de drivers e de executores de sessões diferentes podem se comunicar entre si. Os namespaces não garantem nenhuma segurança entre os pods. O Kubernetes não aceita permissões seletivas em um subconjunto de pods dentro de um determinado namespace.

Execução de sessões em lote

Para enviar um trabalho em lote, use o comando a seguir.

```
}
}' < livy-endpoint > / batches
```

Para monitorar seu trabalho em lote, use o comando a seguir.

```
curl -s -k -H 'Content-Type: application/json' -X GET <livy-endpoint>/batches/my-
session
```

Execução de sessões interativas

Para executar sessões interativas com o Apache Livy, consulte as etapas a seguir.

- 1. Certifique-se de ter acesso a um notebook Jupyter auto-hospedado ou gerenciado, como um notebook SageMaker Al Jupyter. O caderno Jupyter deve ter o sparkmagic instalado.
- 2. Crie um bucket para a configuração spark.kubernetes.file.upload.path do Spark. Certifique-se de que a conta de serviço do Spark tenha acesso de leitura e gravação ao bucket. Para obter mais detalhes sobre como configurar a conta de serviço do Spark, consulte Setting up access permissions with IAM roles for service accounts (IRSA)
- 3. Faça upload do sparkmagic no caderno Jupyter com o comando %load_ext sparkmagic.magics.
- 4. Execute o comando %manage_spark para configurar o endpoint do Livy com o caderno Jupyter. Selecione a guia Adicionar endpoints, escolha o tipo de autenticação configurado, adicione o endpoint do Livy ao caderno e escolha Adicionar endpoint.
- 5. Execute %manage_spark novamente para criar o contexto do Spark e vá para Criar sessão. Escolha o endpoint do Livy, especifique um nome de sessão exclusivo, escolha um idioma e adicione as propriedades a seguir.

```
{
   "conf": {
      "spark.kubernetes.namespace": "livy-namespace",
      "spark.kubernetes.container.image": "public.ecr.aws/emr-on-eks/spark/
emr-7.10.0:latest",
      "spark.kubernetes.authenticate.driver.serviceAccountName": "<spark-service-account>",
      "spark.kubernetes.file.upload.path": "<URI_TO_S3_LOCATION_>"
   }
}
```

Envie a aplicação e espere que ela crie o contexto do Spark.

7. Para monitorar o status da sessão interativa, execute o comando a seguir.

```
curl -s -k -H 'Content-Type: application/json' -X GET livy-endpoint/sessions/my-interactive-session
```

Monitoramento de aplicações do Spark

Para monitorar o progresso das aplicações do Spark com a interface de usuário do Livy, use o link http://<livy-endpoint>/ui.

Desinstalação do Apache Livy com o Amazon EMR no EKS

Siga estas etapas para desinstalar o Apache Livy.

1. Exclua a configuração do Livy usando os nomes do namespace e o nome da aplicação. Neste exemplo, o nome da aplicação é livy-demo e o namespace é livy-ns.

```
helm uninstall livy-demo -n livy-ns
```

- 2. Ao desinstalar, o Amazon EMR no EKS exclui o serviço Kubernetes no Livy, os balanceadores de carga e os AWS grupos-alvo que você criou durante a instalação. A exclusão de recursos pode levar alguns minutos. Certifique-se de que os recursos sejam excluídos antes de instalar o Livy no namespace novamente.
- 3. Exclua o namespace do Spark.

```
kubectl delete namespace spark-ns
```

Segurança para o Apache Livy com o Amazon EMR no EKS

Consulte os tópicos a seguir para saber mais sobre como configurar a segurança do Apache Livy com o Amazon EMR no EKS. Essas opções incluem o uso da segurança da camada de transporte, o controle de acesso baseado em perfis, que é o acesso baseado no perfil de uma pessoa dentro de uma organização, e o uso de perfis do IAM, que fornecem acesso aos recursos de acordo com as permissões concedidas.

Tópicos

Configuração de um endpoint do Apache Livy seguro com TLS/SSL

Desinstalação 1777

- Configuração das permissões de aplicações do Apache Livy e Spark com o controle de acesso por perfil (RBAC)
- Configuração de permissões de acesso com perfis do IAM para contas de serviço (IRSA)

Configuração de um endpoint do Apache Livy seguro com TLS/SSL

Consulte as seções a seguir para saber mais sobre como configurar o Apache Livy para Amazon EMR no EKS com criptografia end-to-end TLS e SSL.

Configuração da criptografia TLS e SSL

Para configurar a criptografia SSL no endpoint do Apache Livy, siga estas etapas.

- Instale o driver CSI do Secrets Store e o AWS Secrets and Configuration Provider (ASCP) o
 driver CSI e o ASCP do Secrets Store armazenam com segurança os certificados e senhas JKS do
 Livy que o pod do servidor Livy precisa para habilitar o SSL. Você também pode instalar apenas o
 driver CSI do Secrets Store e usar qualquer outro provedor de segredos compatível.
- <u>Crie um certificado do ACM</u>: esse certificado é necessário para proteger a conexão entre o cliente e o endpoint do ALB.
- Configure um certificado JKS, uma senha de chave e uma senha de armazenamento de chaves para AWS Secrets Manager — necessários para proteger a conexão entre o endpoint ALB e o servidor Livy.
- Adicione permissões à conta de serviço Livy para recuperar segredos AWS Secrets Manager
 — o servidor Livy precisa dessas permissões para recuperar segredos do ASCP e adicionar as
 configurações Livy para proteger o servidor Livy. Para adicionar permissões do IAM a uma conta
 de serviço, consulte Setting up access permissions with IAM roles for service accounts (IRSA).

Configurando um certificado JKS com uma chave e uma senha do keystore para AWS Secrets Manager

Siga estas etapas para configurar um certificado do JKS com uma chave e uma senha do keystore.

1. Gere um arquivo do keystore para o servidor do Livy.

```
keytool -genkey -alias <host> -keyalg RSA -keysize 2048 -dname
CN=<host>,OU=hw,O=hw,L=<your_location>,ST=<state>,C=<country> -
keypass <keyPassword> -keystore <keystore_file> -storepass <storePassword> --
validity 3650
```

2. Crie um certificado.

```
keytool -export -alias <host> -keystore mykeystore.jks -rfc -
file mycertificate.cert -storepass <storePassword>
```

3. Crie um arquivo da truststore.

```
keytool -import -noprompt -alias <host>-file <cert_file> -
keystore <truststore_file> -storepass <truststorePassword>
```

4. Salve o certificado JKS em AWS Secrets Manager. Substitua livy-jks-secret pelo segredo e fileb://mykeystore.jks pelo caminho para o certificado do JKS do keystore.

```
aws secretsmanager create-secret \
--name livy-jks-secret \
--description "My Livy keystore JKS secret" \
--secret-binary fileb://mykeystore.jks
```

5. Salve o a senha do keystore e da chave no Secrets Manager. Use seus próprios parâmetros.

```
aws secretsmanager create-secret \
--name livy-jks-secret \
--description "My Livy key and keystore password secret" \
--secret-string "{\"keyPassword\":\"<test-key-password>\",\"keyStorePassword\":\"<test-key-store-password>\"}"
```

6. Crie um namespace de servidor do Livy com o comando a seguir.

```
kubectl create ns <livy-ns>
```

7. Crie o objeto ServiceProviderClass para o servidor do Livy que tem o certificado do JKS e as senhas.

```
cat >livy-secret-provider-class.yaml << EOF
apiVersion: secrets-store.csi.x-k8s.io/v1
kind: SecretProviderClass
metadata:
   name: aws-secrets
spec:
   provider: aws
   parameters:
    objects: |</pre>
```

```
- objectName: "livy-jks-secret"
    objectType: "secretsmanager"
- objectName: "livy-passwords"
    objectType: "secretsmanager"

EOF
kubectl apply -f livy-secret-provider-class.yaml -n livy-ns>
```

Conceitos básicos do Apache Livy habilitado por SSL

Depois de habilitar o SSL no servidor do Livy, você deve configurar a serviceAccount para ter acesso aos segredos de keyStore e keyPasswords no AWS Secrets Manager.

Crie o namespace do servidor do Livy.

```
kubectl create namespace <livy-ns>
```

 Configure a conta de serviço do Livy para ter acesso aos segredos no Secrets Manager. Para obter mais informações sobre como configurar o IRSA, consulte <u>Setting up IRSA while installing</u> Apache Livy.

```
aws ecr get-login-password \--region region-id | helm registry login \
--username AWS \
--password-stdin ECR-registry-account.dkr.ecr.region-id.amazonaws.com
```

3. Instale o Livy. Para o parâmetro de versão do chart do Helm, use o rótulo de lançamento do Amazon EMR, como 7.1.0. Você também deve substituir o ID da conta de registro e o ID da região do Amazon ECR pelos seus IDs. Você pode encontrar o ECR-registry-account valor correspondente para suas contas Região da AWS de registro do Amazon ECR por região.

```
helm install livy-app-name \
    oci://895885662937.dkr.ecr.region-id.amazonaws.com/livy \
    --version 7.10.0 \
    --namespace livy-namespace-name \
    --set image=<ECR-registry-account.dkr.ecr > .< region > .amazonaws.com/livy/
emr-7.10.0:latest \
    --set sparkNamespace=spark-namespace \
    --set ssl.enabled=true
    --set ssl.CertificateArn=livy-acm-certificate-arn
    --set ssl.secretProviderClassName=aws-secrets
    --set ssl.keyStoreObjectName=livy-jks-secret
```

```
--set ssl.keyPasswordsObjectName=livy-passwords
--create-namespace
```

4. Continue a partir da etapa 5 em Installing Apache Livy on Amazon EMR on EKS.

Configuração das permissões de aplicações do Apache Livy e Spark com o controle de acesso por perfil (RBAC)

Para implantar o Livy, o Amazon EMR no EKS cria uma conta e perfil de serviço do servidor e uma conta e perfil de serviço do Spark. Esses perfis devem ter as permissões do RBAC necessárias para concluir a configuração e executar as aplicações do Spark.

Permissões do RBAC para o perfil e a conta de serviço do servidor

O Amazon EMR no EKS cria a conta e o perfil de serviço do servidor do Livy para gerenciar sessões do Livy para trabalhos do Spark e rotear o tráfego entre a entrada e outros recursos.

O nome padrão para esta conta de serviço é emr-containers-sa-livy. É necessário ter as permissões a seguir.

```
rules:
- apiGroups:
  resources:
  - "namespaces"
  verbs:
  - "get"
- apiGroups:
  _ ""
  resources:
  - "serviceaccounts"
    "services"
    "configmaps"
    "events"
    "pods"
    "pods/log"
  verbs:
  - "get"
    "list"
    "watch"
    "describe"
    "create"
```

```
"edit"
   "delete"
   "deletecollection"
   "annotate"
   "patch"
   "label"
- apiGroups:
  _ ""
  resources:
  - "secrets"
  verbs:
  - "create"
    "patch"
    "delete"
    "watch"
- apiGroups:
  _ ""
  resources:
  - "persistentvolumeclaims"
  verbs:
  - "get"
    "list"
    "watch"
    "describe"
    "create"
    "edit"
    "delete"
    "annotate"
    "patch"
    "label"
```

Permissões do RBAC para a conta e o perfil de serviço do Spark

Um pod do driver do Spark precisa de uma conta de serviço do Kubernetes no mesmo namespace que o pod. Essa conta de serviço precisa de permissões para gerenciar os pods do executor e todos os recursos exigidos pelo pod do driver. A menos que a conta de serviço padrão no namespace tenha as permissões necessárias, o driver falhará e será encerrado. As permissões RBAC a seguir são necessárias.

```
rules:
- apiGroups:
- ""
    "batch"
```

```
"extensions"
  "apps"
resources:
- "configmaps"
  "serviceaccounts"
  "events"
  "pods"
  "pods/exec"
  "pods/log"
  "pods/portforward"
  "secrets"
  "services"
  "persistentvolumeclaims"
  "statefulsets"
verbs:
- "create"
  "delete"
  "get"
  "list"
  "patch"
  "update"
  "watch"
  "describe"
  "edit"
  "deletecollection"
  "patch"
  "label"
```

Configuração de permissões de acesso com perfis do IAM para contas de serviço (IRSA)

Por padrão, o driver e os executores do servidor Livy e do aplicativo Spark não têm acesso aos recursos. AWS A conta de serviço do servidor e a conta de serviço do Spark controlam o acesso aos AWS recursos do servidor Livy e dos pods do aplicativo Spark. Para conceder acesso, você precisa mapear as contas de serviço com uma função do IAM que tenha as AWS permissões necessárias.

Você pode configurar o mapeamento do IRSA antes de instalar o Apache Livy, durante a instalação ou depois de concluir a instalação.

Configuração do IRSA durante a instalação do Apache Livy (para conta de serviço do servidor)



Note

Esse mapeamento é compatível somente com a conta de serviço do servidor.

- Verifique se você concluiu a configuração do Apache Livy para Amazon EMR no EKS e está no meio da instalação do Apache Livy com o Amazon EMR no EKS.
- Crie um namespace do Kubernetes para o servidor do Livy. Neste exemplo, o nome do 2. namespace é livy-ns.
- Crie uma política do IAM que inclua as permissões Serviços da AWS para as quais você deseja 3. que seus pods acessem. O exemplo a seguir cria uma política do IAM para obter recursos do Amazon S3 para o ponto de entrada do Spark.

```
cat >my-policy.json <<EOF{
"Version": "2012-10-17",
    "Statement": [
"Effect": "Allow",
            "Action": "s3:GetObject",
            "Resource": "arn:aws:s3:::my-spark-entrypoint-bucket"
        }
    ]
}
E0F
aws iam create-policy --policy-name my-policy --policy-document file://my-
policy.json
```

Use o comando a seguir para definir seu Conta da AWS ID como uma variável.

```
account_id=$(aws sts get-caller-identity --query "Account" --output text)
```

5. Defina o provedor de identidade OpenID Connect (OIDC) do cluster como uma variável de ambiente.

```
oidc_provider=$(aws eks describe-cluster --name my-cluster --region $AWS_REGION --
query "cluster.identity.oidc.issuer" --output text | sed -e "s/^https:\/\//")
```

6. Defina variáveis para o namespace e o nome da conta de serviço. Use seus próprios valores.

```
export namespace=default
export service_account=my-service-account
```

7. Crie um arquivo de política de confiança com o comando a seguir. Se quiser conceder acesso ao perfil a todas as contas de serviço em um namespace, copie o comando a seguir e substitua StringEquals por StringLike e \$service_account por *.

```
cat >trust-relationship.json <<EOF</pre>
{
  "Version": "2012-10-17",
  "Statement": [
      "Effect": "Allow",
      "Principal": {
        "Federated": "arn:aws:iam::$account_id:oidc-provider/$oidc_provider"
      },
      "Action": "sts:AssumeRoleWithWebIdentity",
      "Condition": {
        "StringEquals": {
          "$oidc_provider:aud": "sts.amazonaws.com",
          "$oidc_provider:sub": "system:serviceaccount:$namespace:$service_account"
        }
      }
    }
  ]
}
E0F
```

8. Crie a função.

```
aws iam create-role --role-name my-role --assume-role-policy-document file://trust-relationship.json --description "my-role-description"
```

9. Use o comando de instalação do Helm a seguir para configurar o serviceAccount.executionRoleArn para mapear o IRSA. Confira a seguir um exemplo do comando de instalação do Helm. Você pode encontrar o ECR-registry-account valor correspondente para suas contas Região da AWS de registro do Amazon ECR por região.

```
helm install livy-demo \
oci://895885662937.dkr.ecr.us-west-2.amazonaws.com/livy \
```

```
--version 7.10.0 \
--namespace livy-ns \
--set image=ECR-registry-account.dkr.ecr.region-id.amazonaws.com/livy/
emr-7.10.0:latest \
--set sparkNamespace=spark-ns \
--set serviceAccount.executionRoleArn=arn:aws:iam::123456789012:role/my-role
```

Mapeamento do IRSA para uma conta de serviço do Spark

Antes de mapear o IRSA para uma conta de serviço do Spark, conclua os seguintes itens:

- Verifique se você concluiu a configuração do Apache Livy para Amazon EMR no EKS e está no meio da instalação do Apache Livy com o Amazon EMR no EKS.
- Você deve ter um provedor OpenID Connect (OIDC) existente do IAM para o cluster. Para verificar se você já tem um provedor ou saber como criar um, consulte <u>Criar um provedor OIDC do IAM</u> para o cluster.
- Certifique-se de ter instalado a versão 0.171.0 ou posterior da CLI do eksctl ou o AWS CloudShell. Para instalar ou atualizar o eksctl, consulte <u>Instalação</u> na documentação do eksctl.

Siga estas etapas para mapear o IRSA à conta de serviço do Spark:

1. Use o comando a seguir para obter a conta de serviço do Spark.

```
SPARK_NAMESPACE=<spark-ns>
LIVY_APP_NAME=<livy-app-name>
kubectl --namespace $SPARK_NAMESPACE describe sa -1 "app.kubernetes.io/instance=
$LIVY_APP_NAME" | awk '/^Name:/ {print $2}'
```

2. Defina variáveis para o namespace e o nome da conta de serviço.

```
export namespace=default
export service_account=my-service-account
```

3. Use o comando a seguir para criar um arquivo de política de confiança para o perfil do IAM. O exemplo a seguir dá permissão a todas as contas de serviço dentro do namespace para usar o perfil. Para fazer isso, substitua StringEquals por StringLike e \$service_account por *.

```
cat >trust-relationship.json <<EOF</pre>
```

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Principal": {
        "Federated": "arn:aws:iam::$account_id:oidc-provider/$oidc_provider"
      },
      "Action": "sts:AssumeRoleWithWebIdentity",
      "Condition": {
        "StringEquals": {
          "$oidc_provider:aud": "sts.amazonaws.com",
          "$oidc_provider:sub": "system:serviceaccount:$namespace:$service_account"
        }
      }
    }
  ]
}
E0F
```

4. Crie a função.

```
aws iam create-role --role-name my-role --assume-role-policy-document file://trust-relationship.json --description "my-role-description"
```

5. Mapeie o servidor ou a conta de serviço do Spark com o comando eksctl a seguir. Use seus próprios valores.

```
eksctl create iamserviceaccount --name spark-sa \
--namespace spark-namespace --cluster livy-eks-cluster \
--attach-role-arn arn:aws:iam::0123456789012:role/my-role \
--approve --override-existing-serviceaccounts
```

Propriedades de instalação do Apache Livy em versões do Amazon EMR no EKS

A instalação do Apache Livy permite que você selecione uma versão do chart do Helm no Livy. O chart do Helm oferece uma variedade de propriedades para personalizar sua experiência de instalação e configuração. Essas propriedades são compatíveis com o Amazon EMR no EKS nas versões 7.1.0 e superiores.

Tópicos

• Propriedades de instalação do Amazon EMR 7.1.0

Propriedades de instalação do Amazon EMR 7.1.0

Propriedade	Descrição	Padrão
image	O URI de lançamento do Amazon EMR do servidor do Livy. Essa configuração é necessária.	1111
sparkNamespace	Namespace para executar sessões do Livy do Spark. Por exemplo, especifique "livy". Essa configuração é necessári a.	***************************************
nameOverride	Forneça um nome em vez de livy. O nome é definido como um rótulo para todos os recursos do Livy	"livy"
fullnameOverride	Forneça um nome para usar em vez dos nomes completos dos recursos.	7111
ssl.enabled	Habilita o end-to-end SSL do endpoint Livy para o servidor Livy.	FALSE
ssl.certificateArn	Se o SSL estiver habilitado, esse será o ARN do certifica	1111

Propriedade	Descrição	Padrão
	do do ACM para o NLB criado pelo serviço.	
ssl. secretProviderClassNome	Se o SSL estiver habilitad o, esse é o nome da classe do provedor secreto para proteger o NLB na conexão do servidor do Livy com SSL.	1111
ssl. keyStoreObjectNome	Se o SSL estiver habilitado, o nome do objeto do certifica do do keystore na classe do provedor secreto.	***************************************
ssl. keyPasswordsObject Nome	Se o SSL estiver habilitado, o nome do objeto do segredo que tem a senha do keystore e de chaves.	HH
rbac.create	Se verdadeiro, cria recursos do RBAC.	FALSE
serviceAccount.create	Se verdadeiro, cria uma conta de serviço do Livy.	VERDADEIRO
serviceAccount.name	O nome da conta de serviço a ser usada no Livy. Se você não definir essa propriedade e criar uma conta de serviço, o Amazon EMR no EKS gerará automaticamente um nome usando a propriedade de substituição fullname.	"emr-containers-sa-livy"
Conta de serviço. execution RoleArn	O ARN do perfil de execução da conta de serviço do Livy.	1111

Propriedade	Descrição	Padrão
sparkServiceAccount.criar	SE verdadeiro, cria a conta de serviço do Spark no .Release.Namespace	VERDADEIRO
sparkServiceAccount.name	O nome da conta de serviço que será usada no Spark. Se você não definir essa propriedade e criar uma conta de serviço do Spark, o Amazon EMR no EKS gerará automaticamente um nome com a propriedade fullname0verride e o sufixo -spark-livy.	"emr-containers-sa-spark-lívio"
service.name	O nome do serviço do Livy	<pre>"emr-containers-li vy"</pre>
service.annotations	Anotações do serviço do Livy	{}
loadbalancer.enabled	Se um balanceador de carga deve ser criado para o serviço Livy usado para expor o endpoint do Livy fora do cluster do Amazon EKS.	FALSE

Propriedade	Descrição	Padrão
loadbalancer.internal	Se o endpoint do Livy deve ser configurado como interno à VPC ou externo.	FALSE
	Definir essa proprieda de para FALSE expõe o endpoint a fontes fora da VPC. Recomendamos proteger seu endpoint com TLS/SSL. Para obter mais informações, consulte Setting up TLS and SSL encryption.	
imagePullSecrets	A lista de nomes imagePull Secret a serem usados para extrair a imagem do Livy de repositórios privados.	
recursos	As solicitações de recursos e os limites dos contêineres do Livy.	{}
nodeSelector	Os nós para os quais programar os pods do Livy.	0
tolerations	Uma lista contendo as tolerâncias dos pods do Livy a serem definidas.	
affinity	As regras de afinidade dos pods do Livy.	0
persistence.enabled	Se verdadeiro, habilita a persistência dos diretórios de sessões.	FALSE

Propriedade	Descrição	Padrão
persistence.subPath	O subcaminho da VPC a ser montado nos diretórios de sessões.	1111
persistence.existingClaim	A VPC a ser usada em vez de criar uma nova.	0
persistence.storageClass	A classe de armazenamento a ser usada. Para definir esse parâmetro, use o formato storageClassName: <storageclassname: <storageclass=""> . Definir esse parâmetro para "-" desabilita o provisionamento dinâmico. Se você definir esse parâmetro como nulo ou não especificar nada, o Amazon EMR no EKS não definirá a e usará o storageClassName provisionador padrão.</storageclassname:>	
persistence.accessMode	O modo de acesso da PVC.	ReadWriteOnce
persistence.size	O tamanho da VPC.	20 Gi
persistence.annotations	Anotações adicionais para a VPC.	0
env.*	Ambientes adicionais a serem configurados no contêiner do Livy. Para obter mais informações, consulte Inputting your own Livy and Spark configurations while installing Livy.	{

Propriedade	Descrição	Padrão
envFrom.*	Ambientes adicionais a serem definidos para o Livy usando um mapa de configuração ou segredo do Kubernetes.	
livyConf.*	Entradas adicionais do livy.conf a serem definidas usando um mapa de configura ção ou segredo montado do Kubernetes.	{
sparkDefaultsConf.*	Entradas adicionais de spark-defaults.conf a serem definidas usando um mapa de configuração ou segredo montado do Kubernetes.	{}

Solução de erros comuns de formato de variável de ambiente

Quando você insere as configurações do Livy e do Spark, há formatos de variáveis de ambiente que não são compatíveis e podem causar erros. O procedimento guia você por uma série de etapas que ajudam a garantir o uso de formatos corretos.

Como inserir suas próprias configurações do Livy e do Spark ao instalar o Livy

Você pode configurar qualquer variável de ambiente do Apache Livy ou do Apache Spark com a propriedade env.* do Helm. Siga as etapas abaixo para converter o exemplo de configuração example.config.with-dash.withUppercase em um formato de variável de ambiente compatível.

 Substitua as letras maiúsculas por 1 e a forma minúscula da letra. Por exemplo, example.config.with-dash.withUppercase torna-se example.config.withdash.with1uppercase.

- 2. Substitua traços (-) por 0. Por exemplo, example.config.with-dash.with1uppercase se torna example.config.with0dash.with1uppercase
- Substitua pontos (.) por sublinhados (_). Por exemplo, example.config.with0dash.with1uppercase torna-se example_config_with0dash_with1uppercase.
- 4. Substitua todas as letras minúsculas por maiúsculas.
- 5. Adicione o prefixo LIVY_ ao nome da variável.
- 6. Use a variável ao instalar o Livy por meio do gráfico de comando usando o formato --set env. YOUR_VARIABLE_NAME.valor= yourvalue

Por exemplo, para definir as configurações livy.server.recovery.state-store = filesystem e spark.kubernetes.executor.podNamePrefix = my-prefix do Livy e do Spark, use estas propriedades do Helm:

```
-set env.LIVY_LIVY_SERVER_RECOVERY_STATE0STORE.value=filesystem
-set env.LIVY_SPARK_KUBERNETES_EXECUTOR_POD0NAME0PREFIX.value=myprefix
```

Gerenciamento de execuções de trabalhos do Amazon EMR no EKS

As seções a seguir abordam tópicos que ajudam você a gerenciar suas execuções de trabalhos do Amazon EMR no EKS. Isso inclui configurar os parâmetros de execução do trabalho ao usar o AWS CLI, configurar como seus dados de registro são armazenados, executar scripts do Spark SQL para executar consultas, entender os estados de execução do trabalho e saber como monitorar trabalhos. Você pode trabalhar com esses tópicos, geralmente em ordem, se quiser configurar e concluir uma execução de trabalho para processar dados.

Tópicos

- Gerenciando execuções de trabalhos com o AWS CLI
- Executando scripts do Spark SQL por meio da API StartJobRun
- Estados de execução de trabalho
- Visualização de trabalhos no console do Amazon EMR
- Erros comuns ao executar trabalhos

Gerenciando execuções de trabalhos com o AWS CLI

Este tópico aborda como gerenciar execuções de trabalhos com o AWS Command Line Interface (AWS CLI). São detalhadas as propriedades, como parâmetros de segurança, o driver e várias configurações de substituição. Também inclui subtópicos que abrangem várias maneiras de configurar o registro em log.

Tópicos

- Opções para a configuração de uma execução de trabalho
- Configuração de uma execução de trabalho para usar logs do Amazon S3
- Configurar uma execução de trabalho para usar o Amazon CloudWatch Logs
- Listagem de execuções de trabalhos
- Descrição de uma execução de trabalho
- Cancelamento de uma execução de trabalho

Opções para a configuração de uma execução de trabalho

Use as seguintes opções para configurar os parâmetros de execução de trabalho:

- --execution-role-arn: você deve fornecer um perfil do IAM que é usado para a execução de trabalhos. Para obter mais informações, consulte <u>Uso de perfis de execução de trabalho com o</u> <u>Amazon EMR no EKS</u>.
- --release-label: você pode implantar o Amazon EMR no EKS com as versões 5.32.0 e
 6.2.0 e posteriores do Amazon EMR. O Amazon EMR no EKS não é compatível com as versões anteriores do Amazon EMR. Para obter mais informações, consulte Versões do Amazon EMR no EKS.
- --job-driver: o driver de trabalho usado para fornecer entradas sobre o trabalho principal. Este é um campo do tipo união no qual você só pode transferir um dos valores para o tipo de trabalho que deseja executar. Os tipos de trabalho com suporte incluem:
 - Trabalhos do Spark-submit: usados para executar um comando por meio do spark-submit. Você
 pode usar esse tipo de tarefa para executar Scala, SparkR PySpark, SparkSQL e qualquer outra
 tarefa compatível por meio do Spark Submit. Esse tipo de trabalho tem os seguintes parâmetros:
 - Ponto de entrada Essa é a referência HCFS (sistema de arquivos compatível com Hadoop)
 ao jar/py arquivo principal que você deseja executar.

- EntryPointArguments Essa é uma matriz de argumentos que você deseja passar para o jar/py arquivo principal. Você deve realizar a leitura desses parâmetros usando seu código de Entrypoint. Cada argumento na matriz deve ser separado por uma vírgula. EntryPointArguments não pode conter colchetes ou parênteses, como (), {} ou [].
- SparkSubmitParameters Esses são os parâmetros adicionais do Spark que você deseja enviar para o trabalho. Use este parâmetro para substituir as propriedades padrão do Spark, como a memória do driver ou o número de executores, como —conf ou —class. Para obter informações adicionais, consulte Launching Applications with spark-submit.
- Trabalhos do Spark SQL: usados para executar um arquivo de consulta SQL por meio do Spark SQL. Você pode usar esse tipo de trabalho para executar trabalhos do Spark SQL. Esse tipo de trabalho tem os seguintes parâmetros:
 - Entrypoint: esta é a referência do HCFS (sistema de arquivos compatível com Hadoop) para o arquivo de consulta SQL que você deseja executar.
 - Para obter uma lista de parâmetros adicionais do Spark que você pode usar para um trabalho do Spark SQL, consulte Executando scripts do Spark SQL por meio da API StartJobRun .
- --configuration-overrides: você pode substituir as configurações padrão para as aplicações ao fornecer um objeto de configuração. Você pode usar uma sintaxe abreviada para fornecer a configuração ou fazer referência ao objeto de configuração em um arquivo JSON. Os objetos de configuração consistem em uma classificação, propriedades e configurações opcionais aninhadas As propriedades consistem nas configurações que você deseja substituir neste arquivo. Você pode especificar várias classificações para diversas aplicações em um único objeto JSON. As classificações de configuração disponíveis variam de acordo com a versão do Amazon EMR. Para obter uma lista das classificações de configuração disponíveis para cada versão de liberação do Amazon EMR, consulte Versões do Amazon EMR no EKS.

Se você transferir a mesma configuração em uma substituição de aplicação e nos parâmetros do spark-submit, os parâmetros do spark-submit terão precedência. A lista completa de prioridades de configuração é apresentada a seguir, na ordem da prioridade mais alta para a prioridade mais baixa.

- Configuração fornecida ao criar SparkSession.
- Configuração fornecida como parte do sparkSubmitParameters usando —conf.
- Configuração fornecida como parte das substituições de aplicações.
- Configurações otimizadas escolhidas pelo Amazon EMR para a versão.
- Configurações padrão de código aberto para a aplicação.

Para monitorar a execução de trabalhos usando o Amazon CloudWatch ou o Amazon S3, você deve fornecer os detalhes da configuração do. CloudWatch Para obter mais informações, consulte Configuração de uma execução de trabalho para usar logs do Amazon S3 e Configurar uma execução de trabalho para usar o Amazon CloudWatch Logs. Se o bucket ou grupo de CloudWatch logs do S3 não existir, o Amazon EMR o criará antes de fazer o upload dos registros para o bucket.

 Para obter uma lista adicional de opções de configuração do Kubernetes, consulte Spark Properties on Kubernetes.

As configurações do Spark apresentadas a seguir não têm suporte.

- spark.kubernetes.authenticate.driver.serviceAccountName
- spark.kubernetes.authenticate.executor.serviceAccountName
- spark.kubernetes.namespace
- spark.kubernetes.driver.pod.name
- spark.kubernetes.container.image.pullPolicy
- spark.kubernetes.container.image



Note

Você pode usar spark.kubernetes.container.image para imagens do Docker personalizadas. Para obter mais informações, consulte Personalização de imagens do Docker para o Amazon EMR no EKS.

Configuração de uma execução de trabalho para usar logs do Amazon S3

Para poder monitorar o progresso do trabalho e solucionar falhas, você deve configurar seus trabalhos para enviar informações de log para o Amazon S3, o CloudWatch Amazon Logs ou ambos. Este tópico ajuda você a começar a publicar logs de aplicações no Amazon S3 em trabalhos iniciados com o Amazon EMR no EKS.

Política do IAM para logs do S3

Antes que os trabalhos possam enviar dados de log ao Amazon S3, as permissões apresentadas a seguir devem ser inclusas na política de permissões para o perfil de execução do trabalho. Substitua amzn-s3-demo-logging-bucket pelo nome do bucket de registro em log.

JSON

```
"Version": "2012-10-17",
  "Statement": [
      "Effect": "Allow",
      "Action": [
        "s3:PutObject",
        "s3:GetObject",
        "s3:ListBucket"
      ],
      "Resource": [
        "arn:aws:s3:::amzn-s3-demo-bucket",
        "arn:aws:s3:::amzn-s3-demo-bucket/*"
      ],
      "Sid": "AllowS3Putobject"
    }
  ]
}
```

Note

O Amazon EMR no EKS também pode criar um bucket do Amazon S3. Se um bucket do Amazon S3 não estiver disponível, inclua a permissão "s3:CreateBucket" na política do IAM.

Depois de conceder ao perfil de execução as permissões adequadas para enviar logs ao Amazon S3, os dados de log serão enviados para os locais do Amazon S3 apresentados a seguir quando s3MonitoringConfiguration for transmitida na seção monitoringConfiguration de uma solicitação start-job-run, conforme mostrado em Gerenciando execuções de trabalhos com o AWS CLI.

- Registros do remetente -//virtual-cluster-id/jobs/ /containers/logUri/(stderr.gz/stdout.gz)
 job-id pod-name
- Registros do driver -//virtual-cluster-id/jobs/ logUri job-id /containers/ /spark--driver/ (spark-application-idstderr.gz/stdout.gz) job-id

Registros do executor -//virtual-cluster-id/jobs/ logUri job-id /containers///(spark-application-idstderr.gz/stdout.gz) executor-pod-name

Configurar uma execução de trabalho para usar o Amazon CloudWatch Logs

Para monitorar o progresso do trabalho e solucionar falhas, você deve configurar seus trabalhos para enviar informações de log para o Amazon S3, CloudWatch Amazon Logs ou ambos. Este tópico ajuda você a começar a usar CloudWatch registros em seus trabalhos que são lançados com o Amazon EMR no EKS. Para obter mais informações sobre CloudWatch registros, consulte Monitoramento de arquivos de log no Guia CloudWatch do usuário da Amazon.

CloudWatch Política de IAM de registros

Para que seus trabalhos enviem dados de registro para o CloudWatch Logs, as seguintes permissões devem ser incluídas na política de permissões da função de execução do trabalho. Substitua $my_log_group_name$ e $my_log_stream_prefix$ por nomes de seu grupo de CloudWatch registros e nomes de fluxo de registros, respectivamente. O Amazon EMR no EKS cria o grupo de logs e o fluxo de logs se eles não existirem, desde que o ARN do perfil de execução tenha as permissões apropriadas.

JSON

```
"logs:PutLogEvents"
],
    "Resource": [
        "arn:aws:logs:*:*:log-group:my_log_group_name:log-
stream:my_log_stream_prefix/*"
        ],
        "Sid": "AllowLOGSPutlogevents"
    }
]
```

Note

O Amazon EMR no EKS também pode criar um fluxo de logs. Se um fluxo de logs não existir, a política do IAM deverá incluir a permissão "logs:CreateLogGroup".

Depois de conceder à sua função de execução as permissões adequadas, seu aplicativo envia os dados de registro para o CloudWatch cloudWatchMonitoringConfiguration Logs quando são passados na monitoringConfiguration seção de uma start-job-run solicitação, conforme mostrado emGerenciando execuções de trabalhos com o AWS CLI.

Na StartJobRun API, *log_group_name* é o nome do grupo de registros para CloudWatch e *log_stream_prefix* é o prefixo do nome do fluxo de registros para CloudWatch. Você pode visualizar e pesquisar esses logs no AWS Management Console.

- Registros do remetente -logGroup//virtual-cluster-id/jobs/ logStreamPrefix / containers//(stderr/stdoutjob-id) pod-name
- Registros de drivers -logGroup//logStreamPrefixvirtual-cluster-id/jobs/ job-id / containers//spark--driver/(spark-application-idstderrstdout) job-id
- Registros do executor -logGroup//logStreamPrefixvirtual-cluster-id/jobs/ job-id / containers///(spark-application-idstderr/stdout) executor-pod-name

Listagem de execuções de trabalhos

Você pode executar list-job-run para mostrar os estados das execuções de trabalhos, como demonstra o exemplo a seguir.

aws emr-containers list-job-runs --virtual-cluster-id <cluster-id>

Descrição de uma execução de trabalho

Você pode executar describe-job-run para obter mais detalhes sobre o trabalho, como o estado do trabalho, os detalhes do estado e o nome do trabalho, como demonstra o exemplo a seguir.

aws emr-containers describe-job-run --virtual-cluster-id cluster-id --id job-run-id

Cancelamento de uma execução de trabalho

Você pode executar cancel-job-run para cancelar trabalhos em execução, como demonstra o exemplo a seguir.

aws emr-containers cancel-job-run --virtual-cluster-id cluster-id --id job-run-id

Executando scripts do Spark SQL por meio da API StartJobRun

As versões 6.7.0 e posteriores do Amazon EMR no EKS incluem um driver de trabalho do Spark SQL para que você possa executar scripts do Spark SQL por meio da API StartJobRun. Você pode fornecer arquivos de ponto de entrada SQL para executar consultas do Spark SQL diretamente no Amazon EMR no EKS com a API StartJobRun, sem a necessidade de realizar modificações nos scripts do Spark SQL existentes. A tabela a seguir lista os parâmetros do Spark que são compatíveis com as tarefas do Spark SQL por meio da StartJobRun API.

Você pode escolher entre os parâmetros do Spark apresentados a seguir para enviar para um trabalho do Spark SQL. Use esses parâmetros para substituir as propriedades padrão do Spark.

Opção	Descrição
name NAME	Nome da aplicação
jars JARS	Lista separada por vírgulas de arquivos em JARs a serem inclusos no driver e no caminho de classe de execução.
packages	Lista separada por vírgulas de coordenad as do Maven de arquivos em JARs a serem

Descrição
inclusas nos caminhos de classe do driver e do executor.
Lista separada por vírgulas de groupld:a rtifactId, para excluir ao resolver as dependênc ias fornecidas empackages para evitar conflitos de dependência.
Lista separada por vírgulas de repositór ios remotos adicionais para pesquisar as coordenadas do Maven fornecidas com packages.
Lista separada por vírgulas de arquivos a serem colocados no diretório de trabalho de cada executor.
Propriedade de configuração do Spark.
Caminho para um arquivo do qual as proprieda des extras serão carregadas.
Memória para o driver. O padrão é de 1.024 MB.
Opções extras do Java a serem transferidas para o driver.
Entradas extras de caminhos da biblioteca a serem transferidas para o driver.
Entradas extras de caminhos de classe a serem transferidas para o driver.
Memória por executor. O padrão é de 1 GB.
Número de núcleos usados pelo driver.

Opção	Descrição
total-executor-cores NÚMERO	Total de núcleos para todos os executores.
executor-cores NUM	Número de núcleos usados por cada executor.
num-executors NUM	Número de executores para iniciar.
-hivevar <key=value></key=value>	Substituição de variável para a aplicação dos comandos do Hive, por exemplo, -hivevar A=B.
-hiveconf <pre><pre><pre><pre>-hiveconf <pre><pre><pre>property=value></pre></pre></pre></pre></pre></pre></pre>	Valor a ser usado para a propriedade em questão.

Para uma tarefa do Spark SQL, crie um start-job-run-request arquivo.json e especifique os parâmetros necessários para a execução da tarefa, como no exemplo a seguir:

```
{
  "name": "myjob",
  "virtualClusterId": "123456",
  "executionRoleArn": "iam_role_name_for_job_execution",
  "releaseLabel": "emr-6.7.0-latest",
  "jobDriver": {
    "sparkSqlJobDriver": {
      "entryPoint": "entryPoint_location",
       "sparkSqlParameters": "--conf spark.executor.instances=2 --conf
 spark.executor.memory=2G --conf spark.executor.cores=2 --conf spark.driver.cores=1"
    }
  },
  "configurationOverrides": {
    "applicationConfiguration": [
        "classification": "spark-defaults",
        "properties": {
          "spark.driver.memory":"2G"
      }
    "monitoringConfiguration": {
      "persistentAppUI": "ENABLED",
```

```
"cloudWatchMonitoringConfiguration": {
    "logGroupName": "my_log_group",
    "logStreamNamePrefix": "log_stream_prefix"
    },
    "s3MonitoringConfiguration": {
        "logUri": "s3://my_s3_log_location"
    }
    }
}
```

Estados de execução de trabalho

Quando você envia uma execução de trabalho para uma fila de trabalhos do Amazon EMR no EKS, a execução de trabalho entra no estado PENDING. Em seguida, ela passa pelos estados apresentados a seguir até obter êxito (sair com o código 0) ou falhar (sair com um código diferente de zero).

As execuções de trabalhos podem ter os seguintes estados:

- PENDING: o estado inicial do trabalho quando a execução de trabalho é enviada ao Amazon EMR no EKS. O trabalho está aguardando para ser enviado ao cluster virtual, e o Amazon EMR no EKS está trabalhando no envio desse trabalho.
- SUBMITTED: uma execução de trabalho que foi enviada com êxito ao cluster virtual. O programador do cluster tenta executar esse trabalho no cluster.
- RUNNING: uma execução de trabalho em execução no cluster virtual. Em aplicações do Spark, isso significa que o processo do driver do Spark está no estado running.
- FAILED: uma execução de trabalho que falhou ao ser enviada ao cluster virtual ou que foi concluída sem êxito. Examine StateDetails e FailureReason encontre informações adicionais sobre essa falha no trabalho.
- COMPLETED: uma execução de trabalho que foi concluída com êxito.
- CANCEL_PENDING: uma execução de trabalho foi solicitada para cancelamento. O Amazon EMR no EKS está tentando cancelar o trabalho no cluster virtual.
- CANCELLED: uma execução de trabalho que foi cancelada com êxito.

Visualização de trabalhos no console do Amazon EMR

Os dados de execução do trabalho estão disponíveis para exibição, para que você possa monitorar cada trabalho à medida que ele passa pelos estados. Para visualizar trabalhos no console do Amazon EMR, execute as etapas a seguir.

- 1. No menu à esquerda do console do Amazon EMR, em Amazon EMR no EKS, escolha Clusters virtuais.
- 2. Na lista de clusters virtuais, selecione o cluster virtual do qual deseja visualizar os trabalhos.
- 3. Na tabela Execuções de trabalhos, selecione Visualizar logs para visualizar os detalhes de uma execução de trabalho.

Note

O suporte para a experiência de um clique está habilitado por padrão. Ele pode ser desativado ao configurar persistentAppUI para DISABLED na monitoringConfiguration durante o envio do trabalho. Para obter mais informações, consulte View Persistent Application User Interfaces.

Erros comuns ao executar trabalhos

Os erros a seguir podem ocorrer ao executar a API StartJobRun. A tabela lista cada erro e fornece etapas de mitigação para que você possa resolver os problemas rapidamente.

Mensagem de erro	Condição de erro	Próxima etapa recomendada
erro: argumento <i>argument</i> é obrigatório	Os parâmetros obrigatórios estão ausentes.	Adicione os argumentos ausentes à solicitação de API.
Ocorreu um erro (AccessDe niedException) ao chamar a StartJobRun operação: Usuário: não <i>ARN</i> está autorizado a realizar: emrcontainers: StartJobRun	O perfil de execução está ausente.	Consulte Uso de perfis de execução de trabalho com o Amazon EMR no EKS.

Mensagem de erro	Condição de erro	Próxima etapa recomendada
Ocorreu um erro (AccessDe niedException) ao chamar a StartJobRun operação: Usuário: não <i>ARN</i> está autorizado a realizar: emrcontainers: StartJobRun	O chamador não tem permissão para acessar o perfil de execução [formato válido ou inválido] por meio de chaves de condição.	Consulte Uso de perfis de execução de trabalho com o Amazon EMR no EKS.
Ocorreu um erro (AccessDe niedException) ao chamar a StartJobRun operação: Usuário: não <i>ARN</i> está autorizado a realizar: emrcontainers: StartJobRun	O emissor de trabalho e o ARN do perfil de execução são de contas diferentes.	Certifique-se de que o ARN do emissor de trabalho e do perfil de execução sejam da mesma conta da AWS .
1 erro de validação detectado : o valor <i>Role</i> em 'executio nRoleArn' falhou em satisfazer o padrão de expressão regular do ARN: ^arn :(aws [a-zA-z0-9-] *) :iam: :(\ d {12})? : (papel ((\ u002F) (\ u002F [\ u0021-\ u007F] +\ u002F)) [\ w+=, .@-] +)	O chamador tem permissõe s para acessar o perfil de execução por meio de chaves de condição, mas o perfil não atende às restrições do formato de ARN.	Forneça o perfil de execução seguindo o formato de ARN. Consulte <u>Uso de perfis de execução de trabalho com o Amazon EMR no EKS.</u>
Ocorreu um erro (Resource NotFoundException) ao chamar a StartJobRun operação: O cluster virtual <i>Virtual Cluster ID</i> não existe.	O ID do cluster virtual não foi encontrado.	Forneça um ID do cluster virtual registrado no Amazon EMR no EKS.

Mensagem de erro	Condição de erro	Próxima etapa recomendada
Ocorreu um erro (Validati onException) ao chamar a StartJobRun operação: O estado do cluster virtual não <i>state</i> é válido para criar recursos JobRun.	O cluster virtual não está pronto para executar o trabalho.	Consulte Estados de um cluster virtual.
Ocorreu um erro (Resource NotFoundException) ao chamar a StartJobRun operação: A liberação RELEASE não existe.	A versão especificada no envio do trabalho está incorreta.	Consulte <u>Versões do Amazon</u> <u>EMR no EKS</u> .
Ocorreu um erro (AccessDe niedException) ao chamar a StartJobRun operação: Usuário: não ARN está autorizado a realizar: emrcontainers: StartJobRun on resource: ARN com uma negação explícita. Ocorreu um erro (AccessDe niedException) ao chamar a StartJobRun operação: Usuário: não ARN está autorizado a executar: emrcontainers: StartJobRun on resource: ARN	O usuário não está autorizado a ligar StartJobRun.	Consulte Uso de perfis de execução de trabalho com o Amazon EMR no EKS.

Mensagem de erro	Condição de erro	Próxima etapa recomendada
Ocorreu um erro (Validati onException) ao chamar a StartJobRun operação: configurationOverrides.moni toringConfiguration.s3 MonitoringConfigur ation .logURI falhou em satisfazer a restrição: %s	A sintaxe do URI do caminho do S3 não é válida.	O logUri deve estar no formato s3://

Os erros a seguir podem ocorrer ao executar a API DescribeJobRun antes da execução de trabalho.

Mensagem de erro	Condição de erro	Próxima etapa recomendada
StateDetails: falha no JobRun envio.	Os parâmetros em não StartJobRun são válidos.	Consulte <u>Versões do Amazon</u> <u>EMR no EKS</u> .
Classificação <i>classific ation</i> não suportada.		
failureReason: VALIDATIO N_ERROR		
state: FAILED.		
StateDetails: <i>EKS Cluster ID</i> O cluster não existe. failureReason: CLUSTER_U NAVAILABLE state: FAILED	O cluster do EKS não está disponível.	Verifique se o cluster do EKS existe e tem as permissõe s corretas. Para obter mais informações, consulte Configuração do Amazon EMR no EKS.
StateDetails: <i>EKS Cluster ID</i> O cluster não tem permissões suficientes.	O Amazon EMR não tem permissões para acessar o cluster do EKS.	Verifique se as permissõe s estão configuradas para o Amazon EMR no namespace

Mensagem de erro	Condição de erro	Próxima etapa recomendada
failureReason: CLUSTER_U NAVAILABLE state: FAILED		registrado. Para obter mais informações, consulte Configuração do Amazon EMR no EKS.
StateDetails: O cluster não EKS Cluster ID está acessível no momento. failureReason: CLUSTER_U NAVAILABLE state: FAILED	O cluster do EKS não está acessível.	Verifique se o cluster do EKS existe e tem as permissõe s corretas. Para obter mais informações, consulte Configuração do Amazon EMR no EKS.
StateDetails: falha no JobRun envio devido a um erro interno. failureReason: INTERNAL_ ERROR state: FAILED	Ocorreu um erro interno com o cluster do EKS.	N/D
StateDetails: <i>EKS Cluster ID</i> O cluster não tem recursos suficientes. failureReason: USER_ERROR state: FAILED	Não há recursos suficient es no cluster do EKS para executar o trabalho.	Adicione mais capacidad e ao grupo de nós do EKS ou configure o Autoscale r do EKS. Para obter mais informações, consulte Autoscaler do cluster.

Os erros a seguir podem ocorrer ao executar a API DescribeJobRun depois da execução de trabalho.

Mensagem de erro	Condição de erro	Próxima etapa recomendada
StateDetails: Problemas para monitorar seu. JobRun EKS Cluster IDO cluster não existe. failureReason: CLUSTER_U NAVAILABLE state: FAILED	O cluster do EKS não existe.	Verifique se o cluster do EKS existe e tem as permissõe s corretas. Para obter mais informações, consulte Configuração do Amazon EMR no EKS.
StateDetails: Problemas para monitorar seu. JobRun EKS Cluster IDO cluster não tem permissões suficient es. failureReason: CLUSTER_U NAVAILABLE state: FAILED	O Amazon EMR não tem permissões para acessar o cluster do EKS.	Verifique se as permissõe s estão configuradas para o Amazon EMR no namespace registrado. Para obter mais informações, consulte Configuração do Amazon EMR no EKS.
StateDetails: Problemas para monitorar seu. JobRun No momento, <i>EKS Cluster ID</i> o cluster não está acessível. failureReason: CLUSTER_U NAVAILABLE state: FAILED	O cluster do EKS não está acessível.	Verifique se o cluster do EKS existe e tem as permissõe s corretas. Para obter mais informações, consulte Configuração do Amazon EMR no EKS.
StateDetails: Problemas para monitorar seu JobRun devido a um erro interno	Ocorreu um erro interno que está impedindo o JobRun monitoramento.	N/D

Mensagem de erro	Condição de erro	Próxima etapa recomendada
failureReason: INTERNAL_ ERROR		
state: FAILED		

O erro apresentado a seguir pode ocorrer quando um trabalho não pode ser iniciado e aguarda no estado SUBMITTED por 15 minutos. Isso pode ser causado pela falta de recursos do cluster.

Mensagem de erro	Condição de erro	Próxima etapa recomendada
cluster timeout	O trabalho está no estado SUBMITTED há 15 minutos ou mais.	Você pode substituir a configuração padrão de 15 minutos para este parâmetro com a substituição de configuração mostrada abaixo.

Use a configuração a seguir para alterar a configuração de tempo limite do cluster para 30 minutos. Observe que você fornece o novo valor para job-start-timeout em segundos:

```
{
"configurationOverrides": {
    "applicationConfiguration": [{
        "classification": "emr-containers-defaults",
        "properties": {
            "job-start-timeout":"1800"
        }
    }]
}
```

Uso de modelos de trabalho

Um modelo de trabalho armazena valores que podem ser compartilhados entre invocações da API StartJobRun ao iniciar uma execução de trabalho. Ele oferece suporte a dois casos de uso:

• Para evitar valores de solicitação da API StartJobRun repetitivos e recorrentes.

Uso de modelos de trabalho 211

 Para impor uma regra, determinados valores devem ser fornecidos por meio de solicitações da API StartJobRun.

Os modelos de trabalho possibilitam definir um modelo reutilizável para execuções de trabalho com a finalidade de aplicar personalizações adicionais, por exemplo:

- Configuração da capacidade de computação dos executores e dos drivers.
- Definição de propriedades de segurança e de governança, como perfis do IAM.
- Personalização de uma imagem do Docker para usar em diversas aplicações e pipelines de dados.

Os tópicos a seguir fornecem informações detalhadas sobre o uso de modelos, incluindo como usálos para iniciar a execução de um trabalho e como alterar os parâmetros do modelo.

Tópicos

- Criação e uso de um modelo de trabalho para iniciar uma execução de trabalho
- Definição de parâmetros de modelos de trabalhos
- · Controle do acesso aos modelos de trabalhos

Criação e uso de um modelo de trabalho para iniciar uma execução de trabalho

Esta seção descreve a criação de um modelo de trabalho e o uso do modelo para iniciar a execução de um trabalho com o AWS Command Line Interface (AWS CLI).

Criar um modelo de trabalho

 Crie um arquivo create-job-template-request. json e especifique os parâmetros obrigatórios para seu modelo de trabalho, conforme mostrado no arquivo JSON de exemplo a seguir. Para obter informações sobre todos os parâmetros disponíveis, consulte a <u>CreateJobTemplate</u>API.

A maioria dos valores obrigatórios para a API StartJobRun também são obrigatórios para jobTemplateData. Se você quiser usar espaços reservados para qualquer parâmetro e fornecer valores ao invocar StartJobRun usando um modelo de trabalho, consulte a próxima seção sobre parâmetros do modelo de trabalho.

```
{
   "name": "mytemplate",
   "jobTemplateData": {
        "executionRoleArn": "iam_role_arn_for_job_execution",
        "releaseLabel": "emr-6.7.0-latest",
        "jobDriver": {
            "sparkSubmitJobDriver": {
                "entryPoint": "entryPoint_location",
                "entryPointArguments": [ "argument1", "argument2", ...],
                "sparkSubmitParameters": "--class <main_class> --conf
 spark.executor.instances=2 --conf spark.executor.memory=2G --conf
 spark.executor.cores=2 --conf spark.driver.cores=1"
        },
        "configurationOverrides": {
            "applicationConfiguration": [
                {
                    "classification": "spark-defaults",
                    "properties": {
                         "spark.driver.memory":"2G"
                    }
                }
            ],
            "monitoringConfiguration": {
                "persistentAppUI": "ENABLED",
                "cloudWatchMonitoringConfiguration": {
                    "logGroupName": "my_log_group",
                    "logStreamNamePrefix": "log_stream_prefix"
                },
                "s3MonitoringConfiguration": {
                    "logUri": "s3://my_s3_log_location/"
            }
        }
     }
}
```

2. Use o comando create-job-template com um caminho para o arquivo create-job-template-request.json armazenado localmente.

```
aws emr-containers create-job-template \
--cli-input-json file://./create-job-template-request.json
```

Iniciar uma execução de trabalho usando um modelo de trabalho

Forneça o ID do cluster virtual, o ID do modelo de trabalho e o nome do trabalho no comando StartJobRun, conforme mostrado no exemplo a seguir.

```
aws emr-containers start-job-run \
--virtual-cluster-id 123456 \
--name myjob \
--job-template-id 1234abcd
```

Definição de parâmetros de modelos de trabalhos

Os parâmetros de modelos de trabalhos permitem especificar variáveis no modelo de trabalho. Os valores para essas variáveis de parâmetro precisarão ser especificados ao iniciar uma execução de trabalho usando esse modelo de trabalho. Os parâmetros do modelo de trabalho são especificados no formato \${parameterName}. Você pode optar por especificar qualquer valor em um campo jobTemplateData como um parâmetro do modelo de trabalho. Para cada uma das variáveis de parâmetro do modelo de trabalho, especifique seu tipo de dados (STRING ou NUMBER) e, opcionalmente, um valor padrão. O exemplo abaixo mostra como você pode especificar os parâmetros do modelo de trabalho para os valores de local do ponto de entrada, classe principal e local do log do S3.

Para especificar o local do ponto de entrada, a classe principal e o local do log do Amazon S3 como parâmetros do modelo de trabalho

1. Crie um arquivo create-job-template-request. json e especifique os parâmetros obrigatórios para seu modelo de trabalho, conforme mostrado no arquivo JSON de exemplo a seguir. Para obter mais informações sobre os parâmetros, consulte a CreateJobTemplateAPI.

```
"name": "mytemplate",
"jobTemplateData": {
    "executionRoleArn": "iam_role_arn_for_job_execution",
    "releaseLabel": "emr-6.7.0-latest",
    "jobDriver": {
        "sparkSubmitJobDriver": {
            "entryPoint": "${EntryPointLocation}",
            "entryPointArguments": [ "argument1", "argument2",...],
            "sparkSubmitParameters": "--class ${MainClass} --conf
spark.executor.instances=2 --conf spark.executor.memory=2G --conf
spark.executor.cores=2 --conf spark.driver.cores=1"
```

```
}
        },
        "configurationOverrides": {
            "applicationConfiguration": [
                {
                     "classification": "spark-defaults",
                     "properties": {
                          "spark.driver.memory":"2G"
                }
            ],
            "monitoringConfiguration": {
                 "persistentAppUI": "ENABLED",
                "cloudWatchMonitoringConfiguration": {
                     "logGroupName": "my_log_group",
                     "logStreamNamePrefix": "log_stream_prefix"
                },
                "s3MonitoringConfiguration": {
                     "logUri": "${LogS3BucketUri}"
                }
            }
        },
        "parameterConfiguration": {
            "EntryPointLocation": {
                "type": "STRING"
            },
            "MainClass": {
                "type": "STRING",
                "defaultValue":"Main"
            },
            "LogS3BucketUri": {
                 "type": "STRING",
                 "defaultValue": "s3://my_s3_log_location/"
            }
        }
    }
}
```

2. Use o comando create-job-template com um caminho para o arquivo create-job-template-request.json armazenado localmente ou no Amazon S3.

```
aws emr-containers create-job-template \
--cli-input-json file://./create-job-template-request.json
```

Para iniciar uma execução de trabalho usando um modelo de trabalho com parâmetros de modelos de trabalhos

Para iniciar uma execução de trabalho com um modelo de trabalho contendo parâmetros de modelos de trabalhos, especifique o ID do modelo de trabalho, bem como os valores dos parâmetros do modelo de trabalho na solicitação da API StartJobRun, conforme mostrado abaixo.

```
aws emr-containers start-job-run \
--virtual-cluster-id 123456 \
--name myjob \
--job-template-id 1234abcd \
--job-template-parameters '{"EntryPointLocation": "entry_point_location", "MainClass":
"ExampleMainClass", "LogS3BucketUri": "s3://example_s3_bucket/"}'
```

Controle do acesso aos modelos de trabalhos

A política StartJobRun permite impor que um usuário ou um perfil possa somente executar trabalhos usando modelos de trabalhos especificados e não possa executar operações StartJobRun sem usar os modelos de trabalhos especificados. Para conseguir isso, primeiro, certifique-se de conceder ao usuário ou ao perfil uma permissão de leitura para os modelos de trabalhos especificados, conforme mostrado abaixo.

JSON

Para garantir que um usuário ou que um perfil seja capaz de invocar a operação StartJobRun somente ao usar modelos de trabalhos especificados, você pode atribuir a permissão de política StartJobRun apresentada a seguir para um determinado usuário ou perfil.

JSON

```
"Version": "2012-10-17",
  "Statement": [
      "Effect": "Allow",
      "Action": [
        "emr-containers:StartJobRun"
      ],
      "Resource": [
        "arn:aws:emr-containers:*:*:/virtualclusters/virtual_cluster_id"
      ],
      "Condition": {
        "ArnLike": {
          "emr-containers:JobTemplateArn": [
            "arn:aws:emr-containers:*:*:jobtemplate/job_template_1_id",
            "arn:aws:emr-containers:*:*:jobtemplate/job_template_2_id"
          ]
        }
      },
      "Sid": "AllowEMRCONTAINERSStartjobrun"
    }
  ]
}
```

Se o modelo de trabalho especificar um parâmetro de modelo de trabalho dentro do campo ARN do perfil de execução, o usuário poderá fornecer um valor para esse parâmetro e, assim, invocar StartJobRun usando um perfil de execução arbitrário. Para restringir os perfis de execução que o usuário pode fornecer, consulte Controle do acesso ao perfil de execução em Uso de perfis de execução de trabalho com o Amazon EMR no EKS.

Se nenhuma condição for especificada na política de ação StartJobRun acima para um determinado usuário ou perfil, o usuário ou o perfil terá permissão para invocar a ação StartJobRun no cluster virtual especificado usando um modelo de trabalho arbitrário ao qual ele tenha acesso de leitura ou usando um perfil de execução arbitrário.

Uso de modelos de pod

A partir das versões 5.33.0 ou 6.3.0 do Amazon EMR, o Amazon EMR no EKS oferece suporte ao recurso de modelo de pod do Spark. Um pod corresponde a um grupo de um ou mais contêineres, com recursos de rede e armazenamento compartilhados, e uma especificação de como executar os contêineres. Os modelos de pod são especificações que determinam como ocorrerá a execução de cada pod. É possível usar arquivos de modelo de pod para definir as configurações de pods de drivers ou de executores para as quais as configurações do Spark não oferecem suporte. Para obter mais informações sobre o recurso de modelo de pod do Spark, consulte Pod Template.



Note

O recurso de modelo de pod funciona somente com pods de drivers e de executores. Não é possível configurar pods de emissores de trabalhos usando o modelo de pods.

Cenários comuns

Você pode definir como executar trabalhos do Spark em clusters do EKS compartilhados ao usar modelos de pod com o Amazon EMR no EKS, e economizar custos e aprimorar a utilização e a performance dos recursos.

- Para reduzir custos, você pode programar tarefas do driver do Spark para serem executadas nas instâncias EC2 sob demanda da Amazon e programar tarefas do executor do Spark para serem executadas nas instâncias Amazon Spot. EC2
- Para aumentar a utilização de recursos, você pode oferecer suporte a diversas equipes que executam as workloads no mesmo cluster do EKS. Cada equipe receberá um grupo de EC2 nós designado da Amazon para executar suas cargas de trabalho. Você pode usar os modelos de pod para aplicar uma tolerância correspondente à workload.
- Para aprimorar o monitoramento, você pode executar um contêiner de registro em log separado para encaminhar os logs para a aplicação de monitoramento existente.

Por exemplo, o arquivo de modelo de pod a seguir demonstra um cenário de uso comum.

apiVersion: v1 kind: Pod

spec:

Uso de modelos de pod 218

```
volumes:
   - name: source-data-volume
     emptyDir: {}
   - name: metrics-files-volume
     emptyDir: {}
 nodeSelector:
   eks.amazonaws.com/nodegroup: emr-containers-nodegroup
 containers:
 - name: spark-kubernetes-driver # This will be interpreted as driver Spark main
container
   env:
     - name: RANDOM
       value: "random"
   volumeMounts:
     - name: shared-volume
       mountPath: /var/data
     - name: metrics-files-volume
       mountPath: /var/metrics/data
 - name: custom-side-car-container # Sidecar container
   image: <side_car_container_image>
   env:
     - name: RANDOM_SIDECAR
       value: random
   volumeMounts:
     - name: metrics-files-volume
       mountPath: /var/metrics/data
   command:
     - /bin/sh
     - '-c'
     - <command-to-upload-metrics-files>
 initContainers:
 - name: spark-init-container-driver # Init container
   image: <spark-pre-step-image>
   volumeMounts:
     - name: source-data-volume # Use EMR predefined volumes
       mountPath: /var/data
   command:
     - /bin/sh
     - '-c'
      <command-to-download-dependency-jars>
```

O modelo de pod conclui as seguintes tarefas:

Cenários comuns 219

- Adiciona um novo contêiner Init que é executado antes do início do contêiner principal do Spark.
 O contêiner init compartilha o EmptyDirvolume chamado source-data-volume com o contêiner principal do Spark. Você pode fazer com que seu contêiner Init execute as etapas de inicialização, como fazer download de dependências ou gerar dados de entrada. Em seguida, o contêiner principal do Spark consumirá os dados.
- Adiciona outro contêiner sidecar que é executado junto com o contêiner principal do Spark. Os dois contêineres estão compartilhando outro volume EmptyDir chamado metrics-files-volume.
 Seu trabalho do Spark pode gerar métricas, como as métricas do Prometheus. Em seguida, o trabalho do Spark pode colocar as métricas em um arquivo e fazer com que o contêiner sidecar carregue os arquivos em seu próprio sistema de BI para análise futura.
- Adiciona uma nova variável de ambiente ao contêiner principal do Spark. Você pode fazer com que seu trabalho consuma a variável de ambiente.
- Defina um <u>seletor de nó</u>, para que o pod seja programado somente no grupo de nós emrcontainers-nodegroup. Isso ajuda a isolar os recursos de computação entre trabalhos e equipes.

Habilitação de modelos de pod com o Amazon EMR no EKS

Para habilitar o recurso de modelo de pod com o Amazon EMR no EKS, configure as propriedades do Spark spark.kubernetes.driver.podTemplateFile e spark.kubernetes.executor.podTemplateFile para direcionar aos arquivos de modelo de pod no Amazon S3. O Spark fará download do arquivo de modelo de pod e o usará para estruturar pods de drivers e de executores.

Note

O Spark usa o perfil de execução de trabalho para carregar o modelo de pod, portanto, o perfil de execução de trabalho deve ter permissões para acessar o Amazon S3 com a finalidade de carregar os modelos de pod. Para obter mais informações, consulte Criação de um perfil de execução de trabalho.

Você pode usar SparkSubmitParameters para especificar o caminho do Amazon S3 para o modelo de pod, como demonstra o arquivo JSON de execução de trabalho a seguir.

```
{
  "name": "myjob",
```

```
"virtualClusterId": "123456",
  "executionRoleArn": "iam_role_name_for_job_execution",
  "releaseLabel": "release_label",
  "jobDriver": {
    "sparkSubmitJobDriver": {
      "entryPoint": "entryPoint_location",
      "entryPointArguments": ["argument1", "argument2", ...],
       "sparkSubmitParameters": "--class <main_class> \
         --conf
 spark.kubernetes.driver.podTemplateFile=s3://path_to_driver_pod_template \
 spark.kubernetes.executor.podTemplateFile=s3://path_to_executor_pod_template \
         --conf spark.executor.instances=2 \
         --conf spark.executor.memory=2G \
         --conf spark.executor.cores=2 \
         --conf spark.driver.cores=1"
    }
  }
}
```

Como alternativa, você pode usar configuration0verrides para especificar o caminho do Amazon S3 para o modelo de pod, como demonstra o arquivo JSON de execução de trabalho a seguir.

```
{
  "name": "myjob",
  "virtualClusterId": "123456",
  "executionRoleArn": "iam_role_name_for_job_execution",
  "releaseLabel": "release_label",
  "jobDriver": {
    "sparkSubmitJobDriver": {
      "entryPoint": "entryPoint_location",
      "entryPointArguments": ["argument1", "argument2", ...],
       "sparkSubmitParameters": "--class <main_class> \
         --conf spark.executor.instances=2 \
         --conf spark.executor.memory=2G \
         --conf spark.executor.cores=2 \
         --conf spark.driver.cores=1"
    }
  },
  "configurationOverrides": {
    "applicationConfiguration": [
```

Note

- Você precisa seguir as diretrizes de segurança ao usar o recurso de modelo de pod com o Amazon EMR no EKS, por exemplo, como isolar códigos de aplicações não confiáveis. Para obter mais informações, consulte <u>Práticas recomendadas de segurança para o</u> Amazon EMR no EKS.
- 2. Não é possível alterar os nomes dos contêineres principais do Spark ao usar spark.kubernetes.driver.podTemplateContainerName e spark.kubernetes.executor.podTemplateContainerName, porque esses nomes são codificados como spark-kubernetes-driver e spark-kubernetes-executors. Se desejar personalizar o contêiner principal do Spark, deverá especificar o contêiner em um modelo de pod com esses nomes codificados.

Campos do modelo de pod

Considere as restrições para campos apresentadas a seguir ao configurar um modelo de pod com o Amazon EMR no EKS.

 O Amazon EMR no EKS permite somente os campos a seguir em um modelo de pod para habilitar a programação adequada de trabalhos.

Estes são os campos permitidos em nível de pod:

- apiVersion
- kind
- metadata

Campos do modelo de pod 222

- spec.activeDeadlineSeconds
- spec.affinity
- spec.containers
- spec.enableServiceLinks
- spec.ephemeralContainers
- spec.hostAliases
- spec.hostname
- spec.imagePullSecrets
- spec.initContainers
- spec.nodeName
- spec.nodeSelector
- spec.overhead
- spec.preemptionPolicy
- spec.priority
- spec.priorityClassName
- spec.readinessGates
- spec.runtimeClassName
- spec.schedulerName
- spec.subdomain
- spec.terminationGracePeriodSeconds
- spec.tolerations
- spec.topologySpreadConstraints
- spec.volumes

Estes são os campos permitidos em nível de contêiner principal do Spark:

- env
- envFrom
- name
- lifecycle

- resources
- startupProbe
- stdin
- stdinOnce
- terminationMessagePath
- terminationMessagePolicy
- tty
- volumeDevices
- volumeMounts
- workingDir

Quando você usa qualquer campo não permitido no modelo de pod, o Spark gera uma exceção e o trabalho falha. O exemplo a seguir mostra uma mensagem de erro no log do controlador do Spark devido ao uso de campos não permitidos.

```
Executor pod template validation failed. Field container.command in Spark main container not allowed but specified.
```

 O Amazon EMR no EKS define previamente os parâmetros apresentados a seguir em um modelo de pod. Os campos especificados em um modelo de pod não devem se sobrepor a esses campos.

Estes são os nomes de volumes definidos previamente:

- emr-container-communicate
- config-volume
- emr-container-application-log-dir
- emr-container-event-log-dir
- temp-data-dir
- mnt-dir
- home-dir
- emr-container-s3

Estas são as montagens de volume definidas previamente que se aplicam somente ao contêiner principal do Spark:

- Nome:emr-container-application-log-dir; MountPath: /var/log/spark/user
- Nome:emr-container-event-log-dir; MountPath: /var/log/spark/apps
- Nome:mnt-dir; MountPath: /mnt
- Nome:temp-data-dir; MountPath: /tmp
- Nome:home-dir; MountPath: /home/hadoop

Estas são as variáveis de ambiente definidas previamente que se aplicam somente ao contêiner principal do Spark:

- SPARK_CONTAINER_ID
- K8S_SPARK_LOG_URL_STDERR
- K8S_SPARK_LOG_URL_STDOUT
- SIDECAR_SIGNAL_FILE

Note

Você ainda pode usar esses volumes definidos previamente e montá-los em seus contêineres sidecar adicionais. Por exemplo, é possível usar emr-container-application-log-dir e montá-lo em seu próprio contêiner sidecar definido no modelo de pod.

Se os campos especificados entrarem em conflito com qualquer um dos campos definidos previamente no modelo de pod, o Spark gerará uma exceção e o trabalho falhará. O exemplo a seguir mostra uma mensagem de erro no log da aplicação do Spark devido a conflitos com os campos definidos previamente.

Defined volume mount path on main container must not overlap with reserved mount paths: [<reserved-paths>]

Considerações sobre contêineres sidecar

O Amazon EMR controla o ciclo de vida dos pods provisionados pelo Amazon EMR no EKS. Os contêineres sidecar devem seguir um ciclo de vida semelhante ao do contêiner principal do Spark. Se você injetar contêineres sidecar adicionais em seus pods, recomendamos a integração com o

gerenciamento do ciclo de vida do pod definido pelo Amazon EMR para que o contêiner sidecar possa ser interrompido sozinho quando o contêiner principal do Spark sair.

Para reduzir custos, recomendamos implementar um processo que impeça que pods de drivers com contêineres sidecar continuem em execução após a conclusão do seu trabalho. O driver do Spark exclui os pods de executores quando o executor é encerrado. No entanto, quando um programa de driver é concluído, os contêineres sidecar adicionais continuam em execução. O pod é faturado até que o Amazon EMR no EKS limpe o pod do driver. Geralmente, isso ocorre menos de um minuto após a conclusão do contêiner principal do driver do Spark. Para reduzir custos, você pode integrar os contêineres sidecar adicionais ao mecanismo de gerenciamento do ciclo de vida que o Amazon EMR no EKS define para os pods de drivers e de executores, conforme descrito na seção a seguir.

O contêiner principal do Spark nos pods de drivers e de executores envia heartbeat para um arquivo /var/log/fluentd/main-container-terminated a cada dois segundos. Ao adicionar a montagem de volume emr-container-communicate definida previamente do Amazon EMR ao contêiner sidecar, você pode definir um subprocesso do contêiner sidecar para monitorar periodicamente o horário da última modificação desse arquivo. O subprocesso é interrompido sozinho caso descubra que o contêiner principal do Spark interrompeu a heartbeat por um período mais longo.

O exemplo apresentado a seguir demonstra um subprocesso que rastreia o arquivo de pulsação e é interrompido sozinho. *your_volume_mount* Substitua pelo caminho em que você monta o volume predefinido. O script é empacotado dentro da imagem usada pelo contêiner sidecar. Em um arquivo de modelo de pod, você pode especificar um contêiner sidecar com os seguintes comandos: sub_process_script.sh e main_command.

```
MOUNT_PATH="your_volume_mount"

FILE_TO_WATCH="$MOUNT_PATH/main-container-terminated"

INITIAL_HEARTBEAT_TIMEOUT_THRESHOLD=60

HEARTBEAT_TIMEOUT_THRESHOLD=15

SLEEP_DURATION=10

function terminate_main_process() {
    # Stop main process
}

# Waiting for the first heartbeat sent by Spark main container echo "Waiting for file $FILE_TO_WATCH to appear..."

start_wait=$(date +%s)

while ! [[ -f "$FILE_TO_WATCH" ]]; do
```

```
elapsed_wait=$(expr $(date +%s) - $start_wait)
    if [ "$elapsed_wait" -gt "$INITIAL_HEARTBEAT_TIMEOUT_THRESHOLD" ]; then
        echo "File $FILE_TO_WATCH not found after $INITIAL_HEARTBEAT_TIMEOUT_THRESHOLD
 seconds; aborting"
        terminate_main_process
        exit 1
    fi
    sleep $SLEEP_DURATION;
done;
echo "Found file $FILE_TO_WATCH; watching for heartbeats..."
while [[ -f "$FILE_TO_WATCH" ]]; do
    LAST_HEARTBEAT=$(stat -c %Y $FILE_TO_WATCH)
    ELAPSED_TIME_SINCE_AFTER_HEARTBEAT=$(expr $(date +%s) - $LAST_HEARTBEAT)
    if [ "$ELAPSED_TIME_SINCE_AFTER_HEARTBEAT" -gt "$HEARTBEAT_TIMEOUT_THRESHOLD" ];
 then
        echo "Last heartbeat to file $FILE_TO_WATCH was more than
 $HEARTBEAT_TIMEOUT_THRESHOLD seconds ago at $LAST_HEARTBEAT; terminating"
        terminate_main_process
        exit 0
    fi
    sleep $SLEEP_DURATION;
done;
echo "Outside of loop, main-container-terminated file no longer exists"
# The file will be deleted once the fluentd container is terminated
echo "The file $FILE_TO_WATCH doesn't exist any more;"
terminate_main_process
exit 0
```

Uso de políticas de repetição de trabalho

Nas versões 6.9.0 e posteriores do Amazon EMR no EKS, você pode definir uma política de repetição para as execuções de trabalhos. As políticas de repetição fazem com que um pod do driver de trabalho seja reiniciado automaticamente se falhar ou for excluído. Isso torna os trabalhos de transmissão do Spark de execução prolongada mais resistentes a falhas.

Definição de uma política de repetição para um trabalho

Para configurar uma política de repetição, você fornece um RetryPolicyConfiguration campo usando a StartJobRunAPI. Um exemplo de retryPolicyConfiguration é mostrado aqui:

Uso de políticas de repetição 227

```
aws emr-containers start-job-run \
--virtual-cluster-id cluster_id \
--name sample-job-name \
--execution-role-arn execution-role-arn \
--release-label emr-6.9.0-latest \
--job-driver '{
  "sparkSubmitJobDriver": {
    "entryPoint": "local:///usr/lib/spark/examples/src/main/python/pi.py",
    "entryPointArguments": [ "2" ],
    "sparkSubmitParameters": "--conf spark.executor.instances=2 --conf
 spark.executor.memory=2G --conf spark.executor.cores=2 --conf spark.driver.cores=1"
  }
}'\
--retry-policy-configuration '{
    "maxAttempts": 5
  }'\
--configuration-overrides '{
  "monitoringConfiguration": {
    "cloudWatchMonitoringConfiguration": {
      "logGroupName": "my_log_group_name",
      "logStreamNamePrefix": "my_log_stream_prefix"
    },
    "s3MonitoringConfiguration": {
       "logUri": "s3://amzn-s3-demo-logging-bucket"
    }
  }
}'
```

Note

retryPolicyConfigurationsó está disponível a partir da versão AWS CLI 1.27.68. Para atualizar o AWS CLI para a versão mais recente, consulte <u>Instalando ou atualizando a versão</u> mais recente do AWS CLI

Configure o campo maxAttempts com o número máximo de vezes que você deseja que o pod do driver de trabalho seja reiniciado em casos de falhas ou exclusões. O intervalo de execução entre duas tentativas de repetição do driver de trabalho é um intervalo de repetição exponencial de (10 segundos, 20 segundos, 40 segundos...) que é limitado a seis minutos, conforme descrito na documentação do Kubernetes.



Note

Cada execução adicional de um driver de trabalho será cobrada como outra execução de trabalho e estará sujeita aos preços do Amazon EMR no EKS.

Valores de configuração da política de repetição

 Política de repetição padrão para um trabalho: StartJobRun inclui uma política de repetição definida como uma tentativa máxima, por padrão. Você pode configurar a política de repetição conforme desejado.



Note

Se maxAttempts da retryPolicyConfiguration for definido como uma, significa que nenhuma nova tentativa será feita para ativar o pod do driver em caso de falha.

• Desabilitando a política de repetição de um trabalho: para desativar uma política de repetição, defina o valor máximo de tentativas como 1. retryPolicyConfiguration

```
"retryPolicyConfiguration": {
    "maxAttempts": 1
}
```

Definição de maxAttempts para um trabalho dentro do intervalo válido: a chamada StartJobRun falhará se o valor de maxAttempts estiver fora do intervalo válido. O intervalo válido de maxAttempts é de 1 a 2.147.483.647 (número inteiro de 32 bits), que corresponde ao intervalo compatível com a configuração backOffLimit do Kubernetes. Para obter mais informações, consulte Pod backoff failure policy na documentação do Kubernetes. Se o valor de maxAttempts for inválido, a seguinte mensagem de erro será retornada:

```
{
 "message": "Retry policy configuration's parameter value of maxAttempts is invalid"
}
```

Recuperação de um status da política de repetição para um trabalho

Você pode ver o status das tentativas de repetição de um trabalho com <u>ListJobRuns</u>o.

<u>DescribeJobRun</u> APIs Depois de solicitar um trabalho com uma configuração de política de repetição habilitada, as respostas ListJobRun e DescribeJobRun conterão o status da política de repetição no campo RetryPolicyExecution. Além disso, a resposta DescribeJobRun conterá a RetryPolicyConfiguration que foi inserida na solicitação StartJobRun para o trabalho.

Exemplo de respostas

ListJobRuns response

DescribeJobRun response

```
{
    ...
    "retryPolicyConfiguration": {
        "maxAttempts": 5
    },
    "retryPolicyExecution" : {
        "currentAttemptCount": 2
    },
    ...
    ...
}
```

Esses campos não ficarão visíveis quando a política de repetição estiver desabilitada no trabalho, conforme descrito abaixo em Valores de configuração da política de repetição.

Monitoramento de um trabalho com uma política de repetição

Quando você ativa uma política de repetição, um CloudWatch evento é gerado para cada driver de trabalho criado. Para se inscrever nesses eventos, configure uma regra de CloudWatch evento usando o seguinte comando:

```
aws events put-rule \
--name cwe-test \
--event-pattern '{"detail-type": ["EMR Job Run New Driver Attempt"]}'
```

O evento retornará informações sobre o newDriverPodName, o carimbo de data/hora newDriverCreatedAt, a previousDriverFailureMessage e o currentAttemptCount dos drivers de trabalho. Esses eventos não serão criados se a política de repetição estiver desabilitada.

Para obter mais informações sobre como monitorar seu trabalho com CloudWatch eventos, consulteMonitore trabalhos com a Amazon CloudWatch Events.

Descoberta de logs para drivers e executores

Os nomes dos pods de drivers seguem o formato spark-<job id>-driver-<random-suffix>. O mesmo random-suffix é adicionado aos nomes dos pods de executores que o driver gera. Ao usar esse random-suffix, você pode localizar logs de um driver e seus executores associados. O random-suffix estará presente somente se a política de repetição estiver habilitada para o trabalho. Caso contrário, o random-suffix estará ausente.

Para obter mais informações sobre como configurar trabalhos com uma configuração de monitoramento para o registro em log, consulte Execução de uma aplicação do Spark.

Uso da alternância de log de eventos do Spark

Com as versões 6.3.0 e posteriores do Amazon EMR, você pode ativar o recurso de alternância de log de eventos do Spark para o Amazon EMR no EKS. Em vez de gerar um único arquivo de log de eventos, esse recurso alterna o arquivo com base no intervalo de tempo configurado e remove os arquivos de log de eventos mais antigos.

A alternância de logs de eventos do Spark pode ajudar você a evitar possíveis problemas com um grande arquivo de log de eventos do Spark gerado para trabalhos de execução prolongada ou de transmissão. Por exemplo, você inicia um trabalho de execução prolongada do Spark com um log

Monitoramento do trabalho 231

de eventos habilitado com o parâmetro persistentAppUI. O driver do Spark gera um arquivo de log de eventos. Se o trabalho for executado por horas ou por dias, e houver um espaço em disco limitado no nó do Kubernetes, o arquivo de log de eventos poderá consumir todo o espaço em disco disponível. Ativar o recurso de alternância de log de eventos do Spark resolve o problema ao dividir o arquivo de log em vários arquivos e remover os arquivos mais antigos.



Note

Esse recurso funciona somente com o Amazon EMR no EKS. O Amazon EMR executado na Amazon EC2 não oferece suporte à rotação do log de eventos do Spark.

Para ativar o recurso de alternância de log de eventos do Spark, configure os seguintes parâmetros do Spark:

- spark.eventLog.rotation.enabled: ativa a alternância de log. Por padrão, ele está desabilitado no arquivo de configuração do Spark. Defina-o como verdadeiro para ativar esse recurso.
- spark.eventLog.rotation.interval: especifica o intervalo de tempo para a alternância de log. O valor mínimo é 60 segundos. O valor de padrão é de 300 segundos.
- spark.eventLog.rotation.minFileSize: especifica um tamanho mínimo de arquivo para alternar o arquivo de log. O valor mínimo e padrão é de 1 MB.
- spark.eventLog.rotation.maxFilesToRetain: especifica quantos arquivos de log alternados serão mantidos durante a limpeza. O intervalo válido é de 1 a 10. O valor padrão é 2.

Você pode especificar esses parâmetros na seção sparkSubmitParameters da API StartJobRun, como mostra o exemplo a seguir.

```
"sparkSubmitParameters": "--class org.apache.spark.examples.SparkPi --conf
spark.eventLog.rotation.enabled=true --conf spark.eventLog.rotation.interval=300 --
conf spark.eventLog.rotation.minFileSize=1m --conf
spark.eventLog.rotation.maxFilesToRetain=2"
```

Uso da alternância de log do contêiner do Spark

Com as versões 6.11.0 e posteriores do Amazon EMR, você pode ativar o recurso de alternância de log do contêiner do Spark para o Amazon EMR no EKS. Em vez de gerar um arquivo único

de log stdout ou stderr, esse recurso alterna o arquivo com base na dimensão de alternância configurada e remove os arquivos de log mais antigos do contêiner.

A alternância de logs do contêiner do Spark pode ajudar você a evitar possíveis problemas com grandes arquivos de log do Spark gerados para trabalhos de execução prolongada ou de transmissão. Por exemplo, é possível iniciar um trabalho do Spark de execução prolongada e o driver do Spark gerar um arquivo de log do contêiner. Se o trabalho for executado por horas ou por dias, e houver um espaço em disco limitado no nó do Kubernetes, o arquivo de log do contêiner poderá consumir todo o espaço em disco disponível. Ao ativar a alternância de log do contêiner do Spark, você divide o arquivo de log em vários arquivos e remove os arquivos mais antigos.

Para ativar o recurso de alternância de log do contêiner do Spark, configure os seguintes parâmetros do Spark:

containerLogRotationConfiguration

Inclua esse parâmetro em monitoringConfiguration para ativar a alternância de log. O recurso é desabilitado por padrão. Você deve usar containerLogRotationConfiguration além de s3MonitoringConfiguration.

rotationSize

O parâmetro rotationSize especifica o tamanho do arquivo para a alternância de log. O intervalo de valores possíveis é de 2KB a 2GB. A parcela numérica da unidade do parâmetro rotationSize é transferida como um número inteiro. Como não há suporte para valores decimais, você pode especificar um tamanho de rotação de 1,5 GB, por exemplo, com o valor 1500MB.

maxFilesToKeep

O parâmetro maxFilesToKeep especifica o número máximo de arquivos a serem retidos no contêiner após a alternância ter ocorrido. O valor mínimo é de 1 e o valor máximo é de 50.

Você pode especificar esses parâmetros na seção monitoringConfiguration da API StartJobRun, como mostra o exemplo a seguir. Neste exemplo, com rotationSize = "10 MB" e maxFilesToKeep = 3, o Amazon EMR no EKS alterna seus logs em 10 MB, gera um novo arquivo de log e, em seguida, limpa o arquivo de log mais antigo quando o número de arquivos de log atinge três.

```
{
    "name": "my-long-running-job",
```

```
"virtualClusterId": "123456",
  "executionRoleArn": "iam_role_name_for_job_execution",
  "releaseLabel": "emr-6.11.0-latest",
  "jobDriver": {
    "sparkSubmitJobDriver": {
      "entryPoint": "entryPoint_location",
      "entryPointArguments": ["argument1", "argument2", ...],
       "sparkSubmitParameters": "--class main_class --conf spark.executor.instances=2
 --conf spark.executor.memory=2G --conf spark.executor.cores=2 --conf
 spark.driver.cores=1"
    }
  },
  "configurationOverrides": {
    "applicationConfiguration": [
      {
        "classification": "spark-defaults",
        "properties": {
          "spark.driver.memory":"2G"
         }
      }
    ],
    "monitoringConfiguration": {
      "persistentAppUI": "ENABLED",
      "cloudWatchMonitoringConfiguration": {
        "logGroupName": "my_log_group",
        "logStreamNamePrefix": "log_stream_prefix"
      },
      "s3MonitoringConfiguration": {
        "logUri": "s3://my_s3_log_location"
      },
      "containerLogRotationConfiguration": {
        "rotationSize":"10MB",
        "maxFilesToKeep":"3"
      }
    }
  }
}
```

Para iniciar uma execução de trabalho com a alternância de log do contêiner do Spark, inclua um caminho para o arquivo JSON que você configurou com esses parâmetros no comando StartJobRun.

```
aws emr-containers start-job-run \
```

--cli-input-json file://path-to-json-request-file

Uso da escalabilidade automática vertical com trabalhos do Spark no Amazon EMR

A escalabilidade automática vertical do Amazon EMR no EKS ajusta automaticamente os recursos de memória e de CPU para se adaptar às necessidades da workload fornecida para aplicações do Spark no Amazon EMR. Isso simplifica o gerenciamento de recursos.

Para rastrear a utilização histórica e em tempo real dos recursos de suas aplicações do Spark no Amazon EMR, a escalabilidade automática vertical utiliza o <u>Vertical Pod Autoscaler (VPA)</u> do Kubernetes. A funcionalidade de escalabilidade automática vertical usa os dados que o VPA coleta para ajustar automaticamente os recursos de memória e de CPU atribuídos às aplicações do Spark. Este processo simplificado aumenta a confiabilidade e otimiza os custos.

Tópicos

- Configuração da escalabilidade automática vertical para o Amazon EMR no EKS
- Conceitos básicos da escalabilidade automática vertical para o Amazon EMR no EKS
- Configuração da escalabilidade automática vertical para o Amazon EMR no EKS
- Monitoramento da escalabilidade automática vertical para o Amazon EMR no EKS
- Desinstalação do operador de escalabilidade automática vertical do Amazon EMR no EKS

Configuração da escalabilidade automática vertical para o Amazon EMR no EKS

Este tópico ajuda você a preparar o cluster do Amazon EKS para enviar trabalhos do Spark no Amazon EMR com a escalabilidade automática vertical. O processo de configuração requer que você confirme ou conclua as tarefas apresentadas nas seguintes seções:

Tópicos

- Pré-requisitos
- Instalação do Operator Lifecycle Manager (OLM) no cluster do Amazon EKS
- Instalação do operador de escalabilidade automática vertical do Amazon EMR no EKS

Pré-requisitos

Conclua as tarefas apresentadas a seguir antes de instalar o operador do Kubernetes para escalabilidade automática vertical em seu cluster. Se você já completou algum dos pré-requisitos, pode ignorá-los e passar para os próximos.

- Instale ou atualize para a versão mais recente do AWS CLI Se você já instalou o AWS CLI, confirme se você tem a versão mais recente.
- <u>Instale o kubectl</u>: o kubectl é uma ferramenta de linha de comando que você usa para se comunicar com o servidor de API do Kubernetes. Você precisa do kubectl para instalar e monitorar artefatos relacionados à escalabilidade automática vertical no cluster do Amazon EKS.
- <u>Instale o Operator SDK</u>: o Amazon EMR no EKS usa o Operator SDK como um gerenciador de pacotes durante a vida útil do operador de escalabilidade automática vertical instalado no cluster.
- <u>Instale o Docker</u>: você precisa de acesso à CLI do Docker para autenticar e buscar as imagens do Docker relacionadas à escalabilidade automática vertical para instalar no cluster do Amazon EKS.
- <u>Instale o servidor de métricas do Kubernetes</u>: você deve primeiro instalar o servidor de métricas para que o escalador automático vertical do pod possa buscar métricas do servidor da API do Kubernetes.
- <u>Comece a usar o Amazon EKS (eksctl)</u> (versão 1.24 ou superior): o ajuste de escala automático vertical é compatível com o Amazon EKS nas versões 1.24 e superior. Após criar o cluster, registre-o para uso com o Amazon EMR.
- <u>Selecione um URI de imagem base do Amazon EMR</u> (versão 6.10.0 ou superiores): a escalabilidade automática vertical é compatível com as versões 6.10.0 e superiores do Amazon EMR.

Instalação do Operator Lifecycle Manager (OLM) no cluster do Amazon EKS

Use a CLI do Operator SDK para instalar o Operator Lifecycle Manager (OLM) no cluster do Amazon EMR no EKS no qual você deseja configurar a escalabilidade automática vertical, conforme mostrado no exemplo a seguir. Após configurá-lo, será possível usar o OLM para instalar e gerenciar o ciclo de vida do operador de escalabilidade automática vertical do Amazon EMR.

operator-sdk olm install

Para validar a instalação, execute o comando olm status:

Configuração 236

```
operator-sdk olm status
```

Verifique se o comando retorna um resultado com êxito semelhante ao seguinte exemplo de saída:

```
INFO[0007] Successfully got OLM status for version X.XX
```

Se a instalação não ocorrer com êxito, consulte Solução de problemas de escalabilidade automática vertical do Amazon EMR no EKS.

Instalação do operador de escalabilidade automática vertical do Amazon EMR no EKS

Use as seguintes etapas para instalar o operador de escalabilidade automática vertical no cluster do Amazon EKS:

- 1. Configure as seguintes variáveis de ambiente que serão usadas para concluir a instalação:
 - \$REGION direciona para a Região da AWS do seu cluster. Por exemplo, .us-west-2
 - \$ACCOUNT_ID direciona para o ID da conta do Amazon ECR para sua região. Para obter mais informações, consulte Contas de registro do Amazon ECR por região.
 - \$RELEASE direciona para a versão do Amazon EMR que você deseja usar para o cluster.
 Com a escalabilidade automática vertical, você deve usar o Amazon EMR versão 6.10.0 ou superiores.
- 2. Em seguida, obtenha tokens de autenticação para o registro do Amazon ECR do operador.

```
aws ecr get-login-password \
   --region region-id | docker login \
   --username AWS \
   --password-stdin $ACCOUNT_ID.dkr.ecr.region-id.amazonaws.com
```

3. Instale o operador de escalabilidade automática vertical do Amazon EMR no EKS com o seguinte comando:

```
ECR_URL=$ACCOUNT_ID.dkr.ecr.$REGION.amazonaws.com && \
REPO_DEST=dynamic-sizing-k8s-operator-olm-bundle && \
BUNDLE_IMG=emr-$RELEASE-dynamic-sizing-k8s-operator && \
operator-sdk run bundle \
$ECR_URL/$REPO_DEST/$BUNDLE_IMG\:latest
```

Isso criará uma versão do operador de escalabilidade automática vertical no namespace padrão do cluster do Amazon EKS. Use este comando para realizar a instalação em um namespace diferente:

```
operator-sdk run bundle \
$ACCOUNT_ID.dkr.ecr.$REGION.amazonaws.com/dynamic-sizing-k8s-operator-olm-bundle/
emr-$RELEASE-dynamic-sizing-k8s-operator:latest \
-n operator-namespace
```



Se o namespace especificado não existir, o OLM não instalará o operador. Para obter mais informações, consulte <u>Namespace do Kubernetes não encontrado</u>.

4. Verifique se você instalou o operador com êxito usando a ferramenta de linha de comando kubectl do Kubernetes.

```
kubectl get csv -n operator-namespace
```

O comando kubect1 deve retornar o operador de escalabilidade automática vertical recémimplantado com um status de Fase como Com êxito. Se você tiver problemas com a instalação ou com a configuração, consulte Solução de problemas de escalabilidade automática vertical do Amazon EMR no EKS.

Conceitos básicos da escalabilidade automática vertical para o Amazon EMR no EKS

Use o ajuste de escala automático vertical para o Amazon EMR no EKS quando desejar que o ajuste automático de memória e recursos de CPU se adapte à workload da aplicação do Spark no Amazon EMR. Para obter mais informações, consulte <u>Using vertical autoscaling with Amazon EMR Spark</u> jobs.

Envio de um trabalho do Spark com a escalabilidade automática vertical

Ao enviar um trabalho por meio da <u>StartJobRun</u>API, adicione as duas configurações a seguir ao driver do seu trabalho do Spark para ativar o escalonamento automático vertical:

Começar 238

```
"spark.kubernetes.driver.annotation.emr-containers.amazonaws.com/
dynamic.sizing":"true",
"spark.kubernetes.driver.annotation.emr-containers.amazonaws.com/
dynamic.sizing.signature":"YOUR_JOB_SIGNATURE"
```

No código acima, a primeira linha ativa a funcionalidade de escalabilidade automática vertical. A linha seguinte é uma configuração de assinatura obrigatória que permite que você escolha uma assinatura para o seu trabalho.

Para obter mais informações sobre essas configurações e valores de parâmetros aceitáveis, consulte Configuração da escalabilidade automática vertical para o Amazon EMR no EKS. Por padrão, seu trabalho é enviado com a opção somente monitoramento no modo Desativado para a escalabilidade automática vertical. Esse estado de monitoramento permite calcular e visualizar recomendações de recursos sem realizar a escalabilidade automática. Para obter mais informações, consulte Modos da escalabilidade automática vertical.

O seguinte exemplo mostra como concluir um comando start-job-run de exemplo com a escalabilidade automática vertical:

```
aws emr-containers start-job-run \
--virtual-cluster-id $VIRTUAL_CLUSTER_ID \
--name $JOB_NAME \
--execution-role-arn $EMR_ROLE_ARN \
--release-label emr-6.10.0-latest \
--job-driver '{
  "sparkSubmitJobDriver": {
     "entryPoint": "local:///usr/lib/spark/examples/src/main/python/pi.py"
 }'\
--configuration-overrides '{
    "applicationConfiguration": [{
        "classification": "spark-defaults",
        "properties": {
          "spark.kubernetes.driver.annotation.emr-containers.amazonaws.com/
dynamic.sizing": "true",
          "spark.kubernetes.driver.annotation.emr-containers.amazonaws.com/
dynamic.sizing.signature": "test-signature"
    }]
  }'
```

Começar 239

Verificação de funcionalidade da escalabilidade automática vertical

Para verificar se a escalabilidade automática vertical está funcionando corretamente para o trabalho enviado, use kubectl para obter o recurso personalizado verticalpodautoscaler e visualizar suas recomendações de escalabilidade. Por exemplo, o seguinte comando consulta recomendações sobre o trabalho de exemplo da seção Envio de um trabalho do Spark com a escalabilidade automática vertical:

```
kubectl get verticalpodautoscalers --all-namespaces \
-l=emr-containers.amazonaws.com/dynamic.sizing.signature=test-signature
```

A saída desta consulta deve ser semelhante à seguinte:

NAME	MODE	CPU	MEM
PROVIDED AGE ds-jceyefkxnhrvdzw6djum3naf2abm6o63a6dvjkkedqtkhlrf25eq-vpa	Off	3304504865	True
87m			

Se a saída não for semelhante ou tiver um código de erro, consulte <u>Solução de problemas de</u> <u>escalabilidade automática vertical do Amazon EMR no EKS</u> para obter passos que ajudam a resolver o problema.

Configuração da escalabilidade automática vertical para o Amazon EMR no EKS

Você pode configurar o escalonamento automático vertical ao enviar trabalhos do Amazon EMR Spark por meio da API. <u>StartJobRun</u> Defina os parâmetros de configuração relacionados à escalabilidade automática no pod do driver do Spark, conforme mostrado no exemplo em <u>Envio de</u> um trabalho do Spark com a escalabilidade automática vertical.

O operador de escalabilidade automática vertical do Amazon EMR no EKS recebe os pods de drivers que têm escalabilidade automática e, em seguida, configura a integração com o Vertical Pod Autoscaler (VPA) do Kubernetes usando as configurações no pod do driver. Isso facilita o rastreamento de recursos e a escalabilidade automática de pods de executores do Spark.

As seções a seguir descrevem os parâmetros que você pode usar ao configurar a escalabilidade automática vertical para o cluster do Amazon EKS.



Note

Configure o parâmetro de alternância de recursos como um rótulo e configure os parâmetros restantes como anotações no pod do driver do Spark. Os parâmetros de escalabilidade automática pertencem ao domínio emr-containers.amazonaws.com/ e têm o prefixo dynamic.sizing.

Parâmetros necessários

Você deve incluir os dois seguintes parâmetros no driver de trabalho do Spark ao enviar o trabalho:

Chave	Descrição	Valores aceitos	Valor padrão	Tipo	Parâmetro do Spark ¹
dynamic.s izing	Alternância de recursos	true, false	não definido	rótulo	spark.kub ernetes.d river.lab el.emr-co ntainers. amazonaws .com/dyna mic.sizin g
dynamic.s izing.sig nature	Assinatura do trabalho	string	não definido	anotação	spark.kub ernetes.d river.ann otation.e mr-contai ners.amaz onaws.com /dynamic. sizing.si gnature

- ¹ Use esse parâmetro como um SparkSubmitParameter ou uma ConfigurationOverride na API StartJobRun.
- dynamic.sizing: é possível ativar e desativar a escalabilidade automática vertical com o rótulo dynamic.sizing. Para ativar a escalabilidade automática vertical, defina dynamic.sizing como true no pod do driver do Spark. Se você omitir esse rótulo ou defini-lo com qualquer valor diferente de true, a escalabilidade automática vertical será desativada.
- dynamic.sizing.signature: defina a assinatura do trabalho com a anotação dynamic.sizing.signature no pod do driver. A escalabilidade automática vertical agrega dados de uso de recursos em diferentes execuções de trabalhos do Spark no Amazon EMR para derivar recomendações de recursos. Você fornece o identificador exclusivo para unir os trabalhos.

Note

Se o seu trabalho se repetir em um intervalo fixo, como diariamente ou semanalmente, a assinatura do trabalho deverá permanecer a mesma para cada nova instância dele. Isso garante que a escalabilidade automática vertical possa computar e agregar recomendações em diferentes execuções do trabalho.

Parâmetros opcionais

A escalabilidade automática vertical também oferece suporte aos parâmetros opcionais apresentados a seguir. Defina-os como anotações no pod do driver.

Chave	Descrição	Valores aceitos	Valor padrão	Tipo	Parâmetro do Spark ¹
<pre>dynamic.s izing.mod e</pre>	Modo da escalabil idade automática vertical	Off, Initial, Auto	Off	anotação	spark.kub ernetes.d river.ann otation.e mr-contai ners.amaz

¹ Use esse parâmetro como um SparkSubmitParameter ou uma ConfigurationOverride na API StartJobRun.

Chave	Descrição	Valores aceitos	Valor padrão	Tipo	Parâmetro do Spark ¹
					onaws.com /dynamic. sizing.mo de
dynamic.s izing.sca le.memory	Habilitação da escalabil idade de memória	true, false	true	anotação	spark.kub ernetes.d river.ann otation.e mr-contai ners.amaz onaws.com /dynamic. sizing.sc ale.memor y
dynamic.s izing.sca le.cpu	Ativação ou desativação da escalabil idade de CPU	true, false	false	anotação	spark.kub ernetes.d river.ann otation.e mr-contai ners.amaz onaws.com /dynamic. sizing.sc ale.cpu

Chave	Descrição	Valores aceitos	Valor padrão	Tipo	Parâmetro do Spark ¹
dynamic.s izing.sca le.memory .min	Limite mínimo para a escalabil idade de memória	string, K8s resource. Quantity, por exemplo: 1G	não definido	anotação	spark.kub ernetes.d river.ann otation.e mr-contai ners.amaz onaws.com /dynamic. sizing.sc ale.memor y.min
dynamic.s izing.sca le.memory .max	Limite máximo para a escalabil idade de memória	string, K8s resource. Quantity, por exemplo: 4G	não definido	anotação	spark.kub ernetes.d river.ann otation.e mr-contai ners.amaz onaws.com /dynamic. sizing.sc ale.memor y.max

Chave	Descrição	Valores aceitos	Valor padrão	Tipo	Parâmetro do Spark ¹
<pre>dynamic.s izing.sca le.cpu.mi n</pre>	Limite mínimo para a escalabil idade de CPU	string, K8s resource. Quantity, por exemplo: 1	não definido	anotação	spark.kub ernetes.d river.ann otation.e mr-contai ners.amaz onaws.com /dynamic. sizing.sc ale.cpu.m in
dynamic.s izing.sca le.cpu.ma x	Limite máximo para a escalabil idade de CPU	string, K8s resource. Quantity, por exemplo: 2	não definido	anotação	spark.kub ernetes.d river.ann otation.e mr-contai ners.amaz onaws.com /dynamic. sizing.sc ale.cpu.m ax

Modos da escalabilidade automática vertical

O parâmetro mode é mapeado para os diferentes modos de escalabilidade automática compatíveis com o VPA. Use a anotação dynamic.sizing.mode no pod do driver para definir o modo. Os seguintes valores têm suporte para este parâmetro:

 Desativado: um modo de simulação no qual é possível monitorar recomendações, mas a escalabilidade automática não é executada. Este é o modo padrão para a escalabilidade automática vertical. Nesse modo, o recurso de escalabilidade automática vertical do pod associado

computa recomendações e você pode monitorá-las usando ferramentas como o kubectl, o Prometheus e o Grafana.

- Inicial: neste modo, o VPA escala automaticamente os recursos quando o trabalho é iniciado, se houver recomendações disponíveis com base no histórico de execuções do trabalho, como no caso de um trabalho recorrente.
- Automático: neste modo, o VPA realiza a remoção de pods de executores do Spark e os escala automaticamente com as configurações de recursos recomendadas quando o pod do driver do Spark os reinicia. Às vezes, o VPA realiza a remoção de pods de executores do Spark em execução, portanto, isso pode resultar em latência adicional ao tentar novamente o executor interrompido.

Escalabilidade de recursos

Ao configurar a escalabilidade automática vertical, você pode escolher se deseja escalar os recursos de CPU e de memória. Defina as anotações dynamic.sizing.scale.cpu e dynamic.sizing.scale.memory como true ou false. Por padrão, a escalabilidade de CPU é definida como false e a escalabilidade de memória é definida como true.

Mínimos e máximos de recursos (limites)

Como opção, também é possível definir limites para os recursos de CPU e de memória. Escolha valores mínimo e máximo para esses recursos com as anotações dynamic.sizing.[memory/cpu].[min/max] ao ativar a escalabilidade automática. Por padrão, os recursos não têm limitações. Defina as anotações como valores de string que representam uma quantidade de recursos do Kubernetes. Por exemplo, defina dynamic.sizing.memory.max como 4G para representar 4 GB.

Monitoramento da escalabilidade automática vertical para o Amazon EMR no EKS

Você pode usar a ferramenta de linha de comando kubectl do Kubernetes para listar as recomendações ativas relacionadas ao ajuste de escala automático vertical em seu cluster. Também é possível visualizar as assinaturas de trabalho rastreadas e limpar quaisquer recursos desnecessários associados às assinaturas.

Listagem das recomendações de escalabilidade automática vertical para o cluster

Use kubecti para obter o recurso verticalpodautoscaler e visualizar o status atual e as recomendações. O exemplo de consulta a seguir retorna todos os recursos ativos no cluster do Amazon EKS.

```
kubectl get verticalpodautoscalers \
-o custom-columns="NAME:.metadata.name,"\
"SIGNATURE:.metadata.labels.emr-containers\.amazonaws\.com/dynamic\.sizing
\.signature,"\
"MODE:.spec.updatePolicy.updateMode,"\
"MEM:.status.recommendation.containerRecommendations[0].target.memory" \
--all-namespaces
```

A saída desta consulta é semelhante à seguinte:

```
NAME SIGNATURE MODE MEM
ds-example-id-1-vpa job-signature-1 Off none
ds-example-id-2-vpa job-signature-2 Initial 12936384283
```

Consulta e exclusão das recomendações de escalabilidade automática vertical para o cluster

Quando você exclui um recurso de execução de trabalho de escalabilidade automática vertical do Amazon EMR, ele exclui automaticamente o objeto VPA associado que rastreia e armazena recomendações.

O seguinte exemplo usa kubecti para limpar as recomendações para um trabalho identificado por uma assinatura:

```
kubectl delete jobrun -n emr -l=emr-containers\.amazonaws\.com/dynamic\.sizing
\.signature=integ-test
jobrun.dynamicsizing.emr.services.k8s.aws "ds-job-signature" deleted
```

Se você não souber a assinatura específica do trabalho ou desejar limpar todos os recursos do cluster, poderá usar --all ou --all-namespaces em seu comando, em vez do ID exclusivo do trabalho, conforme mostrado no seguinte exemplo:

```
kubectl delete jobruns --all --all-namespaces
```

jobrun.dynamicsizing.emr.services.k8s.aws "ds-example-id" deleted

Desinstalação do operador de escalabilidade automática vertical do Amazon EMR no EKS

Se você desejar remover o operador de escalabilidade automática vertical do cluster do Amazon EKS, use o comando cleanup com a CLI do Operator SDK, conforme mostrado no exemplo a seguir. Isso também excluirá as dependências upstream instaladas com o operador, como o Vertical Pod Autoscaler.

operator-sdk cleanup emr-dynamic-sizing

Se houver trabalhos em execução no cluster quando você excluir o operador, esses trabalhos continuarão em execução sem a escalabilidade automática vertical. Se você enviar trabalhos no cluster após a exclusão o operador, o Amazon EMR no EKS ignorará quaisquer parâmetros relacionados à escalabilidade automática vertical que você possa ter definido durante a configuração.

Desinstalação 248

Execução de workloads interativas no Amazon EMR no EKS

Um endpoint interativo corresponde a um gateway que conecta o Amazon EMR Studio ao Amazon EMR no EKS para que você possa executar workloads interativas. É possível usar endpoints interativos com o EMR Studio para executar analytics interativos com conjuntos de dados em armazenamentos de dados como o Amazon S3 e o Amazon DynamoDB.

Casos de uso

- Criação de um script de ETL com a experiência do IDE do EMR Studio. O IDE ingere dados onpremises e os armazena no Amazon S3 após as transformações para análises posteriores.
- Uso de cadernos para explorar conjuntos de dados e treinar um modelo de machine learning para detectar anomalias nos conjuntos de dados.
- Criação de scripts que geram relatórios diários para aplicações de análise, como painéis de negócios.

Tópicos

- Visão geral dos endpoints interativos
- Pré-requisitos para a criação de um endpoint interativo no Amazon EMR no EKS
- Criação de um endpoint interativo para o cluster virtual
- Definição de configurações para endpoints interativos
- Monitoramento de endpoints interativos
- Uso de cadernos Jupyter de hospedagem própria
- Como obter informações sobre endpoints interativos com comandos da CLI

Visão geral dos endpoints interativos

Um endpoint interativo fornece a funcionalidade para os clientes interativos, como os que usam o Amazon EMR Studio, de se conectarem ao Amazon EMR em clusters do EKS para executar workloads interativas. O endpoint interativo está respaldado pelo Jupyter Enterprise Gateway que fornece a funcionalidade de gerenciamento remoto do ciclo de vida do kernel de que os clientes interativos precisam. Os kernels são processos específicos de linguagem que interagem com o cliente do Amazon EMR Studio baseado em Jupyter para executar workloads interativas.

Os endpoints interativos oferecem suporte aos seguintes kernels:

- Python 3
- PySpark no Kubernetes
- Apache Spark com o Scala



Note

Os preços do Amazon EMR no EKS se aplicam aos endpoints e aos kernels interativos. Para obter mais informações, consulte a página Preços do Amazon EMR no EKS.

As entidades apresentadas a seguir são necessárias para que o EMR Studio se conecte ao Amazon EMR no EKS.

- Cluster virtual do Amazon EMR no EKS: um cluster virtual corresponde a um namespace do Kubernetes no qual você registra o Amazon EMR. O Amazon EMR usa clusters virtuais para executar trabalhos e hospedar endpoints. É possível fazer backup de vários clusters virtuais usando o mesmo cluster físico. No entanto, cada cluster virtual é mapeado para um namespace em um cluster do Amazon EKS. Os clusters virtuais não criam quaisquer recursos ativos que contribuam para o seu faturamento ou que requeiram gerenciamento do ciclo de vida de forma externa ao serviço.
- Endpoint interativo do Amazon EMR no EKS: um endpoint interativo corresponde a um endpoint HTTPS ao qual os usuários do EMR Studio podem conectar um Workspace. É possível acessar os endpoints HTTPS somente ao usar o EMR Studio e criá-los em uma sub-rede privada da Amazon Virtual Private Cloud (Amazon VPC) para o cluster do Amazon EKS.
 - Os kernels Python, PySpark, e Spark Scala usam as permissões definidas na função de execução de tarefas do Amazon EMR no EKS para invocar outras. Serviços da AWS Todos os kernels e os usuários que se conectam ao endpoint interativo utilizam o perfil que você especificou ao criar o endpoint. Recomendamos que você crie endpoints separados para usuários diferentes e que os usuários tenham funções diferentes AWS Identity and Access Management (IAM).
- AWS Controlador do Application Load Balancer O controlador do AWS Application Load Balancer gerencia o Elastic Load Balancing para um cluster Amazon EKS Kubernetes. O controlador provisiona um Application Load Balancer (ALB) quando você cria um recurso Ingress do Kubernetes. Um ALB revela um serviço do Kubernetes, como um endpoint interativo, de forma externa ao cluster do Amazon EKS, mas dentro da mesma Amazon VPC. Ao criar um endpoint interativo, também ocorre a implementação de um recurso Ingress que revela o endpoint

interativo por meio do ALB para os clientes interativos se conectarem. Você só precisa instalar um controlador do AWS Application Load Balancer para cada cluster do Amazon EKS.

O diagrama a seguir descreve a arquitetura dos endpoints interativos no Amazon EMR no EKS. Um cluster do Amazon EKS compreende a computação para executar as workloads de análise e o endpoint interativo. O controlador do Application Load Balancer é executado no namespace kube-system, enquanto as workloads e os endpoints interativos são executados no namespace especificado ao criar o cluster virtual. Quando você cria um endpoint interativo, o ambiente de gerenciamento do Amazon EMR no EKS cria a implantação do endpoint interativo no cluster do Amazon EKS. Além disso, uma instância da entrada do balanceador de carga do aplicativo é criada pelo controlador do balanceador de AWS carga. O Application Load Balancer fornece a interface externa para que os clientes, como o EMR Studio, se conectem ao cluster do Amazon EMR e executem workloads interativas.

Pré-requisitos para a criação de um endpoint interativo no Amazon EMR no EKS

Esta seção descreve os pré-requisitos para configurar um endpoint interativo que o EMR Studio poderá usar para se conectar a um cluster do Amazon EMR no EKS e executar workloads interativas.

AWS CLI

Siga as etapas em <u>Instalar ou atualizar para a versão mais recente do AWS CLI</u> para instalar a versão mais recente do AWS Command Line Interface (AWS CLI).

Instalação do eksctl

Siga as etapas descritas em <u>Instalar o kubectl</u> para instalar a versão mais recente do eksctl. Se você estiver usando a versão 1.22 ou versões posteriores do Kubernetes para o cluster do Amazon EKS, use uma versão do eksctl superior a 0.117.0.

Cluster do Amazon EKS

Crie um cluster do Amazon EKS. Registre o cluster como um cluster virtual com o Amazon EMR no EKS. Confira abaixo os requisitos e as considerações para este cluster:

- O cluster deve estar na mesma Amazon Virtual Private Cloud (VPC) que seu EMR Studio.
- O cluster deve ter, no mínimo, uma sub-rede privada para ativar os endpoints interativos, vincular repositórios baseados em Git e iniciar o Application Load Balancer no modo privado.
- Deve existir, no mínimo, uma sub-rede privada em comum entre o EMR Studio e o cluster do Amazon EKS usado para registrar o cluster virtual. Isso garante que seu endpoint interativo apareça como uma opção nos WorkSpaces do Studio e ativa a conectividade do Studio com o Application Load Balancer.

Existem dois métodos que você pode escolher para conectar o Studio e o cluster do Amazon EKS:

- Criar um cluster do Amazon EKS e associá-lo às sub-redes que pertencem ao seu EMR Studio.
- Como alternativa, crie um EMR Studio e especifique as sub-redes privadas para o cluster do Amazon EKS.
- ARM otimizado para Amazon EKS O Amazon Linux não é AMIs compatível com o Amazon EMR em endpoints interativos EKS.
- Os endpoints interativos funcionam com clusters do Amazon EKS que usam versões do Kubernetes até 1.30.
- Somente grupos de nós gerenciados do Amazon EKS são compatíveis.

Concessão de acesso ao cluster para o Amazon EMR no EKS

Use as etapas em Concessão de acesso ao cluster para o Amazon EMR no EKS para conceder acesso a um namespace específico em seu cluster para o Amazon EMR no EKS.

Ativação de IRSA no cluster do Amazon EKS

Para ativar os perfis do IAM para contas de serviço (IRSA) no cluster do Amazon EKS, siga as etapas em Habilitação de perfis do IAM para contas de serviço (IRSA).

Criação de um perfil de execução de trabalho do IAM

Você deve criar um perfil do IAM para executar workloads em endpoints interativos do Amazon EMR no EKS. Referimo-nos a esse perfil do IAM como perfil de execução de trabalho nesta documentação. Esse perfil do IAM é atribuído ao contêiner de endpoint interativo e aos contêineres de execução reais criados quando você envia trabalhos com o EMR Studio. Você precisará do nome

Concessão de acesso ao cluster 252

do recurso da Amazon (ARN) do seu perfil de execução de trabalho para o Amazon EMR no EKS. Para isso, são necessárias duas etapas:

- Crie um perfil do IAM para execução de trabalhos.
- Atualize a política de confiança do perfil de execução de trabalho.

Concessão de acesso ao Amazon EMR no EKS para os usuários

A entidade do IAM (usuário ou função) que faz a solicitação para criar um endpoint interativo também deve ter as seguintes emr-containers permissões EC2 e Amazon. Siga as etapas descritas em Concessão de acesso ao Amazon EMR no EKS para os usuários para conceder as permissões que permitem que o Amazon EMR no EKS crie, gerencie e exclua os grupos de segurança que limitam o tráfego de entrada para o balanceador de carga do seu endpoint interativo.

As seguintes permissões de emr-containers permitem que o usuário execute operações básicas do endpoint interativo:

```
"ec2:CreateSecurityGroup",
"ec2:DeleteSecurityGroup",
"ec2:AuthorizeSecurityGroupEgress",
"ec2:AuthorizeSecurityGroupIngress",
"ec2:RevokeSecurityGroupEgress",
"ec2:RevokeSecurityGroupIngress"

"emr-containers:CreateManagedEndpoint",
"emr-containers:ListManagedEndpoints",
"emr-containers:DescribeManagedEndpoint",
"emr-containers:DeleteManagedEndpoint"
```

Registro do cluster do Amazon EKS com o Amazon EMR

Configure um cluster virtual e mapeie-o para o namespace no cluster do Amazon EKS no qual você deseja executar os trabalhos. Para clusters AWS Fargate somente, use o mesmo namespace para o Amazon EMR no cluster virtual EKS e para o perfil Fargate.

Para obter informações sobre como configurar um cluster virtual do Amazon EMR no EKS, consulte Registro do cluster do Amazon EKS com o Amazon EMR.

Implemente AWS o Load Balancer Controller no cluster Amazon EKS

É necessário um AWS Application Load Balancer para seu cluster Amazon EKS. Você precisa configurar somente um controlador do Application Load Balancer por cluster do Amazon EKS. Para obter informações sobre como configurar o controlador do AWS Application Load Balancer, consulte Instalando o complemento Load AWS Balancer Controller no Guia do usuário do Amazon EKS.

Criação de um endpoint interativo para o cluster virtual

Este tópico descreve algumas maneiras de criar um endpoint interativo usando a interface de linha de AWS comando (AWS CLI) e inclui detalhes sobre os parâmetros de configuração disponíveis.

Criação de um endpoint interativo com o comando create-managedendpoint

Especifique os parâmetros no comando create-managed-endpoint, conforme apresentado a seguir. O Amazon EMR no EKS oferece suporte à criação de endpoints interativos com as versões 6.7.0 e superiores do Amazon EMR.

```
aws emr-containers create-managed-endpoint \
--type JUPYTER_ENTERPRISE_GATEWAY \
--virtual-cluster-id 1234567890abcdef0xxxxxxxx \
--name example-endpoint-name \
--execution-role-arn arn:aws:iam::444455556666:role/JobExecutionRole \
--release-label emr-6.9.0-latest \
--configuration-overrides '{
    "applicationConfiguration": [{
        "classification": "spark-defaults",
        "properties": {
            "spark.driver.memory": "2G"
        }
    }],
    "monitoringConfiguration": {
        "cloudWatchMonitoringConfiguration": {
            "logGroupName": "log_group_name",
            "logStreamNamePrefix": "log_stream_prefix"
        },
        "persistentAppUI": "ENABLED",
        "s3MonitoringConfiguration": {
            "logUri": "s3://my_s3_log_location"
```

Load Balancer Controller 254

```
}
}'
```

Para obter mais informações, consulte Parâmetros para a criação de um endpoint interativo.

Criação de um endpoint interativo com parâmetros especificados em um arquivo JSON

1. Crie um arquivo create-managed-endpoint-request.json e especifique os parâmetros obrigatórios para o endpoint, conforme mostrado no seguinte arquivo JSON:

```
{
    "name": "MY_TEST_ENDPOINT",
    "virtualClusterId": "MY_CLUSTER_ID",
    "type": "JUPYTER_ENTERPRISE_GATEWAY",
    "releaseLabel": "emr-6.9.0-latest",
    "executionRoleArn": "arn:aws:iam::444455556666:role/JobExecutionRole",
    "configurationOverrides":
    {
        "applicationConfiguration":
        {
                "classification": "spark-defaults",
                "properties":
                     "spark.driver.memory": "8G"
            }
        ],
        "monitoringConfiguration":
        {
            "persistentAppUI": "ENABLED",
            "cloudWatchMonitoringConfiguration":
            {
                "logGroupName": "my_log_group",
                "logStreamNamePrefix": "log_stream_prefix"
            },
            "s3MonitoringConfiguration":
                "logUri": "s3://my_s3_log_location"
            }
```

```
}
}
```

2. Use o comando create-managed-endpoint com um caminho para o arquivo create-managed-endpoint-request.json armazenado localmente ou no Amazon S3.

```
aws emr-containers create-managed-endpoint \
--cli-input-json file://./create-managed-endpoint-request.json --region AWS-Region
```

Saída da criação de endpoint interativo

Você deverá visualizar a saída a seguir no terminal. A saída inclui o nome e o identificador do seu novo endpoint interativo:

```
{
    "id": "1234567890abcdef0",
    "name": "example-endpoint-name",
    "arn": "arn:aws:emr-containers:us-west-2:111122223333:/
virtualclusters/444455556666/endpoints/444455556666",
    "virtualClusterId": "111122223333xxxxxxxxxx"
}
```

A execução de aws emr-containers create-managed-endpoint cria um certificado autoassinado que permite a comunicação HTTPS entre o EMR Studio e o servidor do endpoint interativo.

Se você executar create-managed-endpoint e não tiver concluído os pré-requisitos, o Amazon EMR retornará uma mensagem de erro com as ações que você deve realizar para continuar.

Parâmetros para a criação de um endpoint interativo

Tópicos

- Parâmetros obrigatórios para endpoints interativos
- Parâmetros opcionais para endpoints interativos

Parâmetros obrigatórios para endpoints interativos

Você deve especificar os seguintes parâmetros ao criar um endpoint interativo:

--type

Use JUPYTER_ENTERPRISE_GATEWAY. Este é o único tipo com suporte.

--virtual-cluster-id

O identificador do cluster virtual registrado com o Amazon EMR no EKS.

--name

Um nome descritivo para o endpoint interativo que auxilia os usuários do EMR Studio a selecionálo na lista suspensa.

--execution-role-arn

O nome do recurso da Amazon (ARN) do perfil de execução de trabalho do IAM para o Amazon EMR no EKS que foi criado como parte dos pré-requisitos.

--release-label

O rótulo da versão do Amazon EMR a ser usado para o endpoint. Por exemplo, .emr-6.9.0-latest O Amazon EMR no EKS oferece suporte a endpoints interativos com as versões 6.7.0 e superiores do Amazon EMR.

Parâmetros opcionais para endpoints interativos

Como opção, também é possível especificar os seguintes parâmetros ao criar um endpoint interativo:

--configuration-overrides

Para substituir as configurações padrão das aplicações, forneça um objeto de configuração. Você pode usar uma sintaxe abreviada para fornecer a configuração ou fazer referência ao objeto de configuração em um arquivo JSON.

Os objetos de configuração consistem em uma classificação, propriedades e configurações opcionais aninhadas As propriedades consistem nas configurações que você deseja substituir neste arquivo. Você pode especificar várias classificações para diversas aplicações em um único objeto JSON. As classificações de configuração disponíveis variam de acordo com a versão do Amazon EMR no EKS. Para obter uma lista das classificações de configuração disponíveis para cada versão do Amazon EMR no EKS, consulte Versões do Amazon EMR no EKS. Além das classificações de configuração listadas para cada versão, os endpoints interativos trazem a classificação adicional

jeg-config. Para obter mais informações, consulte <u>Opções de configuração do Jupyter Enterprise</u> Gateway (JEG).

Definição de configurações para endpoints interativos

Esta seção contém uma série de tópicos que abrangem várias configurações para endpoints interativos e configurações de pod. Eles permitem monitorar e solucionar falhas, enviar informações de log para o Amazon S3 ou Amazon CloudWatch Logs para, ou criar endpoints interativos nos quais você especifica modelos de pod personalizados.

Tópicos

- Monitorar trabalhos Spark do
- Especificação de modelos de pod personalizados com endpoints interativos
- Implantação de um pod do JEG em um grupo de nós
- Opções de configuração do Jupyter Enterprise Gateway (JEG)
- Modificando os parâmetros PySpark da sessão
- Imagem de kernel personalizada com o endpoint interativo

Monitorar trabalhos Spark do

Para que você possa monitorar e solucionar falhas, configure seus endpoints interativos para que os trabalhos iniciados com o endpoint possam enviar informações de log para o Amazon S3, Amazon CloudWatch Logs ou ambos. As seções a seguir descrevem como enviar logs de aplicações do Spark para o Amazon S3 para os trabalhos do Spark executados com endpoints interativos do Amazon EMR no EKS.

Configuração da política do IAM para os logs do Amazon S3

Antes que seus kernels possam enviar dados de log ao Amazon S3, a política de permissões para o perfil de execução de trabalho deve incluir as permissões apresentadas a seguir. Substitua amzn-s3-demo-destination-bucket pelo nome do bucket de registro em log.

JSON

{

```
"Version": "2012-10-17",
  "Statement": [
      "Effect": "Allow",
      "Action": [
        "s3:PutObject",
        "s3:GetObject",
        "s3:ListBucket"
      ],
      "Resource": [
        "arn:aws:s3:::amzn-s3-demo-bucket",
        "arn:aws:s3:::amzn-s3-demo-bucket/*"
      ],
      "Sid": "AllowS3Putobject"
    }
  ]
}
```

Note

O Amazon EMR no EKS também pode criar um bucket do S3. Se um bucket do S3 não estiver disponível, inclua a permissão s3:CreateBucket na política do IAM.

Após conceder as permissões necessárias para o envio de logs ao bucket do S3 ao perfil de execução, os dados de log serão enviados para os locais do Amazon S3 apresentados a seguir. Isso acontece quando s3MonitoringConfiguration é transferido na seção monitoringConfiguration de uma solicitação create-managed-endpoint.

- Logs de driver: logUri/virtual-cluster-id/endpoints/endpoint-id/containers/ spark-application-id/spark-application-id-driver/(stderr.gz/stdout.gz).
- Logs de executor: logUri/virtual-cluster-id/endpoints/endpoint-id/containers/ spark-application-id/executor-pod-name-exec-<Number>/(stderr.gz/ stdout.gz).

Note

O Amazon EMR no EKS não faz o upload dos logs do endpoint para o bucket do S3.

Monitorar trabalhos Spark 259

Especificação de modelos de pod personalizados com endpoints interativos

É possível criar endpoints interativos nos quais você especifica modelos de pod personalizados para drivers e executores. Os modelos de pod são especificações que determinam como ocorrerá a execução de cada pod. Você pode usar arquivos de modelo de pod para definir as configurações de pods de drivers ou de executores para as quais as configurações do Spark não oferecem suporte. No momento, os modelos de pod são compatíveis com as versões 6.3.0 e posteriores do Amazon EMR.

Para obter mais informações sobre modelos de pod, consulte <u>Using pod templates</u> no Guia de desenvolvimento do Amazon EMR no EKS.

O seguinte exemplo mostra como criar um endpoint interativo com modelos de pod:

```
aws emr-containers create-managed-endpoint \
    --type JUPYTER_ENTERPRISE_GATEWAY \
    --virtual-cluster-id virtual-cluster-id \
    --name example-endpoint-name \
    --execution-role-arn arn:aws:iam::aws-account-id:role/EKSClusterRole \
    --release-label emr-6.9.0-latest \
    --configuration-overrides '{
        "applicationConfiguration": [
        {
            "classification": "spark-defaults",
            "properties": {
                "spark.kubernetes.driver.podTemplateFile": "path/to/driver/
template.yaml",
                "spark.kubernetes.executor.podTemplateFile": "path/to/executor/
template.yaml"
        }]
    }'
```

Implantação de um pod do JEG em um grupo de nós

O posicionamento do pod do JEG (Jupyter Enterprise Gateway) é um recurso que permite implantar um endpoint interativo em um grupo de nós específico. Com esse recurso, você pode definir configurações, como instance type, para o endpoint interativo.

Associação de um pod do JEG a um grupo de nós gerenciado

A propriedade de configuração apresentada a seguir permite especificar o nome de um grupo de nós gerenciado no cluster do Amazon EKS em que o pod do JEG será implantado.

Um grupo de nós deve ter o rótulo for-use-with-emr-containers-managed-endpointng=*NodeGroupName* do Kubernetes anexado a todos os nós que fazem parte do grupo de nós. Para listar todos os nós de um grupo de nós que têm essa etiqueta, use o seguinte comando:

```
kubectl get nodes --show-labels | grep for-use-with-emr-containers-managed-endpoint-
ng=NodeGroupName
```

Se a saída do comando acima não retornar nós que fazem parte do seu grupo de nós gerenciado, então não há nós no grupo de nós que tenham o rótulo for-use-with-emr-containers-managed-endpoint-ng=NodeGroupName do Kubernetes anexado. Nesse caso, siga as etapas abaixo para anexar esse rótulo aos nós do seu grupo de nós.

1. Use o comando a seguir para adicionar o rótulo for-use-with-emr-containers-managed-endpoint-ng=NodeGroupName do Kubernetes a todos os nós em um grupo de nós gerenciado NodeGroupName:

```
kubectl label nodes --selector eks:nodegroup-name=NodeGroupName for-use-with-emr-
containers-managed-endpoint-ng=NodeGroupName
```

2. Verifique se os nós foram rotulados corretamente usando o seguinte comando:

```
kubectl get nodes --show-labels | grep for-use-with-emr-containers-managed-endpoint-
ng=NodeGroupName
```

Um grupo de nós gerenciado deve estar associado a um grupo de segurança do cluster do Amazon EKS, o que geralmente acontece se você criou o cluster e o grupo de nós gerenciado usando eksctl. Você pode verificar isso no AWS console usando as etapas a seguir.

- Acesse o cluster no console do Amazon EKS.
- 2. Acesse a guia de redes do cluster e anote o grupo de segurança do cluster.
- 3. Acesse a guia de computação do seu cluster e clique no nome do grupo de nós gerenciado.
- 4. Na guia Detalhes do grupo de nós gerenciado, verifique se o grupo de segurança do cluster anotado anteriormente está listado em Grupos de segurança.

Se o grupo de nós gerenciado não estiver anexado ao grupo de segurança do cluster do Amazon EKS, será necessário anexar a etiqueta for-use-with-emr-containers-managed-endpoint-sg=ClusterName/NodeGroupName ao grupo de segurança do grupo de nós. Use as etapas abaixo para anexar essa etiqueta.

- Acesse o EC2 console da Amazon e clique em grupos de segurança no painel de navegação esquerdo.
- 2. Selecione o grupo de segurança do seu grupo de nós gerenciado ao clicar na caixa de seleção.
- 3. Na guia Etiquetas, adicione a etiqueta for-use-with-emr-containers-managed-endpoint-sg=ClusterName/NodeGroupName usando o botão Gerenciar etiquetas.

Associação de um pod do JEG a um grupo de nós autogerenciado

A propriedade de configuração apresentada a seguir permite especificar o nome de um grupo de nós autogerenciado ou não gerenciado no cluster do Amazon EKS em que o pod do JEG será implantado.

O grupo de nós deve ter o rótulo for-use-with-emr-containers-managed-endpointng=*NodeGroupName* do Kubernetes anexado a todos os nós que fazem parte do grupo de nós. Para listar todos os nós de um grupo de nós que têm essa etiqueta, use o seguinte comando:

Se a saída do comando acima não retornar nós que fazem parte do seu grupo de nós autogerenciado, então não há nós no grupo de nós que tenham o rótulo for-use-with-emr-containers-managed-endpoint-ng=NodeGroupName do Kubernetes anexado. Nesse caso, siga as etapas abaixo para anexar esse rótulo aos nós do seu grupo de nós.

1. Se você criou o grupo de nós autogerenciado usando eksctl, use o comando a seguir para adicionar o rótulo for-use-with-emr-containers-managed-endpointng=NodeGroupName do Kubernetes a todos os nós no grupo de nós autogerenciado NodeGroupName de uma vez só.

```
\label nodes \ --selector \ alpha. eksctl.io/nodegroup-name= {\it NodeGroupName} \ for-use-with-emr-containers-managed-endpoint-ng= {\it NodeGroupName} \
```

Caso não tenha usado eksctl para criar o grupo de nós autogerenciado, você precisará substituir o seletor no comando acima por um rótulo diferente do Kubernetes que esteja anexado a todos os nós do grupo de nós.

2. Use o seguinte comando para verificar se os nós foram rotulados corretamente:

```
\label{localization} $$ kubectl get nodes --show-labels \mid grep for-use-with-emr-containers-managed-endpoint-ng=$$ NodeGroupName $$ $$
```

O grupo de segurança do grupo de nós autogerenciado deve ter a etiqueta for-use-with-emr-containers-managed-endpoint-sg=*ClusterName/NodeGroupName* anexada. Use as etapas apresentadas a seguir para anexar a etiqueta ao grupo de segurança do AWS Management Console.

- Navegue até o EC2 console da Amazon. Selecione Grupos de segurança no painel de navegação à esquerda.
- 2. Selecione a caixa de seleção ao lado do grupo de segurança do seu grupo de nós autogerenciado.

3. Na guia Etiquetas, use o botão Gerenciar etiquetas para adicionar a etiqueta for-use-withemr-containers-managed-endpoint-sg=ClusterName/NodeGroupName. Substitua ClusterName e NodeGroupName com os valores apropriados.

Associação de um pod do JEG a um grupo de nós gerenciado com instâncias sob demanda

Você também pode definir rótulos adicionais, conhecidos como seletores de rótulos do Kubernetes, para especificar limitações ou restrições adicionais à execução de um endpoint interativo em um determinado nó ou grupo de nós. O exemplo a seguir mostra como usar EC2 instâncias Amazon sob demanda para um pod JEG.

Note

Você pode usar a propriedade node-labels somente com uma propriedade managed-nodegroup-name ou self-managed-nodegroup-name.

Opções de configuração do Jupyter Enterprise Gateway (JEG)

O Amazon EMR no EKS usa o Jupyter Enterprise Gateway (JEG) para ativar os endpoints interativos. É possível definir os valores a seguir para as configurações do JEG listadas como permitidas ao criar o endpoint.

 RemoteMappingKernelManager.cull_idle_timeout: tempo limite em segundos (inteiro), após o qual um kernel é considerado inativo e pronto para ser descartado. Valores iguais ou inferiores a 0 desativam o descarte. Tempos limite curtos podem resultar no descarte de kernels para usuários com conexões de rede ruins.

• RemoteMappingKernelManager.cull_interval: o intervalo em segundos (inteiro) no qual verificar se há kernels inativos que excedem o valor do tempo limite de descarte.

Modificando os parâmetros PySpark da sessão

Começando com o Amazon EMR no EKS versão 6.9.0, no Amazon EMR Studio você pode ajustar a configuração do Spark associada a uma PySpark sessão executando o comando %%configure mágico na célula do notebook EMR.

O exemplo a seguir mostra uma carga útil de exemplo que você pode usar para modificar a memória, os núcleos e outras propriedades do driver e do executor do Spark. Para as configurações conf, você pode definir qualquer configuração do Spark mencionada na documentação de configuração do Apache Spark.

```
%configure -f
{
   "driverMemory": "16G",
   "driverCores": 4,
   "executorMemory": "32G",
   "executorCores": 2,
   "conf": {
        "spark.dynamicAllocation.maxExecutors": 10,
        "spark.dynamicAllocation.minExecutors": 1
   }
}
```

O exemplo a seguir mostra uma carga útil de exemplo que você pode usar para adicionar arquivos, pyFiles e dependências em JAR a um runtime do Spark.

```
%configure -f
{
   "files": "s3://amzn-s3-demo-bucket-emr-eks/sample_file.txt",
   "pyFiles": : "path-to-python-files",
   "jars" : "path-to-jars
}
```

Imagem de kernel personalizada com o endpoint interativo

Para garantir que você tenha as dependências corretas para sua aplicação ao executar as workloads interativas do Amazon EMR Studio, você pode personalizar as imagens do Docker para os endpoints interativos e executar imagens base de kernel personalizadas. Para criar um endpoint interativo e conectá-lo a uma imagem do Docker personalizada, execute as etapas a seguir.



Note

Você pode substituir somente as imagens base. Não é possível adicionar novos tipos de imagens de kernel.

Crie e publique uma imagem do Docker personalizada. A imagem base contém o runtime do Spark e os kernels do caderno que são executados com ele. Para criar a imagem, você pode seguir as etapas de 1 a 4 descritas em Como personalizar imagens do Docker. Na etapa 1, o URI da imagem base em seu arquivo do Docker deve usar notebook-spark no lugar de spark.

```
ECR-registry-account.dkr.ecr.Region.amazonaws.com/notebook-spark/container-image-
tag
```

Para obter mais informações sobre como selecionar Regiões da AWS e armazenar tags de imagem, consulteDetalhes sobre como selecionar um URI de imagem base.

- 2. Crie um endpoint interativo que possa ser usado com a imagem personalizada.
 - Crie um arquivo JSON custom-image-managed-endpoint. json com o conteúdo apresentado a seguir. Este exemplo usa a versão 6.9.0 do Amazon EMR.

Example

```
{
    "name": "endpoint-name",
    "virtualClusterId": "virtual-cluster-id",
    "type": "JUPYTER_ENTERPRISE_GATEWAY",
    "releaseLabel": "emr-6.9.0-latest",
    "executionRoleArn": "execution-role-arn",
    "configurationOverrides": {
        "applicationConfiguration": [
```

```
{
                "classification": "jupyter-kernel-overrides",
                "configurations": [
                    {
                         "classification": "python3",
                         "properties": {
                             "container-image": "123456789012.dkr.ecr.us-
west-2.amazonaws.com/custom-notebook-python:latest"
                    },
                    {
                         "classification": "spark-python-kubernetes",
                         "properties": {
                             "container-image": "123456789012.dkr.ecr.us-
west-2.amazonaws.com/custom-notebook-spark:latest"
                    }
                ]
            }
        ]
    }
}
```

 b. Crie um endpoint interativo com as configurações especificadas no arquivo JSON, conforme mostrado no exemplo a seguir. Para obter mais informações, consulte <u>Criação de um</u> endpoint interativo com o comando create-managed-endpoint.

```
aws emr-containers create-managed-endpoint --cli-input-json custom-image-managed-endpoint.json
```

 Conecte-se ao endpoint interativo usando o EMR Studio. Para obter mais informações e etapas a serem concluídas, consulte <u>Conectando-se a partir do Studio</u> na seção Amazon EMR no EKS dos documentos do AWS Workshop Studio.

Monitoramento de endpoints interativos

Com o Amazon EMR no EKS versão 6.10 e posterior, os endpoints interativos emitem métricas da CloudWatch Amazon para monitorar e solucionar problemas nas operações do ciclo de vida do kernel. As métricas são acionadas por clientes interativos, como o EMR Studio ou os cadernos Jupyter de hospedagem própria. Cada uma das operações compatíveis com os endpoints interativos tem métricas associadas a elas. As operações são modeladas como dimensões para cada métrica,

conforme mostrado na tabela abaixo. As métricas emitidas por endpoints interativos são visíveis em um namespace personalizado, EMRContainers, em sua conta.

Métrica	Descrição	Unidade
RequestCount	Número cumulativo de solicitações de uma operação processada pelo endpoint interativo.	Contagem
RequestLatency	O horário entre uma solicitaç ão chegar ao endpoint interativo e uma resposta ser enviada pelo endpoint interativ o.	Milissegundo
4 XXError	Emitido quando uma solicitaç ão de uma operação resulta em um erro 4xx durante o processamento.	Contagem
5XXError	Emitido quando uma solicitaç ão de uma operação resulta em um erro 5Xxx no lado do servidor.	Contagem
KernelLaunchSuccess	Aplicável somente para a CreateKernel operação. Indica o número cumulativo de inicializações de kernel que ocorreram com êxito até e incluindo esta solicitação.	Contagem
KernelLaunchFailure	Aplicável somente para a CreateKernel operação. Indica o número cumulativo de falhas de inicialização de kernel até e incluindo esta solicitação.	Contagem

Cada métrica do endpoint interativo tem as seguintes dimensões anexadas a ela:

- ManagedEndpointId: identificador para o endpoint interativo.
- OperationName: a operação acionada pelo cliente interativo.

Os valores possíveis para a dimensão **OperationName** são mostrados na seguinte tabela:

operationName	Descrição da operação
CreateKernel	Solicita que o endpoint interativo inicie um kernel.
ListKernels	Solicita que o endpoint interativo liste os kernels que foram iniciados anteriormente usando o mesmo token de sessão.
GetKernel	Solicita que o endpoint interativo obtenha detalhes sobre um kernel específico que foi iniciado anteriormente.
ConnectKernel	Solicita que o endpoint interativo estabeleça conectividade entre o cliente do caderno e o kernel.
ConfigureKernel	Publica %%configure magic request em um kernel do PySpark.
ListKernelSpecs	Solicita que o endpoint interativo liste as especificações de kernel disponíveis.
GetKernelSpec	Solicita que o endpoint interativo obtenha as especificações de kernel de um kernel que foi iniciado anteriormente.
GetKernelSpecResource	Solicita que o endpoint interativo obtenha recursos específicos associados às específic ações do kernel que foram iniciadas anteriorm ente.

Exemplos

Para acessar o número total de kernels iniciados para um endpoint interativo em um determinado dia:

- 1. Selecione o namespace personalizado: EMRContainers.
- 2. Selecione o ManagedEndpointId e OperationName CreateKernel.
- 3. A métrica RequestCount com a estatística SUM e o período de 1 day fornecerá todas as solicitações de inicialização do kernel realizadas nas últimas 24 horas.
- 4. KernelLaunchSuccess A métrica com estatística SUM e período 1 day fornecerá todas as solicitações bem-sucedidas de inicialização do kernel feitas nas últimas 24 horas.

Para acessar o número de falhas de kernel para um endpoint interativo em um determinado dia:

- 1. Selecione o namespace personalizado: EMRContainers
- Selecione o ManagedEndpointId e OperationName CreateKernel.
- 3. A métrica KernelLaunchFailure com a estatística SUM e o período de 1 day fornecerá todas as solicitações de inicialização do kernel com falha realizadas nas últimas 24 horas. Você também pode selecionar as métricas 4XXError e 5XXError para saber que tipo de falha na inicialização do kernel ocorreu.

Uso de cadernos Jupyter de hospedagem própria

Você pode hospedar e gerenciar o Jupyter ou JupyterLab notebooks em uma EC2 instância da Amazon ou em seu próprio cluster Amazon EKS como um notebook Jupyter auto-hospedado. Em seguida, é possível executar workloads interativas com seus cadernos Jupyter de hospedagem própria. As seções apresentadas a seguir descrevem o processo de configuração e de implantação de um caderno Jupyter de hospedagem própria em um cluster do Amazon EKS.

Criação de um caderno Jupyter de hospedagem própria em um cluster do EKS

- Criar um grupo de segurança
- Criação de um endpoint interativo do Amazon EMR no EKS
- Recuperação do URL do servidor de gateway do endpoint interativo

Exemplos 270

- Recuperação de um token de autenticação para a conexão com o endpoint interativo
- Exemplo: implantar um JupyterLab notebook
- Exclusão de um caderno Jupyter de hospedagem própria

Criar um grupo de segurança

Antes de criar um endpoint interativo e executar um Jupyter ou JupyterLab notebook autohospedado, você deve criar um grupo de segurança para controlar o tráfego entre seu notebook
e o endpoint interativo. Para usar o EC2 console da Amazon ou o Amazon EC2 SDK para criar o
grupo de segurança, consulte as etapas em <u>Criar um grupo de segurança</u> no Guia do EC2 usuário
da Amazon. Você deve criar o grupo de segurança na VPC em que deseja implantar seu servidor de
cadernos.

Para seguir o exemplo deste guia, use a mesma VPC do cluster do Amazon EKS. Se você quiser hospedar seu notebook em uma VPC diferente da VPC do seu cluster Amazon EKS, talvez seja necessário criar uma conexão de emparelhamento entre essas duas. VPCs Para ver as etapas para criar uma conexão de emparelhamento entre duas VPCs, consulte Criar uma conexão de emparelhamento de VPC no Amazon VPC Getting Started Guide.

Você precisará do ID do grupo de segurança para <u>criar um endpoint interativo do Amazon EMR no</u> EKS na próxima etapa.

Criação de um endpoint interativo do Amazon EMR no EKS

Após criar o grupo de segurança para o seu caderno, use as etapas fornecidas em <u>Criação de um</u> <u>endpoint interativo para o cluster virtual</u> para criar um endpoint interativo. Você deve fornecer o ID do grupo de segurança criado para o seu caderno em <u>Criar um grupo de segurança</u>.

Insira o ID de segurança no lugar das *your-notebook-security-group-id* seguintes configurações de substituição de configuração:

```
--configuration-overrides '{
    "applicationConfiguration": [
    {
        "classification": "endpoint-configuration",
        "properties": {
            "notebook-security-group-id": "your-notebook-security-group-id"
        }
```

Criar um grupo de segurança 271

```
}
],
"monitoringConfiguration": {
....'
```

Recuperação do URL do servidor de gateway do endpoint interativo

Após criar um endpoint interativo, recupere o URL do servidor de gateway usando o comando describe-managed-endpoint na AWS CLI. Você precisa desse URL para conectar seu caderno ao endpoint. O URL do servidor de gateway é um endpoint privado.

```
aws emr-containers describe-managed-endpoint \
--region region \
--virtual-cluster-id virtualClusterId \
--id endpointId
```

Inicialmente, seu endpoint está no estado CREATING. Após alguns minutos, ele passa para o estado ACTIVE. Quando o endpoint estiver ACTIVE, ele estará pronto para uso.

Anote o atributo serverUrl que o comando aws emr-containers describe-managedendpoint retorna do endpoint ativo. Você precisa dessa URL para conectar seu notebook ao endpoint ao implantar seu Jupyter ou notebook auto-hospedado. JupyterLab

Recuperação de um token de autenticação para a conexão com o endpoint interativo

Para se conectar a um endpoint interativo a partir de um Jupyter ou JupyterLab notebook, você deve gerar um token de sessão com a API. GetManagedEndpointSessionCredentials O token atua como prova de autenticação para a conexão com o servidor de endpoint interativo.

O comando apresentado a seguir é explicado com mais detalhes com um exemplo de saída abaixo.

```
aws emr-containers get-managed-endpoint-session-credentials \
--endpoint-identifier endpointArn \
--virtual-cluster-identifier virtualClusterArn \
--execution-role-arn executionRoleArn \
--credential-type "TOKEN" \
--duration-in-seconds durationInSeconds \
--region region
```

endpointArn

O ARN do seu endpoint. Você pode descobrir o ARN no resultado de uma chamada describemanaged-endpoint.

virtualClusterArn

O ARN do cluster virtual.

executionRoleArn

O ARN do perfil de execução.

durationInSeconds

A duração, em segundos, pela qual o token é válido. A duração padrão é de 15 minutos (900) e a duração máxima é de 12 horas (43200).

region

A mesma região do seu endpoint.

A saída deve ser semelhante ao exemplo apresentado a seguir. Anote o *session-token* valor que você usará ao implantar seu Jupyter ou notebook auto-hospedado. JupyterLab

```
{
    "id": "credentialsId",
    "credentials": {
        "token": "session-token"
    },
    "expiresAt": "2022-07-05T17:49:38Z"
}
```

Exemplo: implantar um JupyterLab notebook

Depois de concluir as etapas acima, você pode tentar este procedimento de exemplo para implantar um JupyterLab notebook no cluster Amazon EKS com seu endpoint interativo.

- 1. Crie um namespace para executar o servidor de cadernos.
- 2. Crie um arquivo localmente, chamado notebook.yam1, com o conteúdo apresentado a seguir. O conteúdo para o arquivo é descrito abaixo.

```
apiVersion: v1
```

```
kind: Pod
metadata:
  name: jupyter-notebook
  namespace: namespace
spec:
  containers:
  - name: minimal-notebook
    image: jupyter/all-spark-notebook:lab-3.1.4 # open source image
    ports:
    - containerPort: 8888
    command: ["start-notebook.sh"]
   args: ["--LabApp.token=''"]
    env:
    - name: JUPYTER_ENABLE_LAB
      value: "yes"
    - name: KERNEL_LAUNCH_TIMEOUT
     value: "400"
    - name: JUPYTER_GATEWAY_URL
      value: "serverUrl"
    - name: JUPYTER_GATEWAY_VALIDATE_CERT
      value: "false"
    - name: JUPYTER_GATEWAY_AUTH_TOKEN
      value: "session-token"
```

Se você estiver implantando o caderno Jupyter em um cluster somente do Fargate, rotule o pod do Jupyter com um rótulo role, conforme mostrado no seguinte exemplo:

```
metadata:
   name: jupyter-notebook
   namespace: default
   labels:
     role: example-role-name-label
spec:
   ...
```

namespace

O namespace do Kubernetes no qual o caderno é implantado.

serverUrl

O atributo serverUrl que o comando describe-managed-endpoint retornou em Recuperação do URL do servidor de gateway do endpoint interativo .

session-token

O atributo session-token que o comando get-managed-endpoint-sessioncredentials retornou em Recuperação de um token de autenticação para a conexão com o endpoint interativo.

KERNEL_LAUNCH_TIMEOUT

A quantidade de tempo, em segundos, que o endpoint interativo aguarda até que o kernel chegue ao estado RUNNING. Garanta tempo suficiente para a inicialização do kernel ser concluída ao definir o tempo limite de inicialização do kernel para um valor apropriado (máximo de 400 segundos).

KERNEL_EXTRA_SPARK_OPTS

Como opção, você pode transferir configurações adicionais do Spark para os kernels do Spark. Defina essa variável de ambiente com os valores da propriedade de configuração do Spark, conforme mostrado no seguinte exemplo:

3. Implante a especificação do pod no cluster do Amazon EKS:

```
kubectl apply -f notebook.yaml -n namespace
```

Isso iniciará um JupyterLab notebook mínimo conectado ao seu Amazon EMR no endpoint interativo EKS. Aguarde até que o pod esteja RUNNING. Você pode verificar seu status com o seguinte comando:

```
kubectl get pod jupyter-notebook -n namespace
```

Quando o pod estiver pronto, o comando get pod retornará uma saída semelhante a esta:

```
NAME READY STATUS RESTARTS AGE jupyter-notebook 1/1 Running 0 46s
```

- 4. Anexe o grupo de segurança do caderno ao nó em que o caderno está programado.
 - a. Primeiro, identifique o nó em que o pod jupyter-notebook está programado com o comando describe pod.

```
kubectl describe pod jupyter-notebook -n namespace
```

- b. Abra o console do Amazon EKS em https://console.aws.amazon.com/eks/home#/clusters.
- c. Navegue até a guia Computação do cluster do Amazon EKS e selecione o nó identificado pelo comando describe pod. Selecione o ID da instância para o nó.
- d. No menu Ações, selecione Segurança > Alterar grupos de segurança para anexar o grupo de segurança que você criou em Criar um grupo de segurança.
- e. Se você estiver implantando o pod do notebook Jupyter AWS Fargate, crie um SecurityGroupPolicypara aplicar ao pod do notebook Jupyter com o rótulo de função:

```
cat >my-security-group-policy.yaml <<EOF
apiVersion: vpcresources.k8s.aws/v1beta1
kind: SecurityGroupPolicy
metadata:
   name: example-security-group-policy-name
   namespace: default
spec:
   podSelector:
    matchLabels:
      role: example-role-name-label
securityGroups:
   groupIds:
      - your-notebook-security-group-id</pre>
```

E0F

5. Agora, encaminhe a porta para que você possa acessar localmente a JupyterLab interface:

```
kubectl port-forward jupyter-notebook 8888:8888 -n namespace
```

Quando estiver em execução, navegue até o navegador local e visite localhost:8888 para ver a JupyterLab interface:

6. JupyterLabEm, crie um novo notebook Scala. Confira um exemplo de trecho de código que você pode executar para aproximar o valor de Pi:

```
import scala.math.random
import org.apache.spark.sql.SparkSession
/** Computes an approximation to pi */
val session = SparkSession
  .builder
  .appName("Spark Pi")
  .getOrCreate()
val slices = 2
// avoid overflow
val n = math.min(100000L * slices, Int.MaxValue).toInt
val count = session.sparkContext
.parallelize(1 until n, slices)
.map { i =>
 val x = random * 2 - 1
 val y = random * 2 - 1
 if (x*x + y*y <= 1) 1 else 0
}.reduce(_ + _)
println(s"Pi is roughly ${4.0 * count / (n - 1)}")
session.stop()
```

Exclusão de um caderno Jupyter de hospedagem própria

Quando estiver com tudo pronto para excluir seu caderno de hospedagem própria, você também poderá excluir o endpoint interativo e o grupo de segurança. Execute as ações na seguinte ordem:

1. Use o seguinte comando para excluir o pod jupyter-notebook:

```
kubectl delete pod jupyter-notebook -n namespace
```

- 2. Em seguida, exclua o endpoint interativo com o comando delete-managed-endpoint. Para obter as etapas para a exclusão de um endpoint interativo, consulte <u>Exclusão de um endpoint interativo</u>. Inicialmente, seu endpoint estará no estado TERMINATING. Depois que todos os recursos forem limpos, ele transitará para o estado TERMINATED.
- Caso não planeje usar o grupo de segurança de cadernos criado em <u>Criar um grupo de</u> <u>segurança</u> para outras implantações de caderno Jupyter, você poderá excluí-lo. Consulte <u>Excluir</u> <u>um grupo de segurança</u> no Guia EC2 do usuário da Amazon para obter mais informações.

Como obter informações sobre endpoints interativos com comandos da CLI

Este tópico aborda as operações com suporte em um endpoint interativo diferente de <u>create-</u>managed-endpoint.

Busca de detalhes do endpoint interativo

Depois de criar um endpoint interativo, você pode recuperar seus detalhes usando o describemanaged-endpoint AWS CLI comando. Insira seus próprios valores para managed-endpointidvirtual-cluster-id, eregion:

```
aws emr-containers describe-managed-endpoint --id managed-endpoint-id \
--virtual-cluster-id virtual-cluster-id --region region
```

A saída é semelhante à apresentada a seguir, com o endpoint especificado, como o ARN, o ID e o nome.

```
{
  "id": "as3ys2xxxxxxx",
  "name": "endpoint-name",
```

Limpeza 278

```
"arn": "arn:aws:emr-containers:us-east-1:1828xxxxxxxx:/virtualclusters/
lbhl6kwwyoxxxxxxxxxxxxxx/endpoints/as3ysxxxxxxxx",
    "virtualClusterId": "lbhl6kwwyoxxxxxxxxxxxxxxx",
    "type": "JUPYTER_ENTERPRISE_GATEWAY",
    "state": "ACTIVE",
    "releaseLabel": "emr-6.9.0-latest",
   "executionRoleArn": "arn:aws:iam::1828xxxxxxxx:role/RoleName",
    "certificateAuthority": {
        "certificateArn": "arn:aws:acm:us-east-1:1828xxxxxxxx:certificate/zzzzzzz-
e59b-4ed0-aaaa-bbbbbbbbbbbbbb,,
        "certificateData": "certificate-data"
    },
    "configurationOverrides": {
        "applicationConfiguration": [
                "classification": "spark-defaults",
                "properties": {
                    "spark.driver.memory": "8G"
                }
            }
        ],
        "monitoringConfiguration": {
            "persistentAppUI": "ENABLED",
            "cloudWatchMonitoringConfiguration": {
                "logGroupName": "log-group-name",
                "logStreamNamePrefix": "log-stream-name-prefix"
            },
            "s3MonitoringConfiguration": {
                "logUri": "s3-bucket-name"
            }
        }
    },
   "serverUrl": "https://internal-k8s-namespace-ingressa-aaaaaaaaa-
zzzzzzzzz.us-east-1.elb.amazonaws.com:18888 (https://internal-k8s-nspluto-
ingressa-51e860abbd-1620715833.us-east-1.elb.amazonaws.com:18888/)",
    "createdAt": "2022-09-19T12:37:49+00:00",
    "securityGroup": "sg-aaaaaaaaaaaaa",
    "subnetIds": [
        "subnet-1111111111",
        "subnet-222222222",
        "subnet-33333333333"
    "stateDetails": "Endpoint created successfully. It took 3 Minutes 15 Seconds",
    "tags": {}
```

}

Listagem de todos os endpoints interativos associados a um cluster virtual

Use o list-managed-endpoints AWS CLI comando para obter uma lista de todos os endpoints interativos associados a um cluster virtual especificado. Substitua virtual-cluster-id pelo ID do seu cluster virtual.

```
aws emr-containers list-managed-endpoints --virtual-cluster-id virtual-cluster-id
```

A saída do comando list-managed-endpoint é mostrada abaixo:

```
{
    "endpoints": [{
        "id": "as3ys2xxxxxxxx",
        "name": "endpoint-name",
        "arn": "arn:aws:emr-containers:us-east-1:1828xxxxxxxx:/virtualclusters/
lbhl6kwwyoxxxxxxxxxxxxxx/endpoints/as3ysxxxxxxxx",
        "virtualClusterId": "lbhl6kwwyoxxxxxxxxxxxxxxx",
        "type": "JUPYTER_ENTERPRISE_GATEWAY",
        "state": "ACTIVE",
        "releaseLabel": "emr-6.9.0-latest",
        "executionRoleArn": "arn:aws:iam::1828xxxxxxxx:role/RoleName",
        "certificateAuthority": {
            "certificateArn": "arn:aws:acm:us-east-1:1828xxxxxxxx:certificate/zzzzzzz-
e59b-4ed0-aaaa-bbbbbbbbbbbbb",
            "certificateData": "certificate-data"
        },
        "configurationOverrides": {
            "applicationConfiguration": [{
                "classification": "spark-defaults",
                "properties": {
                    "spark.driver.memory": "8G"
                }
            }],
            "monitoringConfiguration": {
                "persistentAppUI": "ENABLED",
                "cloudWatchMonitoringConfiguration": {
                    "logGroupName": "log-group-name",
                    "logStreamNamePrefix": "log-stream-name-prefix"
                "s3MonitoringConfiguration": {
```

```
"logUri": "s3-bucket-name"
                }
            }
        },
        "serverUrl": "https://internal-k8s-namespace-ingressa-aaaaaaaaa-
zzzzzzzzz.us-east-1.elb.amazonaws.com:18888 (https://internal-k8s-nspluto-
ingressa-51e860abbd-1620715833.us-east-1.elb.amazonaws.com:18888/)",
        "createdAt": "2022-09-19T12:37:49+00:00",
        "securityGroup": "sg-aaaaaaaaaaaaa",
        "subnetIds": [
            "subnet-11111111111",
            "subnet-222222222",
            "subnet-33333333333"
        ],
        "stateDetails": "Endpoint created successfully. It took 3 Minutes 15 Seconds",
        "tags": {}
    }]
}
```

Exclusão de um endpoint interativo

Para excluir um endpoint interativo associado a um cluster virtual Amazon EMR no EKS, use delete-managed-endpoint AWS CLI o comando. Quando você exclui um endpoint interativo, o Amazon EMR no EKS remove os grupos de segurança padrão criados para esse endpoint.

Especifique valores para os seguintes parâmetros do comando:

- --id: o identificador do endpoint interativo que você deseja excluir.
- -- virtual-cluster-id O identificador do cluster virtual associado ao endpoint interativo que você deseja excluir. Esse é o mesmo ID do cluster virtual que foi especificado quando o endpoint interativo foi criado.

```
aws emr-containers delete-managed-endpoint --id managed-endpoint-id --virtual-cluster-id virtual-cluster-id
```

O comando retorna uma saída semelhante à seguinte para confirmar que você excluiu o endpoint interativo:

```
{
"id":"8gai4l4exxxxx",
```

}

"virtualClusterId":"0b0qvauoy3ch1nqodxxxxxxxx"

Upload de dados para o Amazon S3 Express One Zone com o Amazon EMR no EKS

Com as versões 7.2.0 e superiores do Amazon EMR, você pode usar o Amazon EMR no EKS com a classe de armazenamento <u>Amazon S3 Express One Zone</u> para melhorar a performance ao executar trabalhos e workloads. O S3 Express One Zone é uma classe de armazenamento de zona única e alta performance do Amazon S3 que oferece acesso consistente a dados de milissegundos de um dígito para a maioria das aplicações sensíveis à latência. Na hora da execução, o S3 Express One Zone oferece o armazenamento de objetos na nuvem com a menor latência e a maior performance do Amazon S3.

Pré-requisitos

Antes de usar o S3 Express One Zone com o Amazon EMR no EKS, você deve ter os seguintes prérequisitos:

- Configuração do Amazon EMR no EKS concluída.
- Depois de configurar o Amazon EMR no EKS, crie um cluster virtual.

Conceitos básicos da classe S3 Express One Zone

Siga estas etapas para começar a usar o S3 Express One Zone

1. Adicione a permissão CreateSession ao perfil de execução do trabalho. Quando o S3 Express One Zone executa inicialmente uma ação como GET, LIST, ou PUT em um objeto do S3, a classe de armazenamento chama CreateSession em seu nome. A seguir está um exemplo de como conceder a permissão CreateSession.

JSON

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
```

Pré-requisitos 283

```
"Resource": [
    "arn:aws:s3express:*:*:bucket/DOC-EXAMPLE-BUCKET"
],
    "Action": [
        "s3express:CreateSession"
],
        "Sid": "AllowS3EXPRESSCreatesession"
}
]
```

2. Você deve usar o conector Apache Hadoop S3A para acessar os buckets do S3 Express, então altere seu Amazon S3 para usar o esquema de uso do conector. URIs s3a Se não usarem o esquema, você poderá alterar a implementação do sistema de arquivos usada nos esquemas s3 e s3n.

Para alterar o esquema do s3, especifique as seguintes configurações de cluster:

Para alterar o esquema do s3n, especifique as seguintes configurações de cluster:

3. Na configuração spark-submit, use o provedor de credenciais de identidade da Web.

Começar 284

"spark.hadoop.fs.s3a.aws.credentials.provider=com.amazonaws.auth.WebIdentityTokenCredentials.provider=com.amazonaws.auth.WebIdentityTokenCredentials.provider=com.amazonaws.auth.WebIdentityTokenCredentials.provider=com.amazonaws.auth.WebIdentityTokenCredentials.provider=com.amazonaws.auth.WebIdentityTokenCredentials.provider=com.amazonaws.auth.WebIdentityTokenCredentials.provider=com.amazonaws.auth.WebIdentityTokenCredentials.provider=com.amazonaws.auth.WebIdentityTokenCredentials.provider=com.amazonaws.auth.WebIdentityTokenCredentials.provider=com.amazonaws.auth.WebIdentityTokenCredentials.provider=com.amazonaws.auth.WebIdentityTokenCredentials.provider=com.amazonaws.auth.WebIdentityTokenCredentials.provider=com.amazonaws.auth.WebIdentityTokenCredentials.provider=com.amazonaws.auth.WebIdentityTokenCredentials.provider=com.amazonaws.auth.WebIdentityTokenCredentials.provider=com.amazonaws.auth.WebIdentityTokenCredentials.provider=com.amazonaws.auth.WebIdentityTokenCredentials.provider=com.amazonaws.auth.WebIdentials.provider=com.amazonaws.auth.We

Começar 285

Monitorar trabalhos

Você pode usar o Amazon CloudWatch Events para rastrear trabalhos executados em um cluster virtual Amazon EMR no EKS. É possível usar os eventos para rastrear a atividade e a integridade de trabalhos executados em um cluster virtual. Os tópicos a seguir mostram maneiras de configurar o monitoramento de forma eficaz para manter a integridade de recursos.

Tópicos

- Monitore trabalhos com a Amazon CloudWatch Events
- Automatize o Amazon EMR no EKS com eventos CloudWatch
- Exemplo: configuração de uma regra que invoque o Lambda
- Monitore o driver pod do trabalho com uma política de repetição usando o Amazon Events
 CloudWatch

Monitore trabalhos com a Amazon CloudWatch Events

O Amazon EMR no EKS emite eventos quando o estado de uma execução de trabalho é alterado. Cada evento fornece informações, como a data e o horário em que o evento ocorreu, em conjunto com mais detalhes sobre o evento, como o ID do cluster virtual e o ID de uma execução de trabalho que foi afetada.

É possível usar os eventos para rastrear a atividade e a integridade de trabalhos executados em um cluster virtual. Você também pode usar o Amazon CloudWatch Events para definir uma ação a ser tomada quando a execução de um trabalho gera um evento que corresponde a um padrão especificado por você. Os eventos são úteis para monitorar uma ocorrência específica durante o ciclo de vida de uma execução de trabalho. Por exemplo, é possível monitorar quando uma execução de trabalho sofre alteração de estado de submitted para running. Para obter mais informações sobre CloudWatch eventos, consulte o Guia EventBridge do usuário da Amazon.

A tabela apresentada a seguir lista os eventos do Amazon EMR no EKS em conjunto com o estado ou a alteração de estado que o evento indica, a severidade do evento e as mensagens do evento. Cada evento é representado como um objeto JSON que é enviado automaticamente a um stream de evento. O objeto JSON inclui mais detalhes sobre o evento. O objeto JSON é particularmente importante quando você configura regras para processamento de CloudWatch eventos usando Eventos, pois as regras buscam corresponder aos padrões no objeto JSON. Para obter mais

informações, consulte <u>Padrões de EventBridge eventos da Amazon</u> e Amazon EMR on EKS Events no Guia EventBridge do usuário da Amazon.

Eventos de alteração de estado em execuções de trabalhos

Estado	Gravidade	Mensagem
SUBMITTED	INFO	Job Run JobRunId (JobRunName) foi enviado com sucesso ao cluster virtual VirtualClusterId no Time UTC.
RUNNING (Em execução)	INFO	Job Run JobRunId (JobRunName) no cluster virtual VirtualClusterId começou a ser executado emTime.
CONCLUÍDO	INFO	Job Run jobRunId (JobRunName) no cluster virtual VirtualClusterId concluído emTime. O Job Run começou a ser executado às Time e levou Num minutos para ser concluído.
CANCELADO	WARN	A solicitação de cancelamento foi bem-suced ida para Job Run <i>JobRunId</i> (<i>JobRunName</i>) no cluster virtual <i>VirtualClusterId</i> em <i>Time</i> e o Job Run agora está cancelado.
COM FALHA	ERRO	Job Run JobRunId (JobRunName) no cluster virtual VirtualClusterId falhou emTime.

Automatize o Amazon EMR no EKS com eventos CloudWatch

Você pode usar o Amazon CloudWatch Events para automatizar seus AWS serviços para responder a eventos do sistema, como problemas de disponibilidade de aplicativos ou alterações de recursos. Os eventos dos AWS serviços são entregues aos CloudWatch Eventos quase em tempo real. É possível criar regras simples para indicar quais eventos são de seu interesse, e quais ações automatizadas devem ser tomadas quando um evento corresponder a uma regra. Ações que podem ser automaticamente acionadas incluem:

- Invocando uma função AWS Lambda
- Invocando o EC2 comando Amazon Run
- Transmitir o evento Amazon Kinesis Data Streams
- Ativando uma máquina de AWS Step Functions estado
- Notificando um tópico Amazon Simple Notification Service (SNS) ou uma fila Amazon Simple Queue Service (SQS)

Alguns exemplos do uso de CloudWatch eventos com o Amazon EMR no EKS incluem o seguinte:

- Ativação de uma função do Lambda quando uma execução de trabalho tiver êxito.
- Notificação de um tópico do Amazon SNS quando uma execução de trabalho falhar.

CloudWatch Eventos para "" detail-type: "" EMR Job Run State Change "são gerados pelo Amazon EMR no EKS paraSUBMITTED, RUNNINGCANCELLED, FAILED e mudanças de COMPLETED estado.

Exemplo: configuração de uma regra que invoque o Lambda

Use as etapas a seguir para configurar uma regra de CloudWatch eventos que invoca o Lambda quando há um evento "EMR Job Run State Change".

```
aws events put-rule \
--name cwe-test \
--event-pattern '{"detail-type": ["EMR Job Run State Change"]}'
```

Adicione a função Lambda que você possui como um novo destino e dê permissão a CloudWatch Events para invocar a função Lambda da seguinte forma. Substitua 123456789012 pelo ID da sua conta.

```
aws events put-targets \
--rule cwe-test \
--targets Id=1,Arn=arn:aws:lambda:us-east-1:123456789012:function:MyFunction
```

```
aws lambda add-permission \
--function-name MyFunction \
--statement-id MyId \
```

```
--action 'lambda:InvokeFunction' \
--principal events.amazonaws.com
```



Note

Não é possível escrever um programa que dependa da ordem ou da existência de eventos de notificação, pois eles podem estar fora de sequência ou ausentes. Os eventos são emitidos com base no melhor esforço.

Monitore o driver pod do trabalho com uma política de repetição usando o Amazon Events CloudWatch

Usando CloudWatch eventos, você pode monitorar os grupos de drivers que foram criados em trabalhos que têm políticas de repetição. Para obter mais informações, consulte Monitoramento de um trabalho com uma política de repetição neste guia.

Gerenciamento de clusters virtuais

Um cluster virtual corresponde a um namespace do Kubernetes no qual o Amazon EMR está registrado. Você pode criar, descrever, listar e excluir clusters virtuais. Eles não consomem quaisquer recursos adicionais em seu sistema. Um único cluster virtual mapeia para um único namespace do Kubernetes. Dado esse relacionamento, você pode modelar clusters virtuais da mesma forma que modela namespaces Kubernetes para atender aos seus requisitos. Confira os possíveis casos de uso na documentação de visão geral dos conceitos do Kubernetes.

Para registrar o Amazon EMR com um namespace do Kubernetes em um cluster do Amazon EKS, você precisa do nome do cluster do EKS e do namespace que foi configurado para executar sua workload. Esses clusters registrados no Amazon EMR são chamados de clusters virtuais porque não gerenciam computação ou armazenamento físicos, mas direcionam para um namespace do Kubernetes no qual sua workload está programada.



Note

Antes de criar um cluster virtual, você deve concluir as etapas de 1 a 8 em Configuração do Amazon EMR no EKS.

Tópicos

- Criação de um cluster virtual
- Listagem de clusters virtuais
- Descrição de um cluster virtual
- Exclusão de um cluster virtual
- Estados de um cluster virtual

Criação de um cluster virtual

Execute o comando apresentado a seguir para criar um cluster virtual ao registrar o Amazon EMR com um namespace em um cluster do EKS. virtual cluster nameSubstitua por um nome que você forneça para seu cluster virtual. eks_cluster_nameSubstitua pelo nome do cluster EKS. Substitua o pelo namespace *namespace_name* com o qual você deseja registrar o Amazon EMR.

Criação de um cluster virtual 290

Como alternativa, você pode criar um arquivo JSON que inclua os parâmetros obrigatórios para o cluster virtual, como demonstra o exemplo a seguir.

Em seguida, execute o comando create-virtual-cluster apresentado a seguir com o caminho para o arquivo JSON.

```
aws emr-containers create-virtual-cluster \
--cli-input-json file://./create-virtual-cluster-request.json
```

Note

Para validar a criação com êxito de um cluster virtual, visualize o status dos clusters virtuais ao executar o comando list-virtual-clusters ou ao acessar a página Clusters virtuais no console do Amazon EMR.

Criação de um cluster virtual 291

Listagem de clusters virtuais

Execute o comando apresentado a seguir para visualizar o status dos clusters virtuais.

aws emr-containers list-virtual-clusters

Descrição de um cluster virtual

Execute o comando apresentado a seguir para obter mais detalhes sobre um cluster virtual, como o namespace, o status e a data de registro. 123456Substitua pelo ID do seu cluster virtual.

aws emr-containers describe-virtual-cluster --id 123456

Exclusão de um cluster virtual

Execute o comando apresentado a seguir para excluir um cluster virtual. 123456Substitua pelo ID do seu cluster virtual.

aws emr-containers delete-virtual-cluster --id 123456

Estados de um cluster virtual

A tabela a seguir descreve os quatro estados possíveis de um cluster virtual.

State	Descrição
RUNNING	O cluster virtual está no estado RUNNING.
TERMINATING	O encerramento solicitado para o cluster virtual está em andamento.
TERMINATED	O encerramento solicitado foi concluído.
ARRESTED	O encerramento solicitado falhou devido a permissões insuficientes.

Listagem de clusters virtuais 292

Tutoriais para o Amazon EMR no EKS

Esta seção descreve casos de uso comuns para quando você trabalha com aplicações do Amazon EMR no EKS. Cada aplicação é especializada e pode executar etapas exclusivas de configuração. Esses tópicos fornecem instruções para o uso de cada aplicação.

Tópicos

- Uso do Delta Lake com o Amazon EMR no EKS
- Uso do Apache Iceberg com o Amazon EMR no EKS
- Usando PyFlink
- Usando AWS Glue com Flink
- Uso do Apache Hudi com o Apache Flink
- Uso do acelerador RAPIDS para Apache Spark com o Amazon EMR no EKS
- Uso da integração do Amazon Redshift para Apache Spark no Amazon EMR no EKS
- Uso do Volcano como um programador personalizado para Apache Spark no Amazon EMR no EKS
- Usando YuniKorn como um programador personalizado para o Apache Spark no Amazon EMR no EKS

Uso do Delta Lake com o Amazon EMR no EKS

O Delta Lake é uma estrutura de armazenamento de código aberto para a criação de uma arquitetura do Lakehouse. O exemplo a seguir mostra como configurar para uso.

Para usar o Delta Lake com aplicações do Amazon EMR no EKS

 Ao iniciar uma execução de trabalho para enviar um trabalho do Spark na configuração da aplicação, inclua os arquivos JAR do Delta Lake:

Uso do Delta Lake 293



Note

As versões 7.0.0 e superiores do Amazon EMR usam o Delta Lake 3.0, que renomeia delta-core.jar para delta-spark.jar. Se você usa o Amazon EMR nas versões 7.0.0 ou superiores, certifique-se de usar o nome de arquivo correto, como no exemplo a seguir:

```
--jars local:///usr/share/aws/delta/lib/delta-spark.jar
```

2. Inclua a configuração adicional do Delta Lake e use o AWS Glue Data Catalog como seu metastore.

```
--configuration-overrides '{
        "applicationConfiguration": [
          "classification" : "spark-defaults",
          "properties" : {
            "spark.sql.extensions" : "io.delta.sql.DeltaSparkSessionExtension",
 "spark.sql.catalog.spark_catalog":"org.apache.spark.sql.delta.catalog.DeltaCatalog",
"spark.hadoop.hive.metastore.client.factory.class":"com.amazonaws.glue.catalog.metastore.AW
       }]}'
```

Uso do Apache Iceberg com o Amazon EMR no EKS

O JAR de runtime para o Iceberg contém as classes necessárias do Iceberg para suporte ao runtime do Spark. O procedimento a seguir mostra como iniciar uma execução de trabalho usando o runtime do Spark do Iceberg.

Usar o Apache Iceberg com as aplicações do Amazon EMR no EKS

1. Ao iniciar uma execução de trabalho para enviar um trabalho do Spark na configuração da aplicação, inclua o arquivo JAR de runtime do Iceberg para o Spark:

```
--job-driver '{"sparkSubmitJobDriver" : {"sparkSubmitParameters" : "--jars
 local:///usr/share/aws/iceberg/lib/iceberg-spark3-runtime.jar"}}'
```

Uso do Iceberg 294

2. Inclua a configuração adicional do Iceberg:

```
--configuration-overrides '{
    "applicationConfiguration": [
    "classification": "spark-defaults",
    "properties": {
        "spark.sql.catalog.dev.warehouse": "s3://amzn-s3-demo-bucket/EXAMPLE-
PREFIX/ ",
        "spark.sql.extensions ":"
    org.apache.iceberg.spark.extensions.IcebergSparkSessionExtensions ",
        "spark.sql.catalog.dev": "org.apache.iceberg.spark.SparkCatalog",
        "spark.sql.catalog.dev.catalog-impl":
    "org.apache.iceberg.aws.glue.GlueCatalog",
        "spark.sql.catalog.dev.io-impl": "org.apache.iceberg.aws.s3.S3FileIO"
    }
}
```

Para saber mais sobre as versões de liberação do Apache Iceberg para o EMR, consulte <u>Iceberg</u> release history.

Configurações de sessão do Spark para integração de catálogos

Configurações de sessão do Spark para integração do catálogo Iceberg AWS Glue

Este exemplo mostra como integrar o Iceberg com: Crawler do AWS Glue

```
spark-sql \
    --conf spark.sql.catalog.rms = org.apache.iceberg.spark.SparkCatalog \
    --conf spark.sql.catalog.rms.type = glue \
    --conf spark.sql.catalog.rms.glue.id = glue RMS catalog ID \
    --conf spark.sql.catalog.rms.glue.account-id = AWS account ID \
    --conf spark.sql.extensions=
    org.apache.iceberg.spark.extensions.IcebergSparkSessionExtensions
```

Veja a seguir um exemplo de consulta:

```
SELECT * FROM rms.rmsdb.table1
```

Configurações de sessão do Spark para integração do catálogo Iceberg REST AWS Glue

Este exemplo mostra como integrar o Iceberg REST com: Crawler do AWS Glue

```
spark-sql \
    --conf spark.sql.catalog.rms = org.apache.iceberg.spark.SparkCatalog \
    --conf spark.sql.catalog.rms.type = rest \
    --conf spark.sql.catalog.rms.warehouse = glue RMS catalog ID \
    --conf spark.sql.catalog.rms.uri = glue endpoint URI/iceberg \
    --conf spark.sql.catalog.rms.rest.sigv4-enabled = true \
    --conf spark.sql.catalog.rms.rest.signing-name = glue \
    --conf spark.sql.catalog.rms.rest.signing-name = glue \
--conf spark.sql.extensions=
org.apache.iceberg.spark.extensions.IcebergSparkSessionExtensions
```

Veja a seguir um exemplo de consulta:

```
SELECT * FROM rms.rmsdb.table1
```

Essa configuração funciona somente para o Redshift Managed Storage. O FGAC para Amazon S3 não é compatível.

Usando PyFlink

Suporte para o Amazon EMR no EKS nas versões 6.15.0 e superiores. PyFlink Se você já tem um PyFlink script, você pode fazer o seguinte:

- · Crie uma imagem personalizada com seu PyFlink script incluído.
- Faça upload do script em um local do Amazon S3

Se você ainda não tem um script, pode usar o exemplo a seguir para iniciar um PyFlink trabalho. Este exemplo recupera o script do S3. Caso esteja usando uma imagem personalizada com o script já incluído na imagem, atualize o caminho do script até o local em que você armazenou o script. Se o script estiver em um local do S3, o Amazon EMR no EKS recuperará o script e o colocará no diretório /opt/flink/usrlib/ no contêiner do Flink.

```
apiVersion: flink.apache.org/v1beta1
```

Usando PyFlink 296

```
kind: FlinkDeployment
metadata:
  name: python-example
spec:
  flinkVersion: v1_17
  flinkConfiguration:
    taskmanager.numberOfTaskSlots: "1"
  executionRoleArn: job-execution-role
  emrReleaseLabel: "emr-6.15.0-flink-latest"
  jobManager:
    highAvailabilityEnabled: false
    replicas: 1
    resource:
      memory: "2048m"
      cpu: 1
  taskManager:
    resource:
      memory: "2048m"
      cpu: 1
  job:
    jarURI: s3://S3 bucket with your script/pyflink-script.py
    entryClass: "org.apache.flink.client.python.PythonDriver"
    args: ["-py", "/opt/flink/usrlib/pyflink-script.py"]
    parallelism: 1
    upgradeMode: stateless
```

Usando AWS Glue com Flink

O Amazon EMR no EKS com Apache Flink, versões 6.15.0 e superiores, oferece suporte ao uso do AWS Glue Data Catalog como um armazenamento de metadados para fluxos de trabalho SQL em lote e streaming.

Primeiro, você deve criar um banco de dados AWS Glue chamado default que sirva como seu catálogo SQL do Flink. Esse catálogo do Flink armazena metadados como bancos de dados, tabelas, partições, visualizações, funções e outras informações necessárias para acessar dados em outros sistemas externos.

```
aws glue create-database \
   --database-input "{\"Name\":\"default\"}"
```

Usando AWS Glue com Flink 297

Para ativar o suporte ao AWS Glue, use uma FlinkDeployment especificação. Esse exemplo de especificação usa um script Python para emitir rapidamente algumas instruções SQL do Flink para interagir com o AWS catálogo Glue.

```
apiVersion: flink.apache.org/v1beta1
kind: FlinkDeployment
metadata:
  name: python-example
spec:
  flinkVersion: v1_17
  flinkConfiguration:
    taskmanager.numberOfTaskSlots: "1"
    aws.glue.enabled: "true"
  executionRoleArn: job-execution-role-arn;
  emrReleaseLabel: "emr-6.15.0-flink-latest"
  jobManager:
    highAvailabilityEnabled: false
    replicas: 1
    resource:
      memory: "2048m"
      cpu: 1
  taskManager:
    resource:
      memory: "2048m"
      cpu: 1
  job:
    jarURI: s3://<S3_bucket_with_your_script/pyflink-glue-script.py
    entryClass: "org.apache.flink.client.python.PythonDriver"
    args: ["-py", "/opt/flink/usrlib/pyflink-glue-script.py"]
    parallelism: 1
    upgradeMode: stateless
```

Veja a seguir um exemplo da aparência do seu PyFlink script.

```
import logging
import sys
from pyflink.datastream import StreamExecutionEnvironment
from pyflink.table import StreamTableEnvironment

def glue_demo():
    env = StreamExecutionEnvironment.get_execution_environment()
    t_env = StreamTableEnvironment.create(stream_execution_environment=env)
    t_env.execute_sql("""
```

Usando AWS Glue com Flink 298

```
CREATE CATALOG glue_catalog WITH (
          'type' = 'hive',
          'default-database' = 'default',
          'hive-conf-dir' = '/glue/confs/hive/conf',
          'hadoop-conf-dir' = '/glue/confs/hadoop/conf'
    t_env.execute_sql("""
          USE CATALOG glue_catalog;
    t_env.execute_sql("""
          DROP DATABASE IF EXISTS eks_flink_db CASCADE;
    t_env.execute_sql("""
          CREATE DATABASE IF NOT EXISTS eks_flink_db WITH ('hive.database.location-
uri'= 's3a://S3-bucket-to-store-metadata/flink/flink-glue-for-hive/warehouse/');
    t_env.execute_sql("""
          USE eks_flink_db;
                  """)
    t_env.execute_sql("""
          CREATE TABLE IF NOT EXISTS eksglueorders (
            order_number BIGINT,
                         DECIMAL(32,2),
            price
            buyer
                         RO first_name STRING, last_name STRING,
            order_time TIMESTAMP(3)
          ) WITH (
            'connector' = 'datagen'
          );
    t_env.execute_sql("""
          CREATE TABLE IF NOT EXISTS eksdestglueorders (
            order_number BIGINT,
            price
                         DECIMAL(32,2),
                         ROW first_name STRING, last_name STRING,
            buyer
            order_time TIMESTAMP(3)
          ) WITH (
            'connector' = 'filesystem',
            'path' = 's3://S3-bucket-to-store-metadata/flink/flink-glue-for-hive/
warehouse/eksdestglueorders',
            'format' = 'json'
          );
                  """)
    t_env.execute_sql("""
```

Usando AWS Glue com Flink 299

```
CREATE TABLE IF NOT EXISTS print_table (
            order_number BIGINT,
                         DECIMAL(32,2),
            price
                         ROW first_name STRING, last_name STRING,
            buyer
                         TIMESTAMP(3)
            order_time
          ) WITH (
            'connector' = 'print'
          );
                """)
    t_env.execute_sql("""
          EXECUTE STATEMENT SET
          BEGIN
          INSERT INTO eksdestglueorders SELECT * FROM eksglueorders LIMIT 10;
          INSERT INTO print_table SELECT * FROM eksdestglueorders;
          END;
            """)
if __name__ == '__main__':
    logging.basicConfig(stream=sys.stdout, level=logging.INFO, format="%(message)s")
    glue_demo()
```

Uso do Apache Hudi com o Apache Flink

O Apache Hudi é uma estrutura de gerenciamento de dados de código aberto com operações em nível de registro, como inserir, atualizar, upsert e excluir, que você pode usar para simplificar o gerenciamento de dados e o desenvolvimento de pipelines de dados. Quando combinado com o gerenciamento eficiente de dados no Amazon S3, o Hudi permite que você ingira e atualize dados em tempo real. O Hudi mantém os metadados de todas as operações que você executa no conjunto de dados, para que todas as ações permaneçam atômicas e consistentes.

O Apache Hudi está disponível no Amazon EMR no EKS com o Apache Flink, com as versões 7.2.0 e superiores do Amazon EMR. Confira as etapas a seguir para saber como começar e enviar trabalhos do Apache Hudi.

Envio de um trabalho do Apache Hudi

Confira as etapas a seguir para saber como enviar um trabalho do Apache Hudi.

Crie um banco de dados AWS Glue chamadodefault.

Uso do Apache Hudi 300

```
aws glue create-database --database-input "{\"Name\":\"default\"}"
```

- 2. Siga as instruções em <u>Flink Kubernetes Operator SQL Example</u> para criar o arquivo flink-sql-runner.jar.
- 3. Crie um script SQL do Hudi como o mostrado a seguir.

```
CREATE CATALOG hudi_glue_catalog WITH (
'type' = 'hudi',
'mode' = 'hms',
'table.external' = 'true',
'default-database' = 'default',
'hive.conf.dir' = '/glue/confs/hive/conf/',
'catalog.path' = 's3://<hudi-example-bucket>/FLINK_HUDI/warehouse/'
);
USE CATALOG hudi_glue_catalog;
CREATE DATABASE IF NOT EXISTS hudi_db;
use hudi_db;
CREATE TABLE IF NOT EXISTS hudi-flink-example-table(
    uuid VARCHAR(20),
    name VARCHAR(10),
    age INT,
    ts TIMESTAMP(3),
    `partition` VARCHAR(20)
)
PARTITIONED BY ('partition')
WITH (
  'connector' = 'hudi',
  'path' = 's3://<hudi-example-bucket>/hudi-flink-example-table',
  'hive_sync.enable' = 'true',
  'hive_sync.mode' = 'glue',
  'hive_sync.table' = 'hudi-flink-example-table',
  'hive_sync.db' = 'hudi_db',
  'compaction.delta_commits' = '1',
  'hive_sync.partition_fields' = 'partition',
  'hive_sync.partition_extractor_class' =
 'org.apache.hudi.hive.MultiPartKeysValueExtractor',
  'table.type' = 'COPY_ON_WRITE'
);
EXECUTE STATEMENT SET
```

```
INSERT INTO hudi-flink-example-table VALUES
    ('id1','Alex',23,TIMESTAMP '1970-01-01 00:00:01','par1'),
    ('id2','Stephen',33,TIMESTAMP '1970-01-01 00:00:02','par1'),
    ('id3','Julian',53,TIMESTAMP '1970-01-01 00:00:03','par2'),
    ('id4','Fabian',31,TIMESTAMP '1970-01-01 00:00:04','par2'),
    ('id5','Sophia',18,TIMESTAMP '1970-01-01 00:00:05','par3'),
    ('id6','Emma',20,TIMESTAMP '1970-01-01 00:00:06','par3'),
    ('id7','Bob',44,TIMESTAMP '1970-01-01 00:00:07','par4'),
    ('id8','Han',56,TIMESTAMP '1970-01-01 00:00:08','par4');
END;
```

- 4. Faça upload do script SQL do Hudi e do arquivo flink-sql-runner.jar em um local do S3.
- 5. No seu arquivo FlinkDeployments do YAML, defina hudi.enabled como true.

```
spec:
  flinkConfiguration:
   hudi.enabled: "true"
```

 Crie um arquivo do YAML para executar sua configuração. Neste exemplo, o nome do arquivo é hudi-write.yaml.

```
apiVersion: flink.apache.org/v1beta1
kind: FlinkDeployment
metadata:
  name: hudi-write-example
spec:
  flinkVersion: v1_18
 flinkConfiguration:
    taskmanager.numberOfTaskSlots: "2"
    hudi.enabled: "true"
  executionRoleArn: "<JobExecutionRole>"
  emrReleaseLabel: "emr-7.10.0-flink-latest"
  jobManager:
    highAvailabilityEnabled: false
    replicas: 1
    resource:
      memory: "2048m"
      cpu: 1
  taskManager:
    resource:
```

```
memory: "2048m"
    cpu: 1
job:
  jarURI: local:///opt/flink/usrlib/flink-sql-runner.jar
  args: ["/opt/flink/scripts/hudi-write.sql"]
  parallelism: 1
  upgradeMode: stateless
podTemplate:
  spec:
    initContainers:
      - name: flink-sql-script-download
        args:
          - s3
          - ср
          - s3://<s3_location>/hudi-write.sql
          - /flink-scripts
        image: amazon/aws-cli:latest
        imagePullPolicy: Always
        resources: {}
        terminationMessagePath: /dev/termination-log
        terminationMessagePolicy: File
        volumeMounts:
          - mountPath: /flink-scripts
            name: flink-scripts
      - name: flink-sql-runner-download
        args:
          - s3
          - s3://<s3_location>/flink-sql-runner.jar
          - /flink-artifacts
        image: amazon/aws-cli:latest
        imagePullPolicy: Always
        resources: {}
        terminationMessagePath: /dev/termination-log
        terminationMessagePolicy: File
        volumeMounts:
          - mountPath: /flink-artifacts
            name: flink-artifact
    containers:
      - name: flink-main-container
        volumeMounts:
          - mountPath: /opt/flink/scripts
            name: flink-scripts
          - mountPath: /opt/flink/usrlib
```

```
name: flink-artifact
volumes:
   - emptyDir: {}
   name: flink-scripts
   - emptyDir: {}
   name: flink-artifact
```

7. Envie um trabalho do Hudi no Flink para o operador Kubernetes do Flink.

```
kubectl apply -f hudi-write.yaml
```

Uso do acelerador RAPIDS para Apache Spark com o Amazon EMR no EKS

Com o Amazon EMR no EKS, você pode executar trabalhos para o acelerador RAPIDS da Nvidia para Apache Spark. Este tutorial aborda como executar trabalhos do Spark usando RAPIDS em tipos de instância de unidade de processamento EC2 gráfico (GPU). O tutorial usa as seguintes versões:

- Amazon EMR no EKS na versão de liberação 6.9.0 e posteriores
- · Apache Spark 3.x

Você pode acelerar o Spark com os tipos de instância de EC2 GPU da Amazon usando o plug-in Nvidia RAPIDS Accelerator for Apache Spark. Ao usar essas tecnologias em conjunto, você acelera os pipelines de ciência de dados sem a necessidade de realizar alterações no código. Isso reduz o tempo de execução necessário para o processamento de dados e para o treinamento de modelos. Ao fazer mais em menos tempo, você gasta menos com o custo da infraestrutura.

Antes de começar, certifique-se de ter os recursos apresentados a seguir.

- Cluster virtual do Amazon EMR no EKS
- Cluster do Amazon EKS com um grupo de nós habilitado para GPU

Um cluster virtual do Amazon EKS corresponde a um manipulador registrado para o namespace do Kubernetes em um cluster do Amazon EKS e é gerenciado pelo Amazon EMR no EKS. O manipulador permite que o Amazon EMR use o namespace do Kubernetes como destino para a execução de trabalhos. Para obter mais informações sobre como configurar um cluster virtual, consulte Configuração do Amazon EMR no EKS neste guia.

Você deve configurar o cluster virtual do Amazon EKS com um grupo de nós que tenha instâncias de GPU. Você deve configurar os nós com um plug-in para dispositivos da Nvidia. Consulte <u>Grupos de nós gerenciados</u> para saber mais.

Para configurar o cluster do Amazon EKS para adicionar grupos de nós habilitados para GPU, execute o seguinte procedimento:

Adicionar grupos de nós habilitados para GPU

 Crie um grupo de nós habilitado para GPU com o comando <u>create-nodegroup</u> apresentado a seguir. Certifique-se de realizar a substituição com os parâmetros corretos para o cluster do Amazon EKS. Use um tipo de instância compatível com RAPIDS para Spark, como P4, P3, G5 ou G4dn.

```
aws eks create-nodegroup \
    --cluster-name EKS_CLUSTER_NAME \
    --nodegroup-name NODEGROUP_NAME \
    --scaling-config minSize=0,maxSize=5,desiredSize=2 CHOOSE_APPROPRIATELY \
    --ami-type AL2_x86_64_GPU \
    --node-role NODE_ROLE \
    --subnets SUBNETS_SPACE_DELIMITED \
    --remote-access ec2SshKey= SSH_KEY \
    --instance-types GPU_INSTANCE_TYPE \
    --disk-size DISK_SIZE \
    --region AWS_REGION
```

2. Instale o plug-in do dispositivo Nvidia em seu cluster para emitir o número de GPUs em cada nó do seu cluster e para executar contêineres habilitados para GPU em seu cluster. Execute o seguinte código para instalar o plug-in:

```
kubectl apply -f https://raw.githubusercontent.com/NVIDIA/k8s-device-plugin/v0.9.0/
nvidia-device-plugin.yml
```

3. Para validar quantos GPUs estão disponíveis em cada nó do seu cluster, execute o seguinte comando:

```
kubectl get nodes "-o=custom-
columns=NAME:.metadata.name,GPU:.status.allocatable.nvidia\.com/gpu"
```

Executar um trabalho do RAPIDS para Spark

 Envie um trabalho do RAPIDS para Spark ao cluster do Amazon EMR no EKS. O código a seguir apresenta um exemplo de comando para iniciar o trabalho. Na primeira vez em que você executar o trabalho, poderá demorar alguns minutos para fazer download da imagem e armazená-la em cache no nó.

```
aws emr-containers start-job-run \
--virtual-cluster-id VIRTUAL_CLUSTER_ID \
--execution-role-arn JOB_EXECUTION_ROLE \
--release-label emr-6.9.0-spark-rapids-latest \
--job-driver '{"sparkSubmitJobDriver": {"entryPoint": "local:///usr/lib/
spark/examples/jars/spark-examples.jar","entryPointArguments": ["10000"],
    "sparkSubmitParameters":"--class org.apache.spark.examples.SparkPi "}}' \
---configuration-overrides '{"applicationConfiguration": [{"classification":
    "spark-defaults","properties": {"spark.executor.instances":
    "2","spark.executor.memory": "2G"}}],"monitoringConfiguration":
    {"cloudWatchMonitoringConfiguration": {"logGroupName": "LOG_GROUP
_NAME"},"s3MonitoringConfiguration": {"logUri": "LOG_GROUP_STREAM"}}}'
```

2. Para validar se o acelerador RAPIDS para Spark está habilitado, verifique os logs do driver do Spark. Esses registros são armazenados no CloudWatch local do S3 que você especifica ao executar o start-job-run comando. O seguinte exemplo apresenta como geralmente é a aparência das linhas de log:

```
22/11/15 00:12:44 INFO RapidsPluginUtils: RAPIDS Accelerator build:
{version=22.08.0-amzn-0, user=release, url=, date=2022-11-03T03:32:45Z, revision=, cudf_version=22.08.0, branch=}

22/11/15 00:12:44 INFO RapidsPluginUtils: RAPIDS Accelerator JNI build:
{version=22.08.0, user=, url=https://github.com/NVIDIA/spark-rapids-jni.git, date=2022-08-18T04:14:34Z, revision=a1b23cd_sample, branch=HEAD}

22/11/15 00:12:44 INFO RapidsPluginUtils: cudf build: {version=22.08.0, user=, url=https://github.com/rapidsai/cudf.git, date=2022-08-18T04:14:34Z, revision=a1b23ce_sample, branch=HEAD}

22/11/15 00:12:44 WARN RapidsPluginUtils: RAPIDS Accelerator 22.08.0-amzn-0 using cudf 22.08.0.

22/11/15 00:12:44 WARN RapidsPluginUtils: spark.rapids.sql.multiThreadedRead.numThreads is set to 20.

22/11/15 00:12:44 WARN RapidsPluginUtils: RAPIDS Accelerator is enabled, to disable GPU support set `spark.rapids.sql.enabled` to false.
```

```
22/11/15 00:12:44 WARN RapidsPluginUtils: spark.rapids.sql.explain is set to `NOT_ON_GPU`. Set it to 'NONE' to suppress the diagnostics logging about the query placement on the GPU.
```

3. Para visualizar as operações que serão executadas em uma GPU, execute as etapas apresentadas a seguir para ativar o registro em log adicional. Observe a configuração "spark.rapids.sql.explain: ALL".

```
aws emr-containers start-job-run \
--virtual-cluster-id VIRTUAL_CLUSTER_ID \
--execution-role-arn JOB_EXECUTION_ROLE \
--release-label emr-6.9.0-spark-rapids-latest \
--job-driver '{"sparkSubmitJobDriver": {"entryPoint": "local:///usr/lib/
spark/examples/jars/spark-examples.jar","entryPointArguments": ["10000"],
    "sparkSubmitParameters":"--class org.apache.spark.examples.SparkPi "}}' \
---configuration-overrides '{"applicationConfiguration":
    [{"classification": "spark-defaults","properties":
    {"spark.rapids.sql.explain":"ALL","spark.executor.instances":
    "2","spark.executor.memory": "2G"}}],"monitoringConfiguration":
    {"cloudWatchMonitoringConfiguration": {"logGroupName":
    "LOG_GROUP_NAME"},"s3MonitoringConfiguration": {"logUri": "LOG_GROUP_STREAM"}}}'
```

O comando anterior é um exemplo de trabalho que usa a GPU. Sua saída seria semelhante ao exemplo abaixo. Consulte esta legenda para obter ajuda para compreender a saída:

- *: marca uma operação que funciona em uma GPU.
- !: marca uma operação que não pode ser executada em uma GPU.
- e: marca uma operação que funciona em uma GPU, mas não será executada porque está em um plano que não pode ser executado em uma GPU.

```
22/11/15 01:22:58 INFO GpuOverrides: Plan conversion to the GPU took 118.64 ms
22/11/15 01:22:58 INFO GpuOverrides: Plan conversion to the GPU took 4.20 ms
22/11/15 01:22:58 INFO GpuOverrides: GPU plan transition optimization took 8.37 ms
22/11/15 01:22:59 WARN GpuOverrides:
    *Exec <ProjectExec> will run on GPU
    *Expression <Alias> substring(cast(date#149 as string), 0, 7) AS month#310
will run on GPU
    *Expression <Substring> substring(cast(date#149 as string), 0, 7) will run on GPU
```

```
*Expression <Cast> cast(date#149 as string) will run on GPU
      *Exec <SortExec> will run on GPU
        *Expression <SortOrder> date#149 ASC NULLS FIRST will run on GPU
        *Exec <ShuffleExchangeExec> will run on GPU
          *Partitioning <RangePartitioning> will run on GPU
            *Expression <SortOrder> date#149 ASC NULLS FIRST will run on GPU
          *Exec <UnionExec> will run on GPU
            !Exec <ProjectExec> cannot run on GPU because not all expressions can
be replaced
             @Expression <AttributeReference> customerID#0 could run on GPU
             @Expression <Alias> Charge AS kind#126 could run on GPU
               @Expression <Literal> Charge could run on GPU
             @Expression <AttributeReference> value#129 could run on GPU
             @Expression <Alias> add_months(2022-11-15, cast(-(cast(_we0#142 as
bigint) + last_month#128L) as int)) AS date#149 could run on GPU
                ! <AddMonths> add_months(2022-11-15, cast(-
(cast(_we0#142 as bigint) + last_month#128L) as int)) cannot run
on GPU because GPU does not currently support the operator class
org.apache.spark.sql.catalyst.expressions.AddMonths
                 @Expression <Literal> 2022-11-15 could run on GPU
                 @Expression <Cast> cast(-(cast(_we0#142 as bigint) +
last_month#128L) as int) could run on GPU
                    @Expression <UnaryMinus> -(cast(_we0#142 as bigint) +
last_month#128L) could run on GPU
                     @Expression <Add> (cast(_we0#142 as bigint) +
last_month#128L) could run on GPU
                       @Expression <Cast> cast(_we0#142 as bigint) could run on
GPU
                         @Expression <AttributeReference> _we0#142 could run on
GPU
                        @Expression <AttributeReference> last_month#128L could run
on GPU
```

Uso da integração do Amazon Redshift para Apache Spark no Amazon EMR no EKS

Com as versões 6.9.0 e posteriores do Amazon EMR, cada imagem de versão inclui um conector entre o <u>Apache Spark</u> e o Amazon Redshift. Dessa forma, você pode usar o Spark no Amazon EMR no EKS para processar dados armazenados no Amazon Redshift. A integração é baseada

Uso do Spark no Redshift 308

no <u>conector de código aberto spark-redshift</u>. Para o Amazon EMR no EKS, a <u>integração do</u> Amazon Redshift para Apache Spark está inclusa como uma integração nativa.

Tópicos

- Inicialização de uma aplicação do Spark usando a integração do Amazon Redshift para Apache
 Spark
- Autenticação com a integração do Amazon Redshift para Apache Spark
- Leitura e gravação de e para o Amazon Redshift
- Considerações e limitações ao usar o conector do Spark

Inicialização de uma aplicação do Spark usando a integração do Amazon Redshift para Apache Spark

Para usar a integração, você deve transferir as dependências obrigatórias do Redshift para Spark com o trabalho do Spark. Você deve usar --jars para incluir as bibliotecas relacionadas ao conector do Redshift. Para visualizar outros locais de arquivo com suporte pela opção --jars, consulte a seção Advanced Dependency Management da documentação do Apache Spark.

```
• spark-redshift.jar
```

- spark-avro.jar
- RedshiftJDBC.jar
- minimal-json.jar

Para iniciar uma aplicação do Spark com a integração do Amazon Redshift para Apache Spark na versão 6.9.0 ou em versões posteriores do Amazon EMR no EKS, use o comando de exemplo a seguir. Observe que os caminhos listados com a opção --conf spark.jars são os caminhos padrão para os arquivos JAR.

```
aws emr-containers start-job-run \
--virtual-cluster-id cluster_id \
--execution-role-arn arn \
--release-label emr-6.9.0-latest\
--job-driver '{
    "sparkSubmitJobDriver": {
        "entryPoint": "s3://script_path",
```

```
"sparkSubmitParameters":
"--conf spark.kubernetes.file.upload.path=s3://upload_path
--conf spark.jars=
    /usr/share/aws/redshift/jdbc/RedshiftJDBC.jar,
    /usr/share/aws/redshift/spark-redshift/lib/spark-redshift.jar,
    /usr/share/aws/redshift/spark-redshift/lib/spark-avro.jar,
    /usr/share/aws/redshift/spark-redshift/lib/minimal-json.jar"
    }
}'
```

Autenticação com a integração do Amazon Redshift para Apache Spark

As seções a seguir mostram as opções de autenticação com o Amazon Redshift durante a integração com o Apache Spark. As seções mostram como recuperar as credenciais de login e detalhes sobre o uso do driver JDBC com a autenticação do IAM.

Use AWS Secrets Manager para recuperar credenciais e conectar-se ao Amazon Redshift

Você pode armazenar credenciais no Secrets Manager para realizar a autenticação com segurança no Amazon Redshift. É possível fazer com que seu trabalho do Spark chame a API GetSecretValue para buscar as credenciais:

```
from pyspark.sql import SQLContextimport boto3

sc = # existing SparkContext
sql_context = SQLContext(sc)

secretsmanager_client = boto3.client('secretsmanager',
    region_name=os.getenv('AWS_REGION'))
secret_manager_response = secretsmanager_client.get_secret_value(
    SecretId='string',
    VersionId='string',
    VersionStage='string'
)
username = # get username from secret_manager_response
password = # get password from secret_manager_response
url = "jdbc:redshift://redshifthost:5439/database?user=" + username + "&password=" + password

# Access to Redshift cluster using Spark
```

Uso da autenticação baseada no IAM com o perfil de execução de trabalho do Amazon EMR no EKS

A partir da versão 6.9.0 do Amazon EMR no EKS, a versão 2.1 ou as versões superiores do driver JDBC do Amazon Redshift são empacotadas no ambiente. Com a versão 2.1 e versões superiores do driver JDBC, é possível especificar o URL do JDBC e não incluir o nome de usuário e a senha brutos. Em vez disso, você pode especificar o esquema jdbc:redshift:iam://. Isso comanda o driver JDBC para usar seu perfil de execução de trabalho do Amazon EMR no EKS para buscar as credenciais automaticamente.

Consulte Configurar uma conexão JDBC ou ODBC para usar credenciais do IAM no Guia de gerenciamento do Amazon Redshift para obter mais informações.

O exemplo de URL a seguir usa um esquema jdbc:redshift:iam://.

```
jdbc:redshift:iam://examplecluster.abc123xyz789.us-west-2.redshift.amazonaws.com:5439/
dev
```

As permissões apresentadas a seguir são obrigatórias para o seu perfil de execução de trabalho quando ele atende às condições fornecidas.

Permissão	Condições para se tornar obrigatória para o perfil de execução de trabalho
<pre>redshift:GetCluste rCredentials</pre>	Obrigatória para que o driver JDBC busque as credenciais do Amazon Redshift.
redshift:DescribeC luster	Obrigatória se você especificar o cluster do Amazon Redshift e a Região da AWS no URL do JDBC em vez do endpoint.
<pre>redshift-serverles s:GetCredentials</pre>	Obrigatória para que o driver JDBC busque as credenciais do Amazon Redshift sem servidor.
redshift-serverles s:GetWorkgroup	Obrigatória se você estiver usando o Amazon Redshift sem servidor e especificar o URL em termos de nome e de região do grupo de trabalho.

Sua política de perfil de execução de trabalho deve ter as permissões apresentadas a seguir.

```
{
            "Effect": "Allow",
            "Action": [
                "redshift:GetClusterCredentials",
                "redshift:DescribeCluster",
                "redshift-serverless:GetCredentials",
                "redshift-serverless:GetWorkgroup"
            ],
            "Resource": [
 "arn:aws:redshift:AWS_REGION:ACCOUNT_ID:dbname:CLUSTER_NAME/DATABASE_NAME",
                "arn:aws:redshift:AWS_REGION:ACCOUNT_ID:dbuser:DATABASE_NAME/USER_NAME"
            ]
        }
```

Autenticação no Amazon Redshift com um driver JDBC

Definição de um nome de usuário e de uma senha no URL do JDBC

Para autenticar um trabalho do Spark em um cluster do Amazon Redshift, você pode especificar o nome e a senha do banco de dados do Amazon Redshift no URL do JDBC.



Note

Se você transferir as credenciais do banco de dados no URL, qualquer pessoa que tenha acesso ao URL também poderá acessar as credenciais. Este método geralmente não é recomendado porque não é uma opção segura.

Se a segurança não for uma preocupação para sua aplicação, você poderá usar o seguinte formato para definir o nome de usuário e a senha no URL do JDBC:

```
jdbc:redshift://redshifthost:5439/database?user=username&password=password
```

Leitura e gravação de e para o Amazon Redshift

Os exemplos de código a seguir são usados PySpark para ler e gravar dados de amostra de e para um banco de dados do Amazon Redshift com uma API de fonte de dados e com o SparkSQL.

Data source API

Use PySpark para ler e gravar dados de amostra de e para um banco de dados do Amazon Redshift com uma API de fonte de dados.

```
import boto3
from pyspark.sql import SQLContext
sc = # existing SparkContext
sql_context = SQLContext(sc)
url = "jdbc:redshift:iam://redshifthost:5439/database"
aws_iam_role_arn = "arn:aws:iam::accountID:role/roleName"
df = sql_context.read \
    .format("io.github.spark_redshift_community.spark.redshift") \
    .option("url", url) \
    .option("dbtable", "tableName") \
    .option("tempdir", "s3://path/for/temp/data") \
    .option("aws_iam_role", "aws_iam_role_arn") \
    .load()
df.write \
    .format("io.github.spark_redshift_community.spark.redshift") \
    .option("url", url) \
    .option("dbtable", "tableName_copy") \
    .option("tempdir", "s3://path/for/temp/data") \
    .option("aws_iam_role", "aws_iam_role_arn") \
    .mode("error") \
    .save()
```

SparkSQL

Use PySpark para ler e gravar dados de amostra de e para um banco de dados do Amazon Redshift usando o SparkSQL.

```
import boto3
import json
import sys
import os
from pyspark.sql import SparkSession

spark = SparkSession \
```

```
.builder \
    .enableHiveSupport() \
    .getOrCreate()
url = "jdbc:redshift:iam://redshifthost:5439/database"
aws_iam_role_arn = "arn:aws:iam::accountID:role/roleName"
bucket = "s3://path/for/temp/data"
tableName = "tableName" # Redshift table name
s = f"""CREATE TABLE IF NOT EXISTS { tableName} (country string, data string)
   USING io.github.spark_redshift_community.spark.redshift
   OPTIONS (dbtable '{tableName}', tempdir '{bucket}', url '{url}', aws_iam_role
 '{aws_iam_role_arn}' ); """
spark.sql(s)
columns = ["country" ,"data"]
data = [("test-country", "test-data")]
df = spark.sparkContext.parallelize(data).toDF(columns)
# Insert data into table
df.write.insertInto(tableName, overwrite=False)
df = spark.sql(f"SELECT * FROM {tableName}")
df.show()
```

Considerações e limitações ao usar o conector do Spark

O conector Spark oferece suporte a várias maneiras de gerenciar credenciais, configurar a segurança e se conectar a outros AWS serviços. Familiarize-se com as recomendações desta lista para configurar uma conexão funcional e resiliente.

- Recomendamos que você ative o SSL para a conexão JDBC do Spark no Amazon EMR para o Amazon Redshift.
- Recomendamos que você gerencie as credenciais do cluster do Amazon Redshift no AWS
 Secrets Manager como uma prática recomendada. Consulte <u>Usando AWS Secrets Manager para</u> recuperar credenciais para se conectar ao Amazon Redshift para ver um exemplo.
- Recomendamos que você transmita um perfil do IAM com o parâmetro aws_iam_role para o parâmetro de autenticação do Amazon Redshift.
- No momento, o parâmetro tempformat não é compatível com o formato Parquet.

Considerações 314

- O URI tempdir aponta para um local do Amazon S3. Esse diretório temporário não é limpo automaticamente e, portanto, pode incorrer em custos adicionais.
- Considere as seguintes recomendações para o Amazon Redshift:
 - Recomendamos bloquear o acesso público ao cluster do Amazon Redshift.
 - Recomendamos ativar o registro em log de auditoria do Amazon Redshift.
 - Recomendamos ativar a criptografia em repouso do Amazon Redshift.
- Considere as seguintes recomendações para o Amazon S3:
 - Recomendamos bloquear o acesso público aos buckets do Amazon S3.
 - Recomendamos usar a <u>criptografia do lado do servidor do Amazon S3</u> para criptografar os buckets do S3 que você usa.
 - Recomendamos usar as políticas de ciclo de vida do Amazon S3 para definir as regras de retenção para o bucket do S3.
 - O Amazon EMR sempre verifica o código importado do código aberto para a imagem. Por motivos de segurança, não oferecemos suporte AWS à codificação de chaves de acesso no tempdir URI como método de autenticação do Spark para o Amazon S3.

Para obter mais informações sobre como usar o conector e os parâmetros compatíveis, consulte os seguintes recursos:

- Integração do Amazon Redshift para Apache Spark no Guia de gerenciamento do Amazon Redshift.
- O repositório da comunidade spark-redshift no GitHub.

Uso do Volcano como um programador personalizado para Apache Spark no Amazon EMR no EKS

Com o Amazon EMR no EKS, você pode usar o operador do Spark ou o spark-submit para executar trabalhos do Spark com programadores personalizados do Kubernetes. Este tutorial aborda como executar trabalhos do Spark com um programador do Volcano em uma fila personalizada.

Visão geral

O <u>Volcano</u> pode ajudar a gerenciar a programação do Spark com funções avançadas, como a programação de filas, a programação de compartilhamento equitativo e a reserva de recursos. Para

Uso do Volcano 315

obter mais informações sobre os benefícios do Volcano, consulte Why Spark chooses Volcano as built-in batch scheduler on Kubernetes? no blog CNCF da Linux Foundation.

Instalação e configuração do Volcano

 Escolha um dos seguintes comandos kubectl para instalar o Volcano, com base em suas necessidades arquitetônicas:

```
# x86_64
kubectl apply -f https://raw.githubusercontent.com/volcano-sh/volcano/v1.5.1/
installer/volcano-development.yaml
# arm64:
kubectl apply -f https://raw.githubusercontent.com/volcano-sh/volcano/v1.5.1/
installer/volcano-development-arm64.yaml
```

 Prepare uma fila de exemplo do Volcano. Uma fila é uma coleção de. PodGroups A fila adota o FIFO e é a base para a divisão de recursos.

```
cat << EOF > volcanoQ.yaml
apiVersion: scheduling.volcano.sh/v1beta1
kind: Queue
metadata:
   name: sparkqueue
spec:
   weight: 4
   reclaimable: false
   capability:
        cpu: 10
        memory: 20Gi
EOF
kubectl apply -f volcanoQ.yaml
```

3. Faça upload de um exemplo de PodGroup manifesto para o Amazon S3. PodGroup é um grupo de grupos com forte associação. Você normalmente usa um PodGroup para agendamento em lote. Envie a amostra PodGroup a seguir para a fila que você definiu na etapa anterior.

```
cat << EOF > podGroup.yaml
apiVersion: scheduling.volcano.sh/v1beta1
kind: PodGroup
spec:
    # Set minMember to 1 to make a driver pod
```

Instalação 316

```
minMember: 1
# Specify minResources to support resource reservation.
# Consider the driver pod resource and executors pod resource.
# The available resources should meet the minimum requirements of the Spark job
# to avoid a situation where drivers are scheduled, but they can't schedule
# sufficient executors to progress.
minResources:
    cpu: "1"
    memory: "1Gi"
# Specify the queue. This defines the resource queue that the job should be
submitted to.
    queue: sparkqueue
EOF

aws s3 mv podGroup.yaml s3://bucket-name
```

Execução de uma aplicação do Spark usando o programador do Volcano com o operador do Spark

- Se ainda não o fez, conclua as etapas apresentadas nas seguintes seções para se preparar para usar:
 - a. Instalação e configuração do Volcano
 - b. Configuração do operador do Spark para o Amazon EMR no EKS
 - c. Instalação do operador do Spark

Inclua os seguintes argumentos ao executar o comando helm install sparkoperator-demo:

```
--set batchScheduler.enable=true
--set webhook.enable=true
```

 Crie um arquivo SparkApplication de definição spark-pi.yaml com batchScheduler configurado.

```
apiVersion: "sparkoperator.k8s.io/v1beta2"
kind: SparkApplication
metadata:
   name: spark-pi
   namespace: spark-operator
```

Envio: operador do Spark 317

```
spec:
 type: Scala
 mode: cluster
 image: "895885662937.dkr.ecr.us-west-2.amazonaws.com/spark/emr-6.10.0:latest"
 imagePullPolicy: Always
 mainClass: org.apache.spark.examples.SparkPi
 mainApplicationFile: "local:///usr/lib/spark/examples/jars/spark-examples.jar"
 sparkVersion: "3.3.1"
  batchScheduler: "volcano"
                              #Note: You must specify the batch scheduler name as
 'volcano'
 restartPolicy:
    type: Never
 volumes:
    - name: "test-volume"
      hostPath:
        path: "/tmp"
        type: Directory
 driver:
    cores: 1
    coreLimit: "1200m"
   memory: "512m"
   labels:
      version: 3.3.1
    serviceAccount: emr-containers-sa-spark
    volumeMounts:
      - name: "test-volume"
        mountPath: "/tmp"
  executor:
    cores: 1
    instances: 1
   memory: "512m"
   labels:
      version: 3.3.1
    volumeMounts:
      - name: "test-volume"
        mountPath: "/tmp"
```

 Envie a aplicação do Spark com o comando apresentado a seguir. Isso também cria um objeto SparkApplication chamado spark-pi:

```
kubectl apply -f spark-pi.yaml
```

4. Verifique os eventos do objeto SparkApplication com o seguinte comando:

Envio: operador do Spark 318

```
kubectl describe pods spark-pi-driver --namespace spark-operator
```

O primeiro evento do pod mostrará que o Volcano programou os pods:

Execução de uma aplicação do Spark usando o programador do Volcano com o **spark-submit**

- 1. Primeiro, conclua as etapas na seção <u>Configuração do spark-submit para o Amazon EMR no EKS</u>. Você deve desenvolver a distribuição do spark-submit com suporte do Volcano. Para obter mais informações, consulte a seção Build de <u>Using Volcano as Customized Scheduler for Spark on Kubernetes</u> na documentação do Apache Spark.
- 2. Defina os valores para as seguintes variáveis de ambiente:

```
export SPARK_HOME=spark-home
export MASTER_URL=k8s://Amazon-EKS-cluster-endpoint
```

3. Envie a aplicação do Spark com o seguinte comando:

```
$SPARK_HOME/bin/spark-submit \
--class org.apache.spark.examples.SparkPi \
--master $MASTER_URL \
--conf spark.kubernetes.container.image=895885662937.dkr.ecr.us-
west-2.amazonaws.com/spark/emr-6.10.0:latest \
--conf spark.kubernetes.authenticate.driver.serviceAccountName=spark \
--deploy-mode cluster \
--conf spark.kubernetes.namespace=spark-operator \
--conf spark.kubernetes.scheduler.name=volcano \
--conf spark.kubernetes.scheduler.volcano.podGroupTemplateFile=/path/to/podgroup-template.yaml \
--conf
spark.kubernetes.driver.pod.featureSteps=org.apache.spark.deploy.k8s.features.VolcanoFeatureSteps=org.apache.spark.deploy.k8s.features.VolcanoFeatureSteps=org.apache.spark.deploy.k8s.features.volcanoFeatureSteps=org.apache.spark.deploy.k8s.features.volcanoFeatureSteps=org.apache.spark.deploy.k8s.features.volcanoFeatureSteps=org.apache.spark.deploy.k8s.features.volcanoFeatureSteps=org.apache.spark.deploy.k8s.features.volcanoFeatureSteps=org.apache.spark.deploy.k8s.features.volcanoFeatureSteps=org.apache.spark.deploy.k8s.features.volcanoFeatureSteps=org.apache.spark.deploy.k8s.features.volcanoFeatureSteps=org.apache.spark.deploy.k8s.features.volcanoFeatureSteps=org.apache.spark.deploy.k8s.features.volcanoFeatureSteps=org.apache.spark.deploy.k8s.features.volcanoFeatureSteps=org.apache.spark.deploy.k8s.features.volcanoFeatureSteps=org.apache.spark.deploy.k8s.features.volcanoFeatureSteps=org.apache.spark.deploy.k8s.features.volcanoFeatureSteps=org.apache.spark.deploy.k8s.features.volcanoFeatureSteps=org.apache.spark.deploy.k8s.features.volcanoFeatureSteps=org.apache.spark.deploy.k8s.features.volcanoFeatureSteps=org.apache.spark.deploy.k8s.features.volcanoFeatureSteps=org.apache.spark.deploy.k8s.features.volcanoFeatureSteps=org.apache.spark.deploy.k8s.features.volcanoFeatureSteps=org.apache.spark.deploy.k8s.features.volcanoFeatureSteps=org.apache.spark.deploy.k8s.features.volcanoFeatureSteps=org.apache.spark.deploy.k8s.features.volcanoFeatureSteps=org.apache.spark.deploy.k8s.features.volcanoFeatureSteps=org.apach
```

```
--conf
spark.kubernetes.executor.pod.featureSteps=org.apache.spark.deploy.k8s.features.VolcanoFea
\
local:///usr/lib/spark/examples/jars/spark-examples.jar 20
```

4. Verifique os eventos do objeto SparkApplication com o seguinte comando:

```
kubectl describe pod spark-pi --namespace spark-operator
```

O primeiro evento do pod mostrará que o Volcano programou os pods:

Usando YuniKorn como um programador personalizado para o Apache Spark no Amazon EMR no EKS

Com o Amazon EMR no EKS, você pode usar o operador do Spark ou o spark-submit para executar trabalhos do Spark com programadores personalizados do Kubernetes. Este tutorial aborda como executar trabalhos do Spark com um YuniKorn agendador em uma fila personalizada e o agendamento de grupos.

Visão geral

O Apache YuniKorn pode ajudar a gerenciar o agendamento do Spark com agendamento compatível com aplicativos, para que você possa ter um controle refinado sobre as cotas e prioridades dos recursos. Com o agendamento em grupo, YuniKorn agenda um aplicativo somente quando a solicitação mínima de recursos para o aplicativo pode ser satisfeita. Para obter mais informações, consulte O que é agendamento de gangues no site de YuniKorn documentação do Apache.

Crie seu cluster e prepare-se para YuniKorn

Use as etapas a seguir para implantar um cluster do Amazon EKS. Você pode alterar a Região da AWS (region) e as zonas de disponibilidade (availabilityZones).

1. Defina o cluster do Amazon EKS:

Usando YuniKorn 320

```
cat <<EOF >eks-cluster.yaml
apiVersion: eksctl.io/v1alpha5
kind: ClusterConfig
metadata:
  name: emr-eks-cluster
  region: eu-west-1
vpc:
  clusterEndpoints:
    publicAccess: true
    privateAccess: true
  withOIDC: true
nodeGroups:
  - name: spark-jobs
    labels: { app: spark }
    instanceType: m5.xlarge
    desiredCapacity: 2
    minSize: 2
    maxSize: 3
    availabilityZones: ["eu-west-1a"]
E0F
```

Crie o cluster:

```
eksctl create cluster -f eks-cluster.yaml
```

3. Crie o namespace spark-job em que você executará o trabalho do Spark:

```
kubectl create namespace spark-job
```

- Em seguida, crie um perfil e uma associação de perfis do Kubernetes. Isso é obrigatório para a conta de serviço usada pela execução de trabalho do Spark.
 - a. Defina a conta de serviço, o perfil e a associação de perfis para os trabalhos do Spark.

```
cat <<EOF >emr-job-execution-rbac.yaml
---
```

Crie seu cluster do 321

```
apiVersion: v1
kind: ServiceAccount
metadata:
  name: spark-sa
  namespace: spark-job
automountServiceAccountToken: false
apiVersion: rbac.authorization.k8s.io/v1
kind: Role
metadata:
  name: spark-role
  namespace: spark-job
rules:
  - apiGroups: ["", "batch", "extensions"]
    resources: ["configmaps", "serviceaccounts", "events", "pods", "pods/
exec", "pods/log", "pods/
portforward", "secrets", "services", "persistentvolumeclaims"]
    verbs: ["create","delete","get","list","patch","update","watch"]
---
apiVersion: rbac.authorization.k8s.io/v1
kind: RoleBinding
metadata:
  name: spark-sa-rb
  namespace: spark-job
roleRef:
  apiGroup: rbac.authorization.k8s.io
  kind: Role
  name: spark-role
subjects:
  - kind: ServiceAccount
    name: spark-sa
    namespace: spark-job
E0F
```

b. Aplique a definição de perfil e de associação de perfis do Kubernetes com o seguinte comando:

```
kubectl apply -f emr-job-execution-rbac.yaml
```

Crie seu cluster do 322

Instalar e configurar YuniKorn

1. Use o seguinte comando kubectl para criar um namespace yunikorn para implantar o programador do Yunikorn:

```
kubectl create namespace yunikorn
```

2. Para instalar o programador, execute os seguintes comandos do Helm:

```
helm repo add yunikorn https://apache.github.io/yunikorn-release
```

helm repo update

helm install yunikorn yunikorn/yunikorn --namespace yunikorn

Execute um aplicativo Spark com YuniKorn agendador com o operador Spark

- Se ainda não o fez, conclua as etapas apresentadas nas seguintes seções para se preparar para usar:
 - a. Crie seu cluster e prepare-se para YuniKorn
 - b. Instalar e configurar YuniKorn
 - c. Configuração do operador do Spark para o Amazon EMR no EKS
 - d. Instalação do operador do Spark

Inclua os seguintes argumentos ao executar o comando helm install sparkoperator-demo:

```
--set batchScheduler.enable=true
--set webhook.enable=true
```

2. Crie um arquivo de definição SparkApplication spark-pi.yaml.

Instalar YuniKorn 323

Para usar YuniKorn como agendador para seus trabalhos, você deve adicionar certas anotações e rótulos à definição do seu aplicativo. As anotações e os rótulos especificam a fila do seu trabalho e a estratégia de programação que você deseja usar.

No exemplo a seguir, a anotação schedulingPolicyParameters determina o agendamento coletivo para a aplicação. Em seguida, o exemplo cria grupos de tarefas, ou "coletivos" de tarefas, para especificar a capacidade mínima que deve estar disponível antes da programação dos pods para iniciar a execução do trabalho. E, por fim, é especificado na definição do grupo de tarefas o uso de grupos de nós com o rótulo "app": "spark", conforme definido na seção Crie seu cluster e prepare-se para YuniKorn.

```
apiVersion: "sparkoperator.k8s.io/v1beta2"
kind: SparkApplication
metadata:
  name: spark-pi
  namespace: spark-job
spec:
  type: Scala
  mode: cluster
  image: "895885662937.dkr.ecr.us-west-2.amazonaws.com/spark/emr-6.10.0:latest"
  imagePullPolicy: Always
  mainClass: org.apache.spark.examples.SparkPi
 mainApplicationFile: "local:///usr/lib/spark/examples/jars/spark-examples.jar"
  sparkVersion: "3.3.1"
  restartPolicy:
    type: Never
  volumes:
    - name: "test-volume"
      hostPath:
        path: "/tmp"
        type: Directory
  driver:
    cores: 1
    coreLimit: "1200m"
   memory: "512m"
   labels:
      version: 3.3.1
    annotations:
      yunikorn.apache.org/schedulingPolicyParameters: "placeholderTimeoutSeconds=30
 gangSchedulingStyle=Hard"
      yunikorn.apache.org/task-group-name: "spark-driver"
```

Envio: operador do Spark 324

```
yunikorn.apache.org/task-groups: |-
      [{
          "name": "spark-driver",
          "minMember": 1,
          "minResource": {
            "cpu": "1200m",
            "memory": "1Gi"
          },
          "nodeSelector": {
            "app": "spark"
          }
        },
          "name": "spark-executor",
          "minMember": 1,
          "minResource": {
            "cpu": "1200m",
            "memory": "1Gi"
          },
          "nodeSelector": {
            "app": "spark"
          }
      }]
  serviceAccount: spark-sa
  volumeMounts:
    - name: "test-volume"
      mountPath: "/tmp"
executor:
 cores: 1
  instances: 1
 memory: "512m"
 labels:
    version: 3.3.1
  annotations:
    yunikorn.apache.org/task-group-name: "spark-executor"
 volumeMounts:
    - name: "test-volume"
      mountPath: "/tmp"
```

3. Envie a aplicação do Spark com o comando apresentado a seguir. Isso também cria um objeto SparkApplication chamado spark-pi:

```
kubectl apply -f spark-pi.yaml
```

Envio: operador do Spark 325

4. Verifique os eventos do objeto SparkApplication com o seguinte comando:

```
kubectl describe sparkapplication spark-pi --namespace spark-job
```

O primeiro evento do pod mostrará que YuniKorn programou os pods:

```
Type
       Reason
                         Age
                               From
                                                              Message
        ____
                                                              _____
Normal Scheduling
                        3m12s yunikorn
                                          spark-operator/org-apache-spark-examples-
sparkpi-2a777a88b98b8a95-driver is queued and waiting for allocation
Normal GangScheduling
                        3m12s yunikorn
                                         Pod belongs to the taskGroup spark-
driver, it will be scheduled as a gang member
Normal Scheduled
                        3m10s yunikorn
                                         Successfully assigned spark
Normal PodBindSuccessful 3m10s yunikorn
                                         Pod spark-operator/
Normal TaskCompleted
                        2m3s yunikorn
                                         Task spark-operator/
Normal Pulling
                        3m10s kubelet
                                         Pulling
```

Execute um aplicativo Spark com YuniKorn agendador com **spark-submit**

- 1. Primeiro, conclua as etapas na seção <u>Configuração do spark-submit para o Amazon EMR no</u> EKS.
- 2. Defina os valores para as seguintes variáveis de ambiente:

```
export SPARK_HOME=spark-home
export MASTER_URL=k8s://Amazon-EKS-cluster-endpoint
```

3. Envie a aplicação do Spark com o seguinte comando:

No exemplo a seguir, a anotação schedulingPolicyParameters determina o agendamento coletivo para a aplicação. Em seguida, o exemplo cria grupos de tarefas, ou "coletivos" de tarefas, para especificar a capacidade mínima que deve estar disponível antes da programação dos pods para iniciar a execução do trabalho. E, por fim, é especificado na definição do grupo de tarefas o uso de grupos de nós com o rótulo "app": "spark", conforme definido na seção Crie seu cluster e prepare-se para YuniKorn.

```
$SPARK_HOME/bin/spark-submit \
--class org.apache.spark.examples.SparkPi \
```

```
--master $MASTER_URL \
 --conf spark.kubernetes.container.image=895885662937.dkr.ecr.us-
west-2.amazonaws.com/spark/emr-6.10.0:latest \
 --conf spark.kubernetes.authenticate.driver.serviceAccountName=spark-sa \
 --deploy-mode cluster \
 --conf spark.kubernetes.namespace=spark-job \
--conf spark.kubernetes.scheduler.name=yunikorn \
 --conf spark.kubernetes.driver.annotation.yunikorn.apache.org/
schedulingPolicyParameters="placeholderTimeoutSeconds=30 gangSchedulingStyle=Hard"
 --conf spark.kubernetes.driver.annotation.yunikorn.apache.org/task-group-
name="spark-driver" \
 --conf spark.kubernetes.executor.annotation.yunikorn.apache.org/task-group-
name="spark-executor" \
 --conf spark.kubernetes.driver.annotation.yunikorn.apache.org/task-groups='[{
            "name": "spark-driver",
            "minMember": 1,
            "minResource": {
              "cpu": "1200m",
              "memory": "1Gi"
            },
            "nodeSelector": {
              "app": "spark"
            }
         },
            "name": "spark-executor",
            "minMember": 1,
            "minResource": {
              "cpu": "1200m",
              "memory": "1Gi"
            },
            "nodeSelector": {
              "app": "spark"
        }1' \
local:///usr/lib/spark/examples/jars/spark-examples.jar 20
```

4. Verifique os eventos do objeto SparkApplication com o seguinte comando:

```
kubectl describe pod spark-driver-pod --namespace spark-job
```

O primeiro evento do pod mostrará que YuniKorn programou os pods:

Type Reason Age From Message ____ --------spark-operator/org-apache-spark-examples-Normal Scheduling 3m12s yunikorn sparkpi-2a777a88b98b8a95-driver is queued and waiting for allocation Normal GangScheduling 3m12s yunikorn Pod belongs to the taskGroup sparkdriver, it will be scheduled as a gang member Normal Scheduled 3m10s yunikorn Successfully assigned spark Normal PodBindSuccessful 3m10s yunikorn Pod spark-operator/ Normal TaskCompleted 2m3s yunikorn Task spark-operator/ Normal Pulling 3m10s kubelet Pulling

Segurança no Amazon EMR no EKS

A segurança na nuvem AWS é a maior prioridade. Como AWS cliente, você se beneficia de data centers e arquiteturas de rede criados para atender aos requisitos das organizações mais sensíveis à segurança.

A segurança é uma responsabilidade compartilhada entre você AWS e você. O <u>modelo de</u> responsabilidade compartilhada descreve isso como segurança da nuvem e segurança na nuvem:

- Segurança da nuvem AWS é responsável por proteger a infraestrutura que executa AWS os serviços na AWS nuvem. AWS também fornece serviços que você pode usar com segurança. Auditores terceirizados testam e verificam regularmente a eficácia de nossa segurança como parte dos Programas de Conformidade Programas de <u>AWS</u> de . Para saber mais sobre os programas de conformidade que se aplicam ao Amazon EMR, consulte <u>AWS Serviços no escopo do programa</u> de conformidade AWS Serviços no escopo do programa .
- Segurança na nuvem Sua responsabilidade é determinada pelo AWS serviço que você usa.
 Você também é responsável por outros fatores, incluindo a confidencialidade de seus dados, os requisitos da empresa e as leis e regulamentos aplicáveis.

Esta documentação ajuda você a entender como aplicar o modelo de responsabilidade compartilhada ao usar o Amazon EMR no EKS. Os tópicos a seguir demonstram como configurar o Amazon EMR no EKS para atender aos seus objetivos de segurança e de conformidade. Você também aprende a usar outros AWS serviços que ajudam a monitorar e proteger seus recursos do Amazon EMR no EKS.

Tópicos

- Práticas recomendadas de segurança para o Amazon EMR no EKS
- Proteção de dados
- Gerenciamento de Identidade e Acesso
- Usando o Amazon EMR no EKS com AWS Lake Formation para controle de acesso refinado
- Registro em log e monitoramento
- Uso do Amazon S3 Access Grants com o Amazon EMR no EKS
- Validação de conformidade para o Amazon EMR no EKS
- Resiliência no Amazon EMR no EKS

- Segurança da infraestrutura no Amazon EMR no EKS
- Análise de configuração e vulnerabilidade
- Conexão com o Amazon EMR no EKS usando um endpoint da VPC de interface
- Configuração do acesso entre contas para o Amazon EMR no EKS

Práticas recomendadas de segurança para o Amazon EMR no **EKS**

O Amazon EMR no EKS disponibiliza diversos recursos de segurança a serem considerados no desenvolvimento e na implementação das suas próprias políticas de segurança. As práticas recomendadas a seguir são diretrizes gerais e não representam uma solução completa de segurança. Como essas práticas recomendadas podem não ser adequadas ou suficientes para o seu ambiente, trate-as como considerações úteis em vez de prescrições.



Note

Para obter mais práticas recomendadas de segurança, consulte Práticas recomendadas de segurança para o Amazon EMR no EKS.

Aplicação do princípio de privilégio mínimo

O Amazon EMR no EKS fornece uma política de acesso granular para aplicações que usam perfis do IAM, como perfis de execução. Esses perfis de execução são mapeados para as contas de serviço do Kubernetes por meio da política de confiança do perfil do IAM. O Amazon EMR no EKS cria pods em um namespace registrado do Amazon EKS que executa o código da aplicação fornecido pelo usuário. Os pods de trabalho que executam o código do aplicativo assumem a função de execução ao se conectar a outros AWS serviços. Recomendamos que os perfis de execução recebam somente o conjunto de privilégios mínimos obrigatórios para o trabalho, como a cobertura da aplicação e o acesso ao destino do log. Também recomendamos auditar regularmente as permissões dos trabalhos e após qualquer alteração no código da aplicação.

Listagem de controle de acesso para endpoints

Os endpoints gerenciados podem ser criados somente para clusters do EKS que foram configurados para usar, no mínimo, uma sub-rede privada em sua VPC. Essa configuração restringe o acesso aos

Práticas recomendadas 330 balanceadores de carga criados por endpoints gerenciados para que eles possam ser acessados somente usando a VPC. Para aumentar ainda mais a segurança, recomendamos configurar grupos de segurança com esses balanceadores de carga para que eles possam restringir o tráfego de entrada a um conjunto selecionado de endereços IP.

Obtenção das atualizações de segurança mais recentes para as imagens personalizadas

Para usar imagens personalizadas com o Amazon EMR no EKS, você pode instalar quaisquer binários e bibliotecas na imagem. Você é responsável pela aplicação de patches de segurança aos binários adicionados à imagem. As imagens do Amazon EMR no EKS têm aplicações regulares dos patches de segurança mais recentes. Para obter a imagem mais recente, você deve criar novamente as imagens personalizadas sempre que houver uma nova versão da imagem base da versão do Amazon EMR. Para obter mais informações, consulte Versões do Amazon EMR no EKS e Detalhes sobre como selecionar um URI de imagem base.

Limitação do acesso à credencial do pod

O Kubernetes oferece suporte a diversos métodos de atribuição de credenciais para um pod. O provisionamento de múltiplos provedores de credenciais pode aumentar a complexidade do seu modelo de segurança. O Amazon EMR no EKS adotou o uso de perfis do IAM para contas de serviços (IRSA) como um provedor de credencial padrão em um namespace do EKS registrado. Outros métodos não são compatíveis, incluindo kube2iam, kiam e o uso de um perfil de EC2 instância da instância em execução no cluster.

Isolamento de código de uma aplicação não confiável

O Amazon EMR no EKS não inspeciona a integridade do código de uma aplicação enviada pelos usuários do sistema. Se você estiver executando um cluster virtual multilocatário configurado usando diversos perfis de execução, que podem ser usados para enviar trabalhos por locatários não confiáveis que executam códigos arbitrários, há o risco de uma aplicação mal-intencionada escalar os privilégios. Nesta situação, considere isolar perfis de execução com privilégios semelhantes em um cluster virtual diferente.

Permissões de controle de acesso por perfil (RBAC)

Os administradores devem controlar rigorosamente as permissões de controle de acesso por perfil (RBAC) para o Amazon EMR em namespaces gerenciados pelo EKS. No mínimo, as permissões

apresentadas a seguir não devem ser concedidas aos emissores de trabalhos no Amazon EMR em namespaces gerenciados pelo EKS.

- Permissões de RBAC do Kubernetes para modificar o configmap: porque o Amazon EMR no EKS usa os configmaps do Kubernetes para gerar modelos de pod gerenciados que têm o nome da conta de serviço gerenciada. Este atributo não deve sofrer mutação.
- Permissões de RBAC do Kubernetes para executar em pods do Amazon EMR no EKS: para evitar conceder acesso a modelos de pod gerenciados que têm o nome da conta de serviço gerenciada. Este atributo não deve sofrer mutação. Essa permissão também pode conceder acesso ao token JWT montado no pod, que pode ser usado para recuperar as credenciais do perfil de execução.
- Permissões RBAC do Kubernetes para criar pods para impedir que os usuários criem pods usando um Kubernetes ServiceAccount que pode ser mapeado para uma função do IAM com mais privilégios do que o usuário. AWS
- Permissões RBAC do Kubernetes para implantar webhook mutante para impedir que os usuários usem o webhook mutante para alterar o nome do Kubernetes para pods criados pelo Amazon EMR no EKS. ServiceAccount
- Permissões de RBAC do Kubernetes para ler segredos do Kubernetes: para evitar que os usuários façam a leitura de dados confidenciais armazenados nesses segredos.

Restrição do acesso às credenciais do perfil do IAM ou do perfil de instância do grupo de nós

- Recomendamos que você atribua AWS permissões mínimas às funções do IAM do nodegroup.
 Isso ajuda a evitar a escalação de privilégios por código que pode ser executado usando credenciais do perfil de instância de nós de processamento do EKS.
- Para bloquear completamente o acesso às credenciais do perfil de instância para todos os pods executados no Amazon EMR em namespaces gerenciados pelo EKS, recomendamos que você execute comandos iptables em nós do EKS. Para obter mais informações, consulte Restringir o acesso às credenciais do perfil da EC2 instância da Amazon. No entanto, é importante definir o escopo adequado para os perfis do IAM da sua conta de serviço para que os pods tenham todas as permissões necessárias. Por exemplo, o perfil do IAM do nó recebe permissões para extrair imagens de contêiner do Amazon ECR. Se essas permissões não forem atribuídas a um pod, o pod não poderá extrair imagens de contêiner do Amazon ECR. O plug-in CNI da VPC também precisa ser atualizado. Para obter mais informações, consulte Walkthrough: Updating the VPC CNI plugin to use IAM roles for service accounts.

Proteção de dados

O <u>modelo de responsabilidade AWS compartilhada</u> se aplica à proteção de dados no Amazon EMR no EKS. Conforme descrito neste modelo, AWS é responsável por proteger a infraestrutura global que executa toda a AWS nuvem. Você é responsável por manter o controle sobre o conteúdo hospedado nessa infraestrutura. Esse conteúdo inclui as tarefas de configuração e gerenciamento de segurança dos serviços da AWS que você usa. Para obter mais informações sobre a privacidade de dados, consulte as <u>Perguntas frequentes sobre privacidade de dados</u>. Para obter informações sobre proteção de dados na Europa, consulte <u>a postagem do blog sobre o Modelo de Responsabilidade AWS Compartilhada e o GDPR</u> no Blog AWS de Segurança.

Para fins de proteção de dados, recomendamos que você proteja as credenciais da AWS conta e configure contas individuais com AWS Identity and Access Management (IAM). Dessa maneira, cada usuário receberá apenas as permissões necessárias para cumprir suas obrigações de trabalho. Recomendamos também que você proteja seus dados das seguintes formas:

- Use uma autenticação multifator (MFA) com cada conta.
- Use SSL/TLS para se comunicar com AWS os recursos. Recomendamos usar o TLS 1.2 ou posterior.
- Configure a API e o registro de atividades do usuário com AWS CloudTrail.
- Use soluções AWS de criptografia, juntamente com todos os controles de segurança padrão nos AWS serviços.
- Use serviços gerenciados de segurança avançada, como o Amazon Macie, que ajuda a localizar e proteger dados pessoais armazenados no Amazon S3.
- Use as opções de criptografia do Amazon EMR no EKS para criptografar dados em repouso e em trânsito.
- Se você precisar de módulos criptográficos validados pelo FIPS 140-2 ao acessar AWS por meio de uma interface de linha de comando ou de uma API, use um endpoint FIPS. Para obter mais informações sobre endpoints do FIPS, consulte <u>Federal Information Processing Standard (FIPS)</u> 140-2.

É altamente recomendável que você nunca coloque informações de identificação confidenciais, como números de conta dos seus clientes, em campos de formato livre, como um campo Nome. Isso inclui quando você trabalha com o Amazon EMR no EKS ou em outros AWS serviços usando o console, a API ou. AWS CLI AWS SDKs Todos os dados inseridos no Amazon EMR no EKS ou em

Proteção de dados 333

outros serviços podem ser coletados para a inclusão em logs de diagnóstico. Ao fornecer um URL para um servidor externo, não inclua informações de credenciais no URL para validar a solicitação a esse servidor.

Criptografia em repouso

A criptografia de dados ajuda a impedir que usuários não autorizados leiam dados em um cluster e em sistemas de armazenamento físico de dados associados. Isso inclui dados salvos em mídias persistentes, conhecidos como dados em repouso, e dados que podem ser interceptados enquanto viajam pela rede, conhecidos como dados em trânsito.

A criptografia de dados requer chaves e certificados. Você pode escolher entre várias opções, incluindo chaves gerenciadas por AWS Key Management Service, chaves gerenciadas pelo Amazon S3 e chaves e certificados de fornecedores personalizados fornecidos por você. Ao usar AWS KMS como seu provedor de chaves, cobranças se aplicam pelo armazenamento e uso de chaves de criptografia. Para obter mais informações, consulte AWS KMS Preço.

Antes de especificar as opções de criptografia, decida quais sistemas de gerenciamento de chaves e de certificados você deseja usar. Em seguida, crie as chaves e os certificados para os provedores personalizados especificados como parte das configurações de criptografia.

Criptografia em repouso para dados do EMRFS no Amazon S3

A criptografia do Amazon S3 funciona com objetos do Sistema de Arquivos do EMR (EMRFS) lidos e gravados no Amazon S3. Você especifica a criptografia do lado do servidor (SSE) ou a criptografia do lado do cliente (CSE) do Amazon S3 como o modo de criptografia padrão ao habilitar a criptografia em repouso. Opcionalmente, você pode especificar diferentes métodos de criptografia para buckets individuais usando Per bucket encryption overrides (Substituições de criptografia por bucket). Independentemente de a criptografia do Amazon S3 estar habilitada, o Transport Layer Security (TLS) criptografa os objetos do EMRFS em trânsito entre os nós do cluster do EMR e o Amazon S3. Para obter informações detalhadas sobre a criptografia do Amazon S3, consulte Proteger dados com criptografia no Guia do desenvolvedor do Amazon Simple Storage Service.



Note

Quando você usa AWS KMS, cobranças são cobradas pelo armazenamento e uso de chaves de criptografia. Para obter mais informações, consulte AWS KMS Preço.

Criptografia em repouso 334

Criptografia do lado do servidor do Amazon S3

Quando você configura a criptografia do lado do servidor do Amazon S3, o Amazon S3 criptografa os dados no nível do objeto à medida que os grava no disco e os descriptografa quando são acessados. Para obter mais informações sobre a SSE, consulte Proteger os dados usando criptografia do lado do servidor no Guia do desenvolvedor do Amazon Simple Storage Service.

Você pode escolher entre dois sistemas de gerenciamento de chaves diferentes ao especificar a SSE no Amazon EMR no EKS:

- SSE-S3: o Amazon S3 gerencia as chaves para você.
- SSE-KMS Você usa um AWS KMS key para configurar políticas adequadas para o Amazon EMR no EKS.

A SSE com chaves fornecidas pelo cliente (SSE-C) não está disponível para o uso com o Amazon EMR no EKS.

Criptografia do lado do cliente do Amazon S3

Com a criptografia do lado do cliente do Amazon S3, a criptografia e a descriptografia do Amazon S3 ocorrem no cliente do EMRFS em seu cluster. Os objetos são criptografados antes de serem carregados no Amazon S3 e descriptografados após serem baixados. O provedor especificado por você fornece a chave de criptografia que o cliente usa. O cliente pode usar chaves fornecidas pelo AWS KMS (CSE-KMS) ou uma classe Java personalizada que fornece a chave raiz do lado do cliente (CSE-C). As especificações de criptografia são ligeiramente diferentes entre a CSE-KMS e a CSE-C, dependendo do provedor especificado e dos metadados do objeto que está sendo descriptografado ou criptografado. Para obter mais informações sobre essas diferenças, consulte Proteger dados usando a criptografia do lado do cliente no Guia do desenvolvedor do Amazon Simple Storage Service.



A CSE do Amazon S3 garante somente que os dados do EMRFS trocados com o Amazon S3 sejam criptografados. Não são todos os dados nos volumes de instâncias do cluster que são criptografados. Além disso, como o Hue não usa o EMRFS, os objetos que o navegador de arquivos do S3 para Hue grava no Amazon S3 não são criptografados.

Criptografia em repouso 335

Criptografia de disco local

O Apache Spark oferece suporte à criptografia de dados temporários gravados em discos locais. Isso abrange arquivos embaralhados, vazamentos embaralhados e blocos de dados armazenados em disco para variáveis de armazenamento em cache e de transmissão. Não abrange a criptografia de dados de saída gerados por aplicativos com APIs saveAsHadoopFile ousaveAsTable. Também é possível que não haja abrangência para arquivos temporários criados explicitamente pelo usuário. Para obter mais informações, consulte Local Storage Encryption na documentação do Spark. O Spark não oferece suporte a dados criptografados no disco local, como dados intermediários gravados em um disco local por um processo executor quando os dados não se encaixam na memória. Os dados que persistem no disco têm como escopo o runtime do trabalho, e a chave usada para criptografar os dados é gerada dinamicamente pelo Spark para cada execução de trabalho. Depois que o trabalho do Spark for encerrado, nenhum outro processo poderá descriptografar os dados.

Para pod de drivers e de executores, você criptografa dados em repouso que persistem no volume montado. Há três opções diferentes de armazenamento AWS nativo que você pode usar com o Kubernetes: EBS, FSx EFS e for Lustre. Todos os três oferecem criptografia em repouso usando uma chave gerenciada pelo serviço ou uma AWS KMS key. Para obter mais informações, consulte EKS Best Practices Guide. Com esta abordagem, todos os dados que persistem no volume montado são criptografados.

Gerenciamento de chaves

Você pode configurar o KMS para alternar automaticamente suas chaves do KMS. Isso alterna suas chaves uma vez por ano, enquanto salva as chaves antigas indefinidamente para que seus dados ainda possam ser descriptografados. Para obter informações adicionais, consulte Rotação. AWS KMS keys

Criptografia em trânsito

Vários mecanismos de criptografia estão habilitados com a criptografia em trânsito. Esses são recursos de código aberto, específicos da aplicação e podem variar de acordo com a versão do Amazon EMR no EKS. Os seguintes recursos de criptografia específicos da aplicação podem ser habilitados com o Amazon EMR no EKS:

- Spark
 - A comunicação de RPC interna entre os componentes do Spark, como o serviço de transferência de blocos e o serviço de shuffle externo, é criptografada usando a criptografia

Criptografia em trânsito 336

AES-256 nas versões 5.9.0 e posteriores do Amazon EMR. Em versões anteriores, a comunicação RPC interna é criptografada usando SASL com DIGEST- MD5 como cifra.

 A comunicação do protocolo HTTP com interfaces de usuário, como o Spark History Server e servidores de arquivos habilitados para HTTPS é criptografada usando a configuração de SSL do Spark. Para obter mais informações, consulte <u>Configuração do SSL</u> na documentação do Spark.

Para obter mais informações, consulte as configurações de segurança do Spark.

- Você deve permitir somente conexões criptografadas via HTTPS (TLS) usando <u>a SecureTransport</u> condição aws: nas políticas do Amazon S3 bucket IAM.
- Os resultados de consultas que transmitem para clientes JDBC ou ODBC s\u00e3o criptografados usando TLS.

Gerenciamento de Identidade e Acesso

AWS Identity and Access Management (IAM) é uma ferramenta AWS service (Serviço da AWS) que ajuda o administrador a controlar com segurança o acesso aos AWS recursos. Os administradores do IAM controlam quem pode ser autenticado (iniciar sessão) e autorizado (ter permissões) para usar os recursos do Amazon EMR no EKS. O IAM é um AWS service (Serviço da AWS) que você pode usar sem custo adicional.

Tópicos

- Público
- Autenticação com identidades
- Gerenciar o acesso usando políticas
- Como o Amazon EMR no EKS funciona com o IAM
- Uso de perfis vinculados ao serviço para o Amazon EMR no EKS
- Políticas gerenciadas para o Amazon EMR no EKS
- Uso de perfis de execução de trabalho com o Amazon EMR no EKS
- Exemplos de políticas baseadas em identidade para o Amazon EMR no EKS
- Políticas para controle de acesso baseado em etiquetas
- Solução de problemas de identidade e acesso do Amazon EMR no EKS

Público

A forma como você usa AWS Identity and Access Management (IAM) difere, dependendo do trabalho que você faz no Amazon EMR no EKS.

Usuário do serviço: se você usar o serviço Amazon EMR no EKS para realizar seu trabalho, o administrador fornecerá as credenciais e as permissões de que você precisa. À medida que você usa mais recursos do Amazon EMR no EKS para realizar o trabalho, talvez sejam necessárias permissões adicionais. Entender como o acesso é gerenciado pode ajudar a solicitar as permissões corretas ao administrador. Se você não conseguir acessar um recurso no Amazon EMR no EKS, consulte Solução de problemas de identidade e acesso do Amazon EMR no EKS.

Administrador do serviço: se você for responsável pelos recursos do Amazon EMR no EKS em sua empresa, provavelmente terá acesso total ao Amazon EMR no EKS. É sua função determinar quais atributos e recursos do Amazon EMR no EKS seus usuários do serviço devem acessar. Envie as solicitações ao administrador do IAM para alterar as permissões dos usuários de serviço. Revise as informações nesta página para compreender os conceitos básicos do IAM. Para saber mais sobre como sua empresa pode usar o IAM com o Amazon EMR no EKS, consulte Como o Amazon EMR no EKS funciona com o IAM.

Administrador do IAM: se você for um administrador do IAM, talvez deseje saber detalhes sobre como escrever políticas para gerenciar o acesso ao Amazon EMR no EKS. Para visualizar exemplos de políticas baseadas em identidade do Amazon EMR no EKS que podem ser usadas no IAM, consulte Exemplos de políticas baseadas em identidade para o Amazon EMR no EKS.

Autenticação com identidades

A autenticação é como você faz login AWS usando suas credenciais de identidade. Você deve estar autenticado (conectado AWS) como o Usuário raiz da conta da AWS, como usuário do IAM ou assumindo uma função do IAM.

Você pode entrar AWS como uma identidade federada usando credenciais fornecidas por meio de uma fonte de identidade. AWS IAM Identity Center Usuários (IAM Identity Center), a autenticação de login único da sua empresa e suas credenciais do Google ou do Facebook são exemplos de identidades federadas. Quando você faz login como identidade federada, o administrador já configurou anteriormente a federação de identidades usando perfis do IAM. Ao acessar AWS usando a federação, você está assumindo indiretamente uma função.

Público 338

Dependendo do tipo de usuário que você é, você pode entrar no AWS Management Console ou no portal de AWS acesso. Para obter mais informações sobre como fazer login em AWS, consulte Como fazer login Conta da AWS no Guia do Início de Sessão da AWS usuário.

Se você acessar AWS programaticamente, AWS fornece um kit de desenvolvimento de software (SDK) e uma interface de linha de comando (CLI) para assinar criptograficamente suas solicitações usando suas credenciais. Se você não usa AWS ferramentas, você mesmo deve assinar as solicitações. Para obter mais informações sobre como usar o método recomendado para designar solicitações por conta própria, consulte Versão 4 do AWS Signature para solicitações de API no Guia do usuário do IAM.

Independente do método de autenticação usado, também pode ser necessário fornecer informações adicionais de segurança. Por exemplo, AWS recomenda que você use a autenticação multifator (MFA) para aumentar a segurança da sua conta. Para saber mais, consulte <u>Autenticação multifator</u> no Guia do usuário do AWS IAM Identity Center e <u>Usar a autenticação multifator da AWS no IAM</u> no Guia do usuário do IAM.

Conta da AWS usuário root

Ao criar uma Conta da AWS, você começa com uma identidade de login que tem acesso completo a todos Serviços da AWS os recursos da conta. Essa identidade é chamada de usuário Conta da AWS raiz e é acessada fazendo login com o endereço de e-mail e a senha que você usou para criar a conta. É altamente recomendável não usar o usuário-raiz para tarefas diárias. Proteja as credenciais do usuário-raiz e use-as para executar as tarefas que somente ele puder executar. Para obter a lista completa das tarefas que exigem login como usuário-raiz, consulte Tarefas que exigem credenciais de usuário-raiz no Guia do Usuário do IAM.

Identidade federada

Como prática recomendada, exija que usuários humanos, incluindo usuários que precisam de acesso de administrador, usem a federação com um provedor de identidade para acessar Serviços da AWS usando credenciais temporárias.

Uma identidade federada é um usuário do seu diretório de usuários corporativo, de um provedor de identidade da web AWS Directory Service, do diretório do Identity Center ou de qualquer usuário que acesse usando credenciais fornecidas Serviços da AWS por meio de uma fonte de identidade. Quando as identidades federadas são acessadas Contas da AWS, elas assumem funções, e as funções fornecem credenciais temporárias.

Autenticação com identidades 339

Para o gerenciamento de acesso centralizado, é recomendável usar o AWS IAM Identity Center. Você pode criar usuários e grupos no IAM Identity Center ou pode se conectar e sincronizar com um conjunto de usuários e grupos em sua própria fonte de identidade para uso em todos os seus Contas da AWS aplicativos. Para obter mais informações sobre o Centro de Identidade do IAM, consulte O que é o Centro de Identidade do IAM? no Guia do Usuário do AWS IAM Identity Center .

Usuários e grupos do IAM

Um <u>usuário do IAM</u> é uma identidade dentro da sua Conta da AWS que tem permissões específicas para uma única pessoa ou aplicativo. Sempre que possível, é recomendável contar com credenciais temporárias em vez de criar usuários do IAM com credenciais de longo prazo, como senhas e chaves de acesso. No entanto, se você tiver casos de uso específicos que exijam credenciais de longo prazo com usuários do IAM, é recomendável alternar as chaves de acesso. Para obter mais informações, consulte <u>Alternar as chaves de acesso regularmente para casos de uso que exijam</u> credenciais de longo prazo no Guia do Usuário do IAM.

Um grupo do IAM é uma identidade que especifica uma coleção de usuários do IAM. Não é possível fazer login como um grupo. É possível usar grupos para especificar permissões para vários usuários de uma vez. Os grupos facilitam o gerenciamento de permissões para grandes conjuntos de usuários. Por exemplo, você pode ter um grupo chamado IAMAdminse conceder a esse grupo permissões para administrar recursos do IAM.

Usuários são diferentes de perfis. Um usuário é exclusivamente associado a uma pessoa ou a uma aplicação, mas um perfil pode ser assumido por qualquer pessoa que precisar dele. Os usuários têm credenciais permanentes de longo prazo, mas os perfis fornecem credenciais temporárias. Para saber mais, consulte <u>Casos de uso para usuários do IAM</u> no Guia do usuário do IAM.

Perfis do IAM

Uma <u>função do IAM</u> é uma identidade dentro da sua Conta da AWS que tem permissões específicas. Ele é semelhante a um usuário do IAM, mas não está associado a uma pessoa específica. Para assumir temporariamente uma função do IAM no AWS Management Console, você pode <u>alternar</u> <u>de um usuário para uma função do IAM (console)</u>. Você pode assumir uma função chamando uma operação de AWS API AWS CLI ou usando uma URL personalizada. Para obter mais informações sobre métodos para usar perfis, consulte Métodos para assumir um perfil no Guia do usuário do IAM.

Perfis do IAM com credenciais temporárias são úteis nas seguintes situações:

 Acesso de usuário federado: para atribuir permissões a identidades federadas, é possível criar um perfil e definir permissões para ele. Quando uma identidade federada é autenticada,

Autenticação com identidades 340

essa identidade é associada ao perfil e recebe as permissões definidas por ele. Para ter mais informações sobre perfis para federação, consulte <u>Criar um perfil para um provedor de identidade de terceiros (federação)</u> no Guia do usuário do IAM. Se usar o Centro de Identidade do IAM, configure um conjunto de permissões. Para controlar o que suas identidades podem acessar após a autenticação, o Centro de Identidade do IAM correlaciona o conjunto de permissões a um perfil no IAM. Para obter informações sobre conjuntos de permissões, consulte <u>Conjuntos de Permissões</u> no Guia do Usuário do AWS IAM Identity Center .

- Permissões temporárias para usuários do IAM: um usuário ou um perfil do IAM pode presumir um perfil do IAM para obter temporariamente permissões diferentes para uma tarefa específica.
- Acesso entre contas: é possível usar um perfil do IAM para permitir que alguém (uma entidade principal confiável) em outra conta acesse recursos em sua conta. Os perfis são a principal forma de conceder acesso entre contas. No entanto, com alguns Serviços da AWS, você pode anexar uma política diretamente a um recurso (em vez de usar uma função como proxy). Para conhecer a diferença entre perfis e políticas baseadas em recurso para acesso entre contas, consulte <u>Acesso a recursos entre contas no IAM</u> no Guia do usuário do IAM.
- Acesso entre serviços Alguns Serviços da AWS usam recursos em outros Serviços da AWS.
 Por exemplo, quando você faz uma chamada em um serviço, é comum que esse serviço execute aplicativos na Amazon EC2 ou armazene objetos no Amazon S3. Um serviço pode fazer isso usando as permissões da entidade principal da chamada, usando um perfil de serviço ou um perfil vinculado ao serviço.
 - Sessões de acesso direto (FAS) Quando você usa um usuário ou uma função do IAM para realizar ações AWS, você é considerado principal. Ao usar alguns serviços, você pode executar uma ação que inicia outra ação em um serviço diferente. O FAS usa as permissões do diretor chamando um AWS service (Serviço da AWS), combinadas com a solicitação AWS service (Serviço da AWS) para fazer solicitações aos serviços posteriores. As solicitações do FAS são feitas somente quando um serviço recebe uma solicitação que requer interações com outros Serviços da AWS ou com recursos para ser concluída. Nesse caso, você precisa ter permissões para executar ambas as ações. Para obter detalhes da política ao fazer solicitações de FAS, consulte Sessões de acesso direto.
 - Perfil de serviço: um perfil de serviço é um perfil do IAM que um serviço assume para executar
 ações em seu nome. Um administrador do IAM pode criar, modificar e excluir um perfil de
 serviço do IAM. Para obter mais informações, consulte Criar um perfil para delegar permissões a
 um AWS service (Serviço da AWS) no Guia do Usuário do IAM.
 - Função vinculada ao serviço Uma função vinculada ao serviço é um tipo de função de serviço vinculada a um. AWS service (Serviço da AWS) O serviço pode presumir o perfil para executar

Autenticação com identidades 341

uma ação em seu nome. As funções vinculadas ao serviço aparecem em você Conta da AWS e são de propriedade do serviço. Um administrador do IAM pode visualizar, mas não editar as permissões para perfis vinculados a serviço.

Aplicativos em execução na Amazon EC2 — Você pode usar uma função do IAM para gerenciar credenciais temporárias para aplicativos que estão sendo executados em uma EC2 instância e fazendo solicitações AWS CLI de AWS API. Isso é preferível a armazenar chaves de acesso na EC2 instância. Para atribuir uma AWS função a uma EC2 instância e disponibilizá-la para todos os aplicativos, você cria um perfil de instância anexado à instância. Um perfil de instância contém a função e permite que programas em execução na EC2 instância recebam credenciais temporárias. Para obter mais informações, consulte <u>Usar uma função do IAM para conceder permissões a aplicativos executados em EC2 instâncias da Amazon</u> no Guia do usuário do IAM.

Gerenciar o acesso usando políticas

Você controla o acesso AWS criando políticas e anexando-as a AWS identidades ou recursos. Uma política é um objeto AWS que, quando associada a uma identidade ou recurso, define suas permissões. AWS avalia essas políticas quando um principal (usuário, usuário raiz ou sessão de função) faz uma solicitação. As permissões nas políticas determinam se a solicitação será permitida ou negada. A maioria das políticas é armazenada AWS como documentos JSON. Para obter mais informações sobre a estrutura e o conteúdo de documentos de políticas JSON, consulte Visão geral das políticas JSON no Guia do usuário do IAM.

Os administradores podem usar políticas AWS JSON para especificar quem tem acesso ao quê. Ou seja, qual entidade principal pode executar ações em quais recursos e em que condições.

Por padrão, usuários e perfis não têm permissões. Para conceder permissão aos usuários para executar ações nos recursos que eles precisam, um administrador do IAM pode criar políticas do IAM. O administrador pode então adicionar as políticas do IAM aos perfis e os usuários podem assumir os perfis.

As políticas do IAM definem permissões para uma ação independentemente do método usado para executar a operação. Por exemplo, suponha que você tenha uma política que permite a ação iam: GetRole. Um usuário com essa política pode obter informações de função da AWS Management Console AWS CLI, da ou da AWS API.

Políticas baseadas em identidade

As políticas baseadas em identidade são documentos de políticas de permissões JSON que você pode anexar a uma identidade, como usuário, grupo de usuários ou perfil do IAM. Essas políticas controlam quais ações os usuários e perfis podem realizar, em quais recursos e em que condições. Para saber como criar uma política baseada em identidade, consulte <u>Definir permissões</u> personalizadas do IAM com as políticas gerenciadas pelo cliente no Guia do Usuário do IAM.

As políticas baseadas em identidade podem ser categorizadas como políticas em linha ou políticas gerenciadas. As políticas em linha são anexadas diretamente a um único usuário, grupo ou perfil. As políticas gerenciadas são políticas autônomas que você pode associar a vários usuários, grupos e funções em seu Conta da AWS. As políticas AWS gerenciadas incluem políticas gerenciadas e políticas gerenciadas pelo cliente. Para saber como escolher entre uma política gerenciada ou uma política em linha, consulte Escolher entre políticas gerenciadas e políticas em linha no Guia do usuário do IAM.

Políticas baseadas em recursos

Políticas baseadas em recursos são documentos de políticas JSON que você anexa a um recurso. São exemplos de políticas baseadas em recursos as políticas de confiança de perfil do IAM e as políticas de bucket do Amazon S3. Em serviços compatíveis com políticas baseadas em recursos, os administradores de serviço podem usá-las para controlar o acesso a um recurso específico. Para o atributo ao qual a política está anexada, a política define quais ações uma entidade principal especificado pode executar nesse atributo e em que condições. Você deve especificar uma entidade principal em uma política baseada em recursos. Os diretores podem incluir contas, usuários, funções, usuários federados ou. Serviços da AWS

Políticas baseadas em recursos são políticas em linha localizadas nesse serviço. Você não pode usar políticas AWS gerenciadas do IAM em uma política baseada em recursos.

Listas de controle de acesso (ACLs)

As listas de controle de acesso (ACLs) controlam quais diretores (membros da conta, usuários ou funções) têm permissões para acessar um recurso. ACLs são semelhantes às políticas baseadas em recursos, embora não usem o formato de documento de política JSON.

O Amazon S3 e o AWS WAF Amazon VPC são exemplos de serviços que oferecem suporte. ACLs Para saber mais ACLs, consulte a <u>visão geral da lista de controle de acesso (ACL)</u> no Guia do desenvolvedor do Amazon Simple Storage Service.

Outros tipos de política

AWS oferece suporte a tipos de políticas adicionais menos comuns. Esses tipos de política podem definir o máximo de permissões concedidas a você pelos tipos de política mais comuns.

- Limites de permissões: um limite de permissões é um recurso avançado no qual você define o máximo de permissões que uma política baseada em identidade pode conceder a uma entidade do IAM (usuário ou perfil do IAM). É possível definir um limite de permissões para uma entidade. As permissões resultantes são a interseção das políticas baseadas em identidade de uma entidade com seus limites de permissões. As políticas baseadas em recurso que especificam o usuário ou o perfil no campo Principal não são limitadas pelo limite de permissões. Uma negação explícita em qualquer uma dessas políticas substitui a permissão. Para obter mais informações sobre limites de permissões, consulte Limites de permissões para identidades do IAM no Guia do usuário do IAM.
- Políticas de controle de serviço (SCPs) SCPs são políticas JSON que especificam as permissões máximas para uma organização ou unidade organizacional (OU) em AWS Organizations. AWS Organizations é um serviço para agrupar e gerenciar centralmente vários Contas da AWS que sua empresa possui. Se você habilitar todos os recursos em uma organização, poderá aplicar políticas de controle de serviço (SCPs) a qualquer uma ou a todas as suas contas. O SCP limita as permissões para entidades nas contas dos membros, incluindo cada uma Usuário raiz da conta da AWS. Para obter mais informações sobre Organizations e SCPs, consulte Políticas de controle de serviços no Guia AWS Organizations do Usuário.
- Políticas de controle de recursos (RCPs) RCPs são políticas JSON que você pode usar para definir o máximo de permissões disponíveis para recursos em suas contas sem atualizar as políticas do IAM anexadas a cada recurso que você possui. O RCP limita as permissões para recursos nas contas dos membros e pode afetar as permissões efetivas para identidades, incluindo a Usuário raiz da conta da AWS, independentemente de pertencerem à sua organização. Para obter mais informações sobre Organizations e RCPs, incluindo uma lista Serviços da AWS desse suporte RCPs, consulte Políticas de controle de recursos (RCPs) no Guia AWS Organizations do usuário.
- Políticas de sessão: são políticas avançadas que você transmite como um parâmetro quando cria de forma programática uma sessão temporária para um perfil ou um usuário federado. As permissões da sessão resultante são a interseção das políticas baseadas em identidade do usuário ou do perfil e das políticas de sessão. As permissões também podem ser provenientes de uma política baseada em recursos. Uma negação explícita em qualquer uma dessas políticas

substitui a permissão. Para obter mais informações, consulte <u>Políticas de sessão</u> no Guia do usuário do IAM.

Vários tipos de política

Quando vários tipos de política são aplicáveis a uma solicitação, é mais complicado compreender as permissões resultantes. Para saber como AWS determinar se uma solicitação deve ser permitida quando vários tipos de políticas estão envolvidos, consulte <u>Lógica de avaliação de políticas</u> no Guia do usuário do IAM.

Como o Amazon EMR no EKS funciona com o IAM

Antes de usar o IAM para gerenciar o acesso ao Amazon EMR no EKS, saiba quais recursos do IAM estão disponíveis para uso com o Amazon EMR no EKS.

Recursos do IAM que podem ser usados com o Amazon EMR no EKS

Recurso do IAM	Compatível com o Amazon EMR no EKS
Políticas baseadas em identidade	Sim
Políticas baseadas em recurso	Não
Ações de políticas	Sim
Recursos de políticas	Sim
Chaves de condição de políticas	Sim
ACLs	Não
ABAC (tags em políticas)	Sim
Credenciais temporárias	Sim
Permissões de entidade principal	Sim
Perfis de serviço	Não
Funções vinculadas ao serviço	Sim

Para ter uma visão de alto nível de como o Amazon EMR no EKS e em AWS outros serviços funcionam com a maioria dos recursos do IAM, <u>AWS consulte os serviços que funcionam com</u> o IAM no Guia do usuário do IAM.

Políticas baseadas em identidade para o Amazon EMR no EKS

Compatível com políticas baseadas em identidade: sim

As políticas baseadas em identidade são documentos de políticas de permissões JSON que você pode anexar a uma identidade, como usuário do IAM, grupo de usuários ou perfil. Essas políticas controlam quais ações os usuários e perfis podem realizar, em quais recursos e em que condições. Para saber como criar uma política baseada em identidade, consulte <u>Definir permissões</u> personalizadas do IAM com as políticas gerenciadas pelo cliente no Guia do Usuário do IAM.

Com as políticas baseadas em identidade do IAM, é possível especificar ações e recursos permitidos ou negados, assim como as condições sob as quais as ações são permitidas ou negadas. Você não pode especificar a entidade principal em uma política baseada em identidade porque ela se aplica ao usuário ou perfil ao qual ela está anexada. Para saber mais sobre todos os elementos que podem ser usados em uma política JSON, consulte Referência de elemento de política JSON do IAM no Guia do usuário do IAM.

Exemplos de políticas baseadas em identidade para o Amazon EMR no EKS

Para visualizar exemplos de políticas baseadas em identidade do Amazon EMR no EKS, consulte Exemplos de políticas baseadas em identidade para o Amazon EMR no EKS.

Políticas baseadas em recursos para o Amazon EMR no EKS

Compatibilidade com políticas baseadas em recursos: não

Políticas baseadas em recursos são documentos de políticas JSON que você anexa a um recurso. São exemplos de políticas baseadas em recursos as políticas de confiança de perfil do IAM e as políticas de bucket do Amazon S3. Em serviços compatíveis com políticas baseadas em recursos, os administradores de serviço podem usá-las para controlar o acesso a um recurso específico. Para o atributo ao qual a política está anexada, a política define quais ações uma entidade principal especificado pode executar nesse atributo e em que condições. Você deve especificar uma entidade principal em uma política baseada em recursos. Os diretores podem incluir contas, usuários, funções, usuários federados ou. Serviços da AWS

Para permitir o acesso entre contas, você pode especificar uma conta inteira ou as entidades do IAM em outra conta como a entidade principal em uma política baseada em recursos. Adicionar uma entidade principal entre contas à política baseada em recurso é apenas metade da tarefa de estabelecimento da relação de confiança. Quando o principal e o recurso são diferentes Contas da AWS, um administrador do IAM na conta confiável também deve conceder permissão à entidade principal (usuário ou função) para acessar o recurso. Eles concedem permissão ao anexar uma política baseada em identidade para a entidade. No entanto, se uma política baseada em recurso conceder acesso a uma entidade principal na mesma conta, nenhuma política baseada em identidade adicional será necessária. Consulte mais informações em Acesso a recursos entre contas no IAM no Guia do usuário do IAM.

Ações de políticas para o Amazon EMR no EKS

Compatível com ações de políticas: sim

Os administradores podem usar políticas AWS JSON para especificar quem tem acesso ao quê. Ou seja, qual entidade principal pode executar ações em quais recursos e em que condições.

O elemento Action de uma política JSON descreve as ações que podem ser usadas para permitir ou negar acesso em uma política. As ações de política geralmente têm o mesmo nome da operação de AWS API associada. Existem algumas exceções, como ações somente de permissão, que não têm uma operação de API correspondente. Algumas operações também exigem várias ações em uma política. Essas ações adicionais são chamadas de ações dependentes.

Incluem ações em uma política para conceder permissões para executar a operação associada.

Para visualizar uma lista de ações do Amazon EMR no EKS, consulte <u>Ações, recursos e chaves de</u> condição para o Amazon EMR no EKS na Referência de autorização do serviço.

As ações de política no Amazon EMR no EKS usam o seguinte prefixo antes da ação:

```
emr-containers
```

Para especificar várias ações em uma única declaração, separe-as com vírgulas.

```
"Action": [
    "emr-containers:action1",
    "emr-containers:action2"
]
```

Para visualizar exemplos de políticas baseadas em identidade do Amazon EMR no EKS, consulte Exemplos de políticas baseadas em identidade para o Amazon EMR no EKS.

Recursos de políticas para o Amazon EMR no EKS

Compatível com recursos de políticas: sim

Os administradores podem usar políticas AWS JSON para especificar quem tem acesso ao quê. Ou seja, qual entidade principal pode executar ações em quais recursos e em que condições.

O elemento de política JSON Resource especifica o objeto ou os objetos aos quais a ação se aplica. As instruções devem incluir um elemento Resource ou NotResource. Como prática recomendada, especifique um recurso usando seu <u>nome do recurso da Amazon (ARN)</u>. Isso pode ser feito para ações que oferecem compatibilidade com um tipo de recurso específico, conhecido como permissões em nível de recurso.

Para ações que não oferecem compatibilidade com permissões em nível de recurso, como operações de listagem, use um curinga (*) para indicar que a instrução se aplica a todos os recursos.

```
"Resource": "*"
```

Para ver uma lista dos tipos de recursos do Amazon EMR no EKS e seus ARNs, consulte <u>Recursos definidos pelo Amazon EMR no EKS na Referência</u> de autorização de serviço. Para aprender com quais ações você pode especificar o ARN de cada recurso, consulte <u>Ações, recursos e chaves de condição para o Amazon EMR no EKS.</u>

Para visualizar exemplos de políticas baseadas em identidade do Amazon EMR no EKS, consulte Exemplos de políticas baseadas em identidade para o Amazon EMR no EKS.

Chaves de condição de políticas para o Amazon EMR no EKS

Compatível com chaves de condição de política específicas de serviço: sim

Os administradores podem usar políticas AWS JSON para especificar quem tem acesso ao quê. Ou seja, qual entidade principal pode executar ações em quais recursos e em que condições.

O elemento Condition (ou bloco Condition) permite que você especifique condições nas quais uma instrução estiver em vigor. O elemento Condition é opcional. É possível criar expressões condicionais que usem <u>agentes de condição</u>, como "igual a" ou "menor que", para fazer a condição da política corresponder aos valores na solicitação.

Se você especificar vários elementos de Condition em uma declaração ou várias chaves em um único elemento de Condition, a AWS os avaliará usando uma operação lógica AND. Se você especificar vários valores para uma única chave de condição, AWS avalia a condição usando uma OR operação lógica. Todas as condições devem ser atendidas antes que as permissões da instrução sejam concedidas.

Você também pode usar variáveis de espaço reservado ao especificar condições. Por exemplo, é possível conceder a um usuário do IAM permissão para acessar um recurso somente se ele estiver marcado com seu nome de usuário do IAM. Para obter mais informações, consulte <u>Elementos da política do IAM: variáveis e tags no Guia do usuário do IAM.</u>

AWS suporta chaves de condição globais e chaves de condição específicas do serviço. Para ver todas as chaves de condição AWS globais, consulte as chaves de contexto de condição AWS global no Guia do usuário do IAM.

Para visualizar uma lista de chaves de condição do Amazon EMR no EKS e aprender quais ações e recursos você pode usar uma chave de condição, consulte <u>Ações, recursos e chaves de condição</u> para o Amazon EMR no EKS na Referência de autorização do serviço.

Para visualizar exemplos de políticas baseadas em identidade do Amazon EMR no EKS, consulte Exemplos de políticas baseadas em identidade para o Amazon EMR no EKS.

Listas de controle de acesso (ACLs) no Amazon EMR no EKS

Suportes ACLs: Não

As listas de controle de acesso (ACLs) controlam quais diretores (membros da conta, usuários ou funções) têm permissões para acessar um recurso. ACLs são semelhantes às políticas baseadas em recursos, embora não usem o formato de documento de política JSON.

Controle de acesso por atributo (ABAC) com o Amazon EMR no EKS

Oferece compatibilidade com ABAC (tags em Sim políticas)

O controle de acesso por atributo (ABAC) é uma estratégia de autorização que define as permissões com base em atributos. Em AWS, esses atributos são chamados de tags. Você pode anexar tags a entidades do IAM (usuários ou funções) e a vários AWS recursos. Marcar de entidades e atributos é

a primeira etapa do ABAC. Em seguida, você cria políticas de ABAC para permitir operações quando a tag da entidade principal corresponder à tag do recurso que ela estiver tentando acessar.

O ABAC é útil em ambientes que estão crescendo rapidamente e ajuda em situações em que o gerenciamento de políticas se torna um problema.

Para controlar o acesso baseado em tags, forneça informações sobre as tags no <u>elemento de condição</u> de uma política usando as aws:ResourceTag/key-name, aws:RequestTag/key-name ou chaves de condição aws:TagKeys.

Se um serviço for compatível com as três chaves de condição para cada tipo de recurso, o valor será Sim para o serviço. Se um serviço for compatível com as três chaves de condição somente para alguns tipos de recursos, o valor será Parcial

Para obter mais informações sobre o ABAC, consulte <u>Definir permissões com autorização do ABAC</u> no Guia do usuário do IAM. Para visualizar um tutorial com etapas para configurar o ABAC, consulte <u>Usar controle de acesso baseado em atributos (ABAC)</u> no Guia do usuário do IAM.

Uso de credencial temporária com o Amazon EMR no EKS

Compatível com credenciais temporárias: sim

Alguns Serviços da AWS não funcionam quando você faz login usando credenciais temporárias. Para obter informações adicionais, incluindo quais Serviços da AWS funcionam com credenciais temporárias, consulte Serviços da AWS "Trabalhe com o IAM" no Guia do usuário do IAM.

Você está usando credenciais temporárias se fizer login AWS Management Console usando qualquer método, exceto um nome de usuário e senha. Por exemplo, quando você acessa AWS usando o link de login único (SSO) da sua empresa, esse processo cria automaticamente credenciais temporárias. Você também cria automaticamente credenciais temporárias quando faz login no console como usuário e, em seguida, alterna perfis. Para obter mais informações sobre como alternar funções, consulte Alternar para um perfil do IAM (console) no Guia do usuário do IAM.

Você pode criar manualmente credenciais temporárias usando a AWS API AWS CLI ou. Em seguida, você pode usar essas credenciais temporárias para acessar AWS. AWS recomenda que você gere credenciais temporárias dinamicamente em vez de usar chaves de acesso de longo prazo. Para obter mais informações, consulte Credenciais de segurança temporárias no IAM.

Permissões de entidades principais entre serviços para o Amazon EMR no EKS

Compatibilidade com o recurso de encaminhamento de sessões de acesso (FAS): sim

Quando você usa um usuário ou uma função do IAM para realizar ações AWS, você é considerado principal. Ao usar alguns serviços, você pode executar uma ação que inicia outra ação em um serviço diferente. O FAS usa as permissões do diretor chamando um AWS service (Serviço da AWS), combinadas com a solicitação AWS service (Serviço da AWS) para fazer solicitações aos serviços posteriores. As solicitações do FAS são feitas somente quando um serviço recebe uma solicitação que requer interações com outros Serviços da AWS ou com recursos para ser concluída. Nesse caso, você precisa ter permissões para executar ambas as ações. Para obter detalhes da política ao fazer solicitações de FAS, consulte Sessões de acesso direto.

Perfis de serviço para o Amazon EMR no EKS

Oferece suporte a perfis de serviço	Não
-------------------------------------	-----

Perfis vinculados ao serviço para o Amazon EMR no EKS

Oferece suporte a perfis vinculados ao serviço Sim

Para obter detalhes sobre como criar ou gerenciar perfis vinculados a serviços, consulte <u>Serviços da AWS que funcionam com o IAM</u>. Encontre um serviço na tabela que inclua um Yes na coluna Perfil vinculado ao serviço. Escolha o link Sim para visualizar a documentação do perfil vinculado a serviço desse serviço.

Uso de perfis vinculados ao serviço para o Amazon EMR no EKS

O Amazon EMR no EKS usa funções vinculadas a <u>serviços AWS Identity and Access Management</u> (IAM). Um perfil vinculado ao serviço é um tipo exclusivo de perfil do IAM vinculado diretamente ao Amazon EMR no EKS. As funções vinculadas ao serviço são predefinidas pelo Amazon EMR no EKS e incluem todas as permissões que o serviço exige para chamar outros AWS serviços em seu nome.

Um perfil vinculado ao serviço facilita a configuração do Amazon EMR no EKS porque você não precisa adicionar manualmente as permissões necessárias. O Amazon EMR no EKS define as permissões dos perfis vinculados ao serviço e, a menos que definido de outra forma, somente o Amazon EMR no EKS pode assumir os perfis. As permissões definidas incluem a política de confiança e a política de permissões, que não pode ser anexada a nenhuma outra entidade do IAM.

Um perfil vinculado ao serviço poderá ser excluído somente após excluir seus atributos relacionados. Isso protege os recursos do Amazon EMR no EKS, pois você não pode remover inadvertidamente a permissão de acesso aos recursos.

Para obter informações sobre outros serviços suportados por perfis vinculados a serviços, consulte Serviços da AWS Suportados pelo IAM e procure os serviços que apresentarem Sim na coluna Função Vinculada a Serviço.. Escolha um Sim com um link para visualizar a documentação do perfil vinculado para esse serviço.

Permissões de perfis vinculados ao serviço para o Amazon EMR no EKS

O Amazon EMR no EKS usa a função vinculada ao serviço chamada. AWSServiceRoleForAmazonEMRContainers

O perfil vinculado ao serviço AWSServiceRoleForAmazonEMRContainers confia nos seguintes serviços para aceitar o perfil:

emr-containers.amazonaws.com

A política de permissões de perfil AmazonEMRContainersServiceRolePolicy permite que o Amazon EMR no EKS conclua um conjunto de ações nos recursos especificados, como demonstra a instrução de política apresentada a seguir.



Note

O conteúdo da política gerenciada muda, então a política mostrada aqui pode ser out-of-date. Veja a maior parte da documentação sobre up-to-date políticas na Amazon EMRContainers ServiceRolePolicy no Guia de referência de políticas AWS gerenciadas.

```
{
  "Version": "2012-10-17",
  "Statement": [
      "Effect": "Allow",
      "Action": [
        "eks:DescribeCluster",
        "eks:ListNodeGroups",
```

```
"eks:DescribeNodeGroup",
    "ec2:DescribeRouteTables",
    "ec2:DescribeSubnets",
    "ec2:DescribeSecurityGroups",
    "elasticloadbalancing:DescribeInstanceHealth",
    "elasticloadbalancing:DescribeLoadBalancers",
    "elasticloadbalancing:DescribeTargetGroups",
    "elasticloadbalancing:DescribeTargetHealth",
    "eks:ListPodIdentityAssociations",
    "eks:DescribePodIdentityAssociation"
  ],
  "Resource": [
    11 * 11
  ],
  "Sid": "AllowEKSDescribecluster"
},
  "Effect": "Allow",
  "Action": [
    "acm:ImportCertificate",
    "acm:AddTagsToCertificate"
  ],
  "Resource": [
    11 * 11
  ],
  "Condition": {
    "StringEquals": {
      "aws:RequestTag/emr-container:endpoint:managed-certificate": "true"
    }
  },
  "Sid": "AllowACMImportcertificate"
},
{
  "Effect": "Allow",
  "Action": [
    "acm:DeleteCertificate"
  ],
  "Resource": [
  ],
  "Condition": {
    "StringEquals": {
      "aws:ResourceTag/emr-container:endpoint:managed-certificate": "true"
    }
```

```
},
    "Sid": "AllowACMDeletecertificate"
}
]
```

Você deve configurar permissões para que uma entidade do IAM (por exemplo, um usuário, grupo ou função) crie, edite ou exclua um perfil vinculado a serviço. Para mais informações, consulte Permissões de perfil vinculado ao serviço no Guia do usuário do IAM.

Criação de um perfil vinculado ao serviço para o Amazon EMR no EKS

Não é necessário criar manualmente um perfil vinculado ao serviço. Ao criar um cluster virtual, o Amazon EMR no EKS cria o perfil vinculado ao serviço para você.

Se excluir esse perfil vinculado ao serviço e precisar criá-lo novamente, será possível usar esse mesmo processo para recriar o perfil em sua conta. Ao criar um cluster virtual, o Amazon EMR no EKS cria o perfil vinculado ao serviço para você novamente.

Você também pode usar o console do IAM para criar um perfil vinculado ao serviço com o caso de uso do Amazon EMR no EKS. Na AWS CLI ou na AWS API, crie uma função vinculada ao serviço com o nome do emr-containers.amazonaws.com serviço. Para obter mais informações, consulte Criar uma função vinculada ao serviço no Guia do usuário do IAM. Se você excluir essa função vinculada ao serviço, será possível usar esse mesmo processo para criar a função novamente.

Edição de um perfil vinculado ao serviço para o Amazon EMR no EKS

O Amazon EMR no EKS não permite editar o perfil vinculado ao serviço AWSServiceRoleForAmazonEMRContainers. Depois que criar um perfil vinculado ao serviço, você não poderá alterar o nome do perfil, pois várias entidades podem fazer referência a ele. No entanto, será possível editar a descrição do perfil usando o IAM. Para obter mais informações, consulte Editar uma função vinculada a serviço no Guia do usuário do IAM.

Exclusão de um perfil vinculado ao serviço do Amazon EMR no EKS

Se você não precisar mais usar um recurso ou serviço que requer um perfil vinculado ao serviço, é recomendável excluí-lo. Dessa forma, você não tem uma entidade não utilizada que não seja monitorada ativamente ou mantida. No entanto, você deve limpar os recursos de seu perfil vinculado ao serviço antes de excluí-lo manualmente.



Note

Se o serviço Amazon EMR no EKS estiver usando o perfil quando você tentar excluir os recursos, a exclusão poderá falhar. Se isso acontecer, espere alguns minutos e tente a operação novamente.

Excluir recursos do Amazon EMR no EKS usados pelos

AWSServiceRoleForAmazonEMRContainers

- Abra o console do Amazon EMR.
- 2. Escolha um cluster virtual.
- 3. Na página Virtual Cluster, escolha Excluir.
- 4. Repita este procedimento para quaisquer outros clusters virtuais em sua conta.

Como excluir manualmente o perfil vinculado ao serviço usando o IAM

Use o console do IAM AWS CLI, o ou a AWS API para excluir a função AWSServiceRoleForAmazonEMRContainers vinculada ao serviço. Para obter mais informações, consulte Excluir um perfil vinculado ao serviço no Guia do usuário do IAM.

Regiões com suporte para perfis vinculados ao serviço do Amazon EMR no EKS

O Amazon EMR no EKS oferece suporte ao uso de perfis vinculados ao serviço em todas as regiões nas quais o serviço está disponível. Para obter mais informações, consulte Cotas e endpoints de serviço do Amazon EMR no EKS.

Políticas gerenciadas para o Amazon EMR no EKS

Veja detalhes sobre as atualizações das políticas AWS gerenciadas do Amazon EMR no EKS desde 1º de março de 2021.

Alteração	Descrição	Data
AmazonEMRContainer sServiceRolePolicy - Foram adicionadas permissões de leitura	As seguintes permissões são adicionad as à política:eks:ListPodIdentit yAssociations ,eks:Descr ibePodIdentityAssociation .	3 de fevereiro de 2023

Alteração	Descrição	Data
para listar associações de identidade de pod EKS em um cluster e outra permissão de leitura para retornar informações descritivas sobre associaçõ es de identidade de pod em um cluster. Para obter mais informações, consulte Amazon EMRContainers ServiceRolePolicy.		
AmazonEMRContainer sServiceRolePolicy : permissões adicionad as para descrever e listar grupos de nós do Amazon EKS, descrever grupos de destino do balancead or de carga e descrever a integridade do destino do balanceador de carga.	As seguintes permissões são adicionad as à política: eks:ListNodeGroups , eks:DescribeNodeGroup , elasticloadbalancing:DescribeTargetGroups e elasticlo adbalancing:DescribeTargetH ealth .	13 de março de 2023
AmazonEMRContainer sServiceRolePolicy - Permissões adicionad as para importar e excluir certificados em AWS Certificate Manager.	As seguintes permissões são adicionad as à política: acm:ImportCertific ate ,acm:AddTagsToCertificate eacm:DeleteCertificate .	3 de dezembro de 2021
O Amazon EMR no EKS começou a rastrear alterações	O Amazon EMR no EKS começou a monitorar as mudanças em suas políticas AWS gerenciadas.	1.º de março de 2021

Uso de perfis de execução de trabalho com o Amazon EMR no EKS

Para usar o comando StartJobRun para enviar uma execução de trabalho em um cluster do EKS, primeiro é necessário integrar um perfil de execução de trabalho para ser usado com um cluster virtual. Para obter mais informações, consulte <u>Criação de um perfil de execução de trabalho</u> em <u>Configuração do Amazon EMR no EKS</u>. Você também pode seguir as instruções na seção <u>Create</u> IAM Role for job execution do Amazon EMR on EKS Workshop.

As permissões apresentadas a seguir devem ser inclusas na política de confiança para o perfil de execução de trabalho.

JSON

```
"Version": "2012-10-17",
  "Statement": [
    {
      "Effect": "Allow",
      "Action": [
        "sts:AssumeRoleWithWebIdentity"
      ],
      "Resource": "*",
      "Condition": {
        "StringLike": {
          "aws:userid": "system:serviceaccount:NAMESPACE:emr-containers-sa-*-
*-AWS_ACCOUNT_ID-BASE36_ENCODED_ROLE_NAME"
        }
      },
      "Sid": "AllowSTSAssumerolewithwebidentity"
  ]
}
```

A política de confiança no exemplo anterior concede permissões somente para uma conta de serviço do Kubernetes gerenciada pelo Amazon EMR com um nome que corresponda ao padrão emr-containers-sa-*-*-*AWS_ACCOUNT_ID-BASE36_ENCODED_ROLE_NAME*. As contas de serviço com esse padrão serão criadas automaticamente no envio do trabalho e terão como escopo o namespace no qual você envia o trabalho. Esta política de confiança permite que estas contas de serviço assumam o perfil de execução e obtenham as credenciais temporárias para o perfil de

execução. As contas de serviço de um cluster diferente do Amazon EKS ou de um namespace diferente no mesmo cluster do EKS estão impedidas de assumir o perfil de execução.

Você pode executar o comando apresentado a seguir para atualizar automaticamente a política de confiança no formato fornecido acima.

Controle do acesso ao perfil de execução

Um administrador do cluster do Amazon EKS pode criar um cluster virtual multilocatário do Amazon EMR no EKS ao qual um administrador do IAM pode adicionar diversos perfis de execução. Como os locatários não confiáveis podem usar esses perfis de execução para enviar trabalhos que executam códigos arbitrários, você pode desejar restringir esses locatários para que eles não possam executar códigos que concedem as permissões atribuídas a um ou mais desses perfis de execução. Para restringir a política do IAM anexada a uma identidade do IAM, o administrador do IAM pode usar a chave de condição opcional emr-containers: ExecutionRoleArn do nome do recurso da Amazon (ARN). Essa condição aceita uma lista de funções de execução ARNs que têm permissões para o cluster virtual, conforme demonstra a política de permissões a seguir.

```
]
        }
      },
      "Sid": "AllowEMRCONTAINERSStartjobrun"
  ]
}
```

Se desejar permitir todos os perfis de execução que começam com um prefixo específico, como MyRole, você poderá substituir o operador de condição ArnEquals pelo operador ArnLike, e poderá substituir o valor de execution_role_arn na condição por um caractere curinga *. Por exemplo, .arn:aws:iam::AWS_ACCOUNT_ID:role/MyRole* Todas as outras chaves de condição de ARN também são compatíveis.



Note

Com o Amazon EMR no EKS, não é possível conceder permissões para perfis de execução com base em etiquetas ou atributos. O Amazon EMR no EKS não oferece suporte ao controle de acesso por etiqueta (TBAC) ou ao controle de acesso por atributo (ABAC) para perfis de execução.

Exemplos de políticas baseadas em identidade para o Amazon EMR no **EKS**

Por padrão, os usuários e os perfis não têm permissão para criar ou modificar os recursos do Amazon EMR no EKS. Eles também não podem realizar tarefas usando a AWS API AWS Management Console, AWS Command Line Interface (AWS CLI) ou. Para conceder permissão aos usuários para executar ações nos recursos que eles precisam, um administrador do IAM pode criar políticas do IAM. O administrador pode então adicionar as políticas do IAM aos perfis e os usuários podem assumir os perfis.

Para aprender a criar uma política baseada em identidade do IAM ao usar esses documentos de política em JSON de exemplo, consulte Criar políticas do IAM (console) no Guia do usuário do IAM.

Para obter detalhes sobre ações e tipos de recursos definidos pelo Amazon EMR no EKS, incluindo o formato de cada um dos tipos de recursos, consulte Ações, recursos e chaves de condição ARNs para o Amazon EMR no EKS na Referência de autorização de serviço.

Tópicos

- Práticas recomendadas de política
- Uso do console do Amazon EMR no EKS
- Permitir que os usuários visualizem suas próprias permissões

Práticas recomendadas de política

As políticas baseadas em identidade determinam se alguém pode criar, acessar ou excluir recursos do Amazon EMR no EKS em sua conta. Essas ações podem incorrer em custos para sua Conta da AWS. Ao criar ou editar políticas baseadas em identidade, siga estas diretrizes e recomendações:

- Comece com as políticas AWS gerenciadas e avance para as permissões de privilégios mínimos

 Para começar a conceder permissões aos seus usuários e cargas de trabalho, use as políticas
 AWS gerenciadas que concedem permissões para muitos casos de uso comuns. Eles estão disponíveis no seu Conta da AWS. Recomendamos que você reduza ainda mais as permissões definindo políticas gerenciadas pelo AWS cliente que sejam específicas para seus casos de uso.

 Para obter mais informações, consulte Políticas gerenciadas pela AWS ou Políticas gerenciadas pela AWS para funções de trabalho no Guia do usuário do IAM.
- Aplique permissões de privilégio mínimo: ao definir permissões com as políticas do IAM, conceda apenas as permissões necessárias para executar uma tarefa. Você faz isso definindo as ações que podem ser executadas em recursos específicos sob condições específicas, também conhecidas como permissões de privilégio mínimo. Para obter mais informações sobre como usar o IAM para aplicar permissões, consulte Políticas e permissões no IAM no Guia do usuário do IAM.
- Use condições nas políticas do IAM para restringir ainda mais o acesso: você pode adicionar uma condição às políticas para limitar o acesso a ações e recursos. Por exemplo, você pode escrever uma condição de política para especificar que todas as solicitações devem ser enviadas usando SSL. Você também pode usar condições para conceder acesso às ações de serviço se elas forem usadas por meio de uma ação específica AWS service (Serviço da AWS), como AWS CloudFormation. Para obter mais informações, consulte Elementos da política JSON do IAM: condição no Guia do usuário do IAM.
- Use o IAM Access Analyzer para validar suas políticas do IAM a fim de garantir permissões seguras e funcionais: o IAM Access Analyzer valida as políticas novas e existentes para que elas sigam a linguagem de política do IAM (JSON) e as práticas recomendadas do IAM. O IAM Access Analyzer oferece mais de cem verificações de política e recomendações práticas para ajudar a criar políticas seguras e funcionais. Para obter mais informações, consulte <u>Validação de políticas</u> do IAM Access Analyzer no Guia do Usuário do IAM.

 Exigir autenticação multifator (MFA) — Se você tiver um cenário que exija usuários do IAM ou um usuário root, ative Conta da AWS a MFA para obter segurança adicional. Para exigir MFA quando as operações de API forem chamadas, adicione condições de MFA às suas políticas. Para obter mais informações, consulte <u>Configuração de acesso à API protegido por MFA</u> no Guia do Usuário do IAM.

Para obter mais informações sobre as práticas recomendadas do IAM, consulte <u>Práticas</u> recomendadas de segurança no IAM no Guia do usuário do IAM.

Uso do console do Amazon EMR no EKS

Para acessar o console do Amazon EMR no EKS, você deve ter um conjunto mínimo de permissões. Essas permissões devem permitir que você liste e visualize detalhes sobre recursos do Amazon EMR no EKS em sua Conta da AWS. Caso crie uma política baseada em identidade mais restritiva que as permissões mínimas necessárias, o console não funcionará como pretendido para entidades (usuários ou perfis) com essa política.

Você não precisa permitir permissões mínimas do console para usuários que estão fazendo chamadas somente para a API AWS CLI ou para a AWS API. Em vez disso, permita o acesso somente a ações que correspondam à operação de API que estiverem tentando executar.

Para garantir que usuários e funções ainda possam usar o Amazon EMR no console EKS, anexe também o Amazon EMR no EKS ConsoleAccess ou a política ReadOnly AWS gerenciada às entidades. Para obter informações, consulte <u>Adicionar permissões a um usuário</u> no Guia do usuário do IAM.

Permitir que os usuários visualizem suas próprias permissões

Este exemplo mostra como criar uma política que permita que os usuários do IAM visualizem as políticas gerenciadas e em linha anexadas a sua identidade de usuário. Essa política inclui permissões para concluir essa ação no console ou programaticamente usando a API AWS CLI ou AWS.

```
"Action": [
                "iam:GetUserPolicy",
                "iam:ListGroupsForUser",
                "iam:ListAttachedUserPolicies",
                "iam:ListUserPolicies",
                "iam:GetUser"
            ],
            "Resource": ["arn:aws:iam::*:user/${aws:username}"]
        },
        {
            "Sid": "NavigateInConsole",
            "Effect": "Allow",
            "Action": [
                "iam:GetGroupPolicy",
                "iam:GetPolicyVersion",
                "iam:GetPolicy",
                "iam:ListAttachedGroupPolicies",
                "iam:ListGroupPolicies",
                "iam:ListPolicyVersions",
                "iam:ListPolicies",
                "iam:ListUsers"
            ],
            "Resource": "*"
        }
    ]
}
```

Políticas para controle de acesso baseado em etiquetas

É possível aplicar condições em sua política baseada em identidade para controlar o acesso a clusters virtuais e execuções de trabalhos com base em etiquetas. Para obter mais informações sobre marcação, consulte Marcação de recursos do Amazon EMR no EKS.

Os exemplos a seguir demonstram diferentes cenários e maneiras de usar operadores de condição com chaves de condição do Amazon EMR no EKS. Estas instruções de política do IAM são destinadas somente para fins de demonstração e não devem ser usadas em ambientes de produção. Há várias maneiras de combinar declarações de políticas para conceder e negar permissões de acordo com seus requisitos. Para obter mais informações sobre como planejar e testar políticas do IAM, consulte o Guia do usuário do IAM.



M Important

Recusar, explicitamente, permissões para ações de uso de tags é uma consideração importante. Isso evita que os usuários façam a marcação de um recurso e, assim, concedam a si mesmos permissões que você não pretendia conceder. Se as ações de marcação de um recurso não forem negadas, um usuário poderá modificar as etiquetas e driblar a intenção das políticas baseadas em etiquetas. Para obter um exemplo de política que nega ações de marcação, consulte Negação de acesso para adição ou remoção de etiquetas.

Os exemplos abaixo demonstram políticas de permissões baseadas em identidade usadas para controlar as ações que são permitidas com clusters virtuais do Amazon EMR no EKS.

Ações permitidas somente em recursos com valores de etiquetas específicos

No exemplo de política a seguir, o operador de StringEquals condição tenta combinar dev com o valor do departamento de tag. Se a etiqueta "Department" não tiver sido adicionada ao cluster virtual ou não contiver o valor "dev", a política não se aplicará e as ações não serão permitidas por esta política. Se nenhuma outra instrução de política permitir as ações, o usuário poderá trabalhar somente com clusters virtuais que tenham essa etiqueta com este valor.

```
"Version": "2012-10-17",
"Statement": [
  {
    "Effect": "Allow",
    "Action": [
      "emr-containers:DescribeVirtualCluster"
    ],
    "Resource": [
      11 * 11
    ],
    "Condition": {
      "StringEquals": {
        "aws:ResourceTag/department": "dev"
      }
    },
    "Sid": "AllowEMRCONTAINERDescribevirtualcluster"
```

```
}
]
}
```

Você também pode especificar vários valores de tag usando um operador de condição. Por exemplo, para permitir ações em clusters virtuais nos quais a etiqueta department contém o valor dev ou test, você pode substituir o bloco condicional no exemplo anterior pelo apresentado a seguir.

```
"Condition": {
     "StringEquals": {
         "aws:ResourceTag/department": ["dev", "test"]
      }
}
```

Marcação obrigatória na criação de um recurso

No exemplo abaixo, a etiqueta precisa ser aplicada ao criar o cluster virtual.

```
"Version": "2012-10-17",
  "Statement": [
      "Effect": "Allow",
      "Action": [
        "emr-containers:CreateVirtualCluster"
      ],
      "Resource": [
        11 * 11
      ],
      "Condition": {
        "StringEquals": {
          "aws:RequestTag/department": "dev"
        }
      },
      "Sid": "AllowEMRCONTAINERSCreatevirtualcluster"
    }
  ]
}
```

A instrução de política apresentada a seguir permite que um usuário crie um cluster virtual somente se o cluster tiver uma etiqueta department, que pode conter qualquer valor.

JSON

```
"Version": "2012-10-17",
  "Statement": [
      "Effect": "Allow",
      "Action": [
        "emr-containers:CreateVirtualCluster"
      ],
      "Resource": [
      ],
      "Condition": {
        "Null": {
          "aws:RequestTag/department": "false"
        }
      },
      "Sid": "AllowEMRCONTAINERSCreatevirtualcluster"
  ]
}
```

Negação de acesso para adição ou remoção de etiquetas

O efeito desta política é negar a um usuário a permissão para adicionar ou remover quaisquer etiquetas em clusters virtuais que estejam marcados com uma etiqueta department que contenha o valor dev.

```
{
    "Version": "2012-10-17",
    "Statement": [
        {
            "Effect": "Deny",
            "Action": [
```

Solução de problemas de identidade e acesso do Amazon EMR no EKS

Use as informações apresentadas a seguir para ajudar a diagnosticar e corrigir problemas comuns que você possa encontrar ao trabalhar com o Amazon EMR no EKS e com o IAM.

Tópicos

- Não tenho autorização para executar uma ação no Amazon EMR no EKS
- · Não estou autorizado a realizar iam: PassRole
- Quero permitir que pessoas fora da minha AWS conta acessem meus recursos do Amazon EMR no EKS

Não tenho autorização para executar uma ação no Amazon EMR no EKS

Se isso AWS Management Console indicar que você não está autorizado a realizar uma ação, entre em contato com o administrador para obter ajuda. O administrador é a pessoa que forneceu o seu nome de usuário e senha.

O erro do exemplo a seguir ocorre quando o usuário mateojackson tenta usar o console para visualizar detalhes sobre um recurso do my-example-widget fictício, mas não tem as permissões fictícias do emr-containers: GetWidget.

Solução de problemas 366

```
User: arn:aws:iam::123456789012:user/mateojackson is not authorized to perform: emr-containers:GetWidget on resource: my-example-widget
```

Neste caso, Mateo pede ao administrador para atualizar suas políticas e permitir o acesso ao recurso my-example-widget usando a ação emr-containers: GetWidget.

Não estou autorizado a realizar iam: PassRole

Se você receber uma mensagem de erro informando que não tem autorização para executar a ação iam: PassRole, suas políticas devem ser atualizadas para permitir a transmissão de um perfil ao Amazon EMR no EKS.

Alguns Serviços da AWS permitem que você passe uma função existente para esse serviço em vez de criar uma nova função de serviço ou uma função vinculada ao serviço. Para fazer isso, é preciso ter permissões para passar o perfil para o serviço.

O erro exemplificado a seguir ocorre quando uma usuária do IAM chamada marymajor tenta usar o console para executar uma ação no Amazon EMR no EKS. No entanto, a ação exige que o serviço tenha permissões concedidas por um perfil de serviço. Mary não tem permissões para passar o perfil para o serviço.

```
User: arn:aws:iam::123456789012:user/marymajor is not authorized to perform:
   iam:PassRole
```

Nesse caso, as políticas de Mary devem ser atualizadas para permitir que ela realize a ação iam: PassRole.

Se precisar de ajuda, entre em contato com seu AWS administrador. Seu administrador é a pessoa que forneceu suas credenciais de login.

Quero permitir que pessoas fora da minha AWS conta acessem meus recursos do Amazon EMR no EKS

Você pode criar um perfil que os usuários de outras contas ou pessoas fora da organização podem usar para acessar seus recursos. É possível especificar quem é confiável para assumir o perfil. Para serviços que oferecem suporte a políticas baseadas em recursos ou listas de controle de acesso (ACLs), você pode usar essas políticas para conceder às pessoas acesso aos seus recursos.

Para saber mais, consulte:

Solução de problemas 367

- Para saber se o Amazon EMR no EKS oferece suporte a esses recursos, consulte Como o Amazon EMR no EKS funciona com o IAM.
- Para saber como fornecer acesso aos seus recursos em todos os Contas da AWS que você
 possui, consulte Como fornecer acesso a um usuário do IAM em outro Conta da AWS que você
 possui no Guia do usuário do IAM.
- Para saber como fornecer acesso aos seus recursos a terceiros Contas da AWS, consulte Como fornecer acesso Contas da AWS a terceiros no Guia do usuário do IAM.
- Para saber como conceder acesso por meio da federação de identidades, consulte <u>Conceder</u>
 <u>acesso a usuários autenticados externamente</u> (federação de identidades) no Guia do usuário do
 IAM.
- Para saber a diferença entre perfis e políticas baseadas em recurso para acesso entre contas, consulte Acesso a recursos entre contas no IAM no Guia do usuário do IAM.

Usando o Amazon EMR no EKS com AWS Lake Formation para controle de acesso refinado

Com o Amazon EMR versão 7.7 e superior, você pode aproveitar o AWS Lake Formation para aplicar controles de acesso refinados nas tabelas do AWS Glue Data Catalog que são apoiadas por buckets do Amazon S3. Esse recurso permite que você configure controles de acesso em nível de tabela, linha, coluna e célula para consultas de leitura em seu Amazon EMR no EKS Spark Jobs.

Tópicos

- Como o Amazon EMR no EKS funciona com AWS o Lake Formation
- · Habilite o Lake Formation com o Amazon EMR no EKS
- Considerações e limitações
- Solução de problemas

Como o Amazon EMR no EKS funciona com AWS o Lake Formation

Usar o Amazon EMR no EKS com o Lake Formation permite aplicar uma camada de permissões em cada Spark Job para aplicar o controle de permissão do Lake Formation quando o Amazon EMR no EKS executa trabalhos. O Amazon EMR no EKS usa perfis de recursos do Spark para criar dois perfis para executar tarefas com eficiência. O perfil do usuário executa o código fornecido pelo usuário, enquanto o perfil do sistema aplica as políticas do Lake Formation. Cada Job habilitado para

Lake Formation utiliza dois drivers Spark, um para o perfil do usuário e outro para o perfil do sistema. Para obter mais informações, consulte O que é AWS Lake Formation.

A seguir, uma visão geral de alto nível de como o Amazon EMR no EKS obtém acesso aos dados protegidos pelas políticas de segurança do Lake Formation.

As etapas a seguir descrevem esse processo:

- 1. Um usuário envia um Spark Job para um Amazon EMR habilitado para AWS Lake Formation no cluster virtual EKS.
- 2. O serviço Amazon EMR on EKS configura o driver do usuário e executa o trabalho no perfil do usuário. O User Driver executa uma versão simplificada do Spark que não tem a capacidade de executar tarefas, solicitar executores, acessar o Amazon S3 ou o Glue Data Catalog. Ele apenas cria um plano de Job.
- 3. O serviço Amazon EMR no EKS configura um segundo driver chamado System Driver e o executa no Perfil do Sistema (com uma identidade privilegiada). O Amazon EKS configura um canal TLS criptografado entre os dois drivers para comunicação. O driver do usuário usa o canal para enviar os planos de trabalho para o driver do sistema. O driver do sistema não executa o código enviado pelo usuário. Ele executa o Spark completo e se comunica com o Amazon S3 e o catálogo de dados para acesso aos dados. Ele solicita executores e compila o Job Plan em uma sequência de estágios de execução.
- 4. O Amazon EMR no serviço EKS então executa as etapas nos executores. O código do usuário em qualquer estágio é executado exclusivamente nos executores do perfil do usuário.
- 5. Os estágios que lêem dados das tabelas do Catálogo de Dados protegidas pelo Lake Formation ou aqueles que aplicam filtros de segurança são delegados aos executores do sistema.

Habilite o Lake Formation com o Amazon EMR no EKS

Com o Amazon EMR versão 7.7 e superior, você pode aproveitar o AWS Lake Formation para aplicar controles de acesso refinados nas tabelas do catálogo de dados que são apoiadas pelo Amazon S3. Esse recurso permite configurar controles de acesso em nível de tabela, linha, coluna e célula para consultas de leitura em seu Amazon EMR no EKS Spark Jobs.

Esta seção aborda como criar uma configuração de segurança e configurar o Lake Formation para trabalhar com o Amazon EMR. Também descreve como criar um cluster virtual com a Configuração

de Segurança que você criou para o Lake Formation. Essas seções devem ser concluídas em sequência.

Etapa 1: configurar permissões em nível de coluna, linha ou célula com base no Lake Formation

Primeiro, para aplicar permissões em nível de linha e coluna com o Lake Formation, o administrador do data lake do Lake Formation deve definir a tag de LakeFormationAuthorizedCallersessão. O Lake Formation usa essa etiqueta de sessão para autorizar os chamadores e fornecer acesso ao data lake.

Navegue até o console do AWS Lake Formation e selecione a opção Configurações de integração de aplicativos na seção Administração na barra lateral. Em seguida, marque a caixa Permitir que mecanismos externos filtrem dados em locais do Amazon S3 registrados no Lake Formation. Adicione a AWS conta em IDs que os Spark Jobs seriam executados e os valores da tag Session.

Observe que a tag de LakeFormationAuthorizedCallersessão passada aqui é passada SecurityConfigurationposteriormente quando você configura as funções do IAM, na seção 3.

Etapa 2: Configurar as permissões do EKS RBAC

Em segundo lugar, você configura permissões para controle de acesso baseado em funções.

Forneça permissões de cluster EKS para o Amazon EMR no serviço EKS

O Amazon EMR no EKS Service deve ter permissões de função de cluster do EKS para que possa criar permissões entre namespaces para que o driver do sistema separe os executores do usuário no namespace do usuário.

Criar função de cluster

Esse exemplo define permissões para uma coleção de recursos.

```
vim emr-containers-cluster-role.yaml
---
apiVersion: rbac.authorization.k8s.io/v1
kind: ClusterRole
metadata:
   name: emr-containers
rules:
```

```
- apiGroups: [""]
    resources: ["namespaces"]
    verbs: ["get"]
  - apiGroups: [""]
    resources: ["serviceaccounts", "services", "configmaps", "events", "pods", "pods/
log"]
    verbs: ["get", "list", "watch", "describe", "create", "edit", "delete",
 "deletecollection", "annotate", "patch", "label"]
  - apiGroups: [""]
    resources: ["secrets"]
    verbs: ["create", "patch", "delete", "watch"]
  - apiGroups: ["apps"]
    resources: ["statefulsets", "deployments"]
    verbs: ["get", "list", "watch", "describe", "create", "edit", "delete", "annotate",
 "patch", "label"]
  - apiGroups: ["batch"]
    resources: ["jobs"]
    verbs: ["get", "list", "watch", "describe", "create", "edit", "delete", "annotate",
 "patch", "label"]
  - apiGroups: ["extensions", "networking.k8s.io"]
    resources: ["ingresses"]
    verbs: ["get", "list", "watch", "describe", "create", "edit", "delete", "annotate",
 "patch", "label"]
  - apiGroups: ["rbac.authorization.k8s.io"]
    resources: ["clusterroles", "clusterrolebindings", "roles", "rolebindings"]
    verbs: ["get", "list", "watch", "describe", "create", "edit", "delete",
 "deletecollection", "annotate", "patch", "label"]
  - apiGroups: [""]
    resources: ["persistentvolumeclaims"]
    verbs: ["get", "list", "watch", "describe", "create", "edit", "delete",
 "deletecollection", "annotate", "patch", "label"]
```

```
kubectl apply -f emr-containers-cluster-role.yaml
```

Criar vinculações de funções de cluster

```
vim emr-containers-cluster-role-binding.yaml
---
apiVersion: rbac.authorization.k8s.io/v1
kind: ClusterRoleBinding
metadata:
   name: emr-containers
```

```
subjects:
- kind: User
  name: emr-containers
  apiGroup: rbac.authorization.k8s.io
roleRef:
  kind: ClusterRole
  name: emr-containers
  apiGroup: rbac.authorization.k8s.io
---
```

```
kubectl apply -f emr-containers-cluster-role-binding.yaml
```

Forneça acesso ao namespace ao Amazon EMR no serviço EKS

Crie dois namespaces do Kubernetes, um para driver e executores do usuário e outro para drivers e executores do sistema, e habilite o acesso ao serviço Amazon EMR on EKS para enviar trabalhos nos namespaces do usuário e do sistema. Siga o guia existente para fornecer acesso a cada namespace, que está disponível em Habilitar o acesso ao cluster usando. aws-auth

Etapa 3: Configurar funções do IAM para componentes de perfil de usuário e sistema

Em terceiro lugar, você configura funções para componentes específicos. Um Spark Job habilitado para Lake Formation tem dois componentes, Usuário e Sistema. O driver e os executores do usuário são executados no namespace do usuário e estão vinculados ao JobExecutionRole que é passado na API. StartJobRun O driver e os executores do sistema são executados no namespace System e estão vinculados à função. QueryEngine

Configurar a função do Query Engine

A QueryEngine função está vinculada aos componentes do System Space e teria permissões para assumir a tag JobExecutionRolewith LakeFormationAuthorizedCallerSession. A política de permissões do IAM do papel do Query Engine é a seguinte:

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Sid": "AssumeJobRoleWithSessionTagAccessForSystemDriver",
```

```
"Effect": "Allow",
      "Action": [
        "sts:AssumeRole",
        "sts:TagSession"
      ],
      "Resource": [
        "arn:aws:iam::*:role/JobExecutionRole"
      ],
      "Condition": {
        "StringLike": {
          "aws:RequestTag/LakeFormationAuthorizedCaller": "EMR on EKS Engine"
        }
      }
    },
      "Sid": "AssumeJobRoleWithSessionTagAccessForSystemExecutor",
      "Effect": "Allow",
      "Action": [
        "sts:AssumeRole"
      ],
      "Resource": [
        "arn:aws:iam::*:role/JobExecutionRole"
      1
    },
      "Sid": "CreateCertificateAccessForTLS",
      "Effect": "Allow",
      "Action": [
        "emr-containers:CreateCertificate"
      ],
      "Resource": [
      ]
    }
  ]
}
```

Configure a política de confiança do papel do Query Engine para confiar no namespace do sistema Kubernetes.

```
aws emr-containers update-role-trust-policy \
```

```
--cluster-name eks cluster \
--namespace eks system namespace \
--role-name query_engine_iam_role_name
```

Para obter mais informações, consulte <u>Atualização da política de confiança da função</u>.

Configurar a função de execução de trabalhos

As permissões do Lake Formation controlam o acesso aos recursos do AWS Glue Data Catalog, aos locais do Amazon S3 e aos dados subjacentes nesses locais. As permissões do IAM controlam o acesso ao Lake Formation, ao AWS Glue APIs e aos recursos. Embora você possa ter a permissão do Lake Formation para acessar uma tabela no Catálogo de Dados (SELECT), sua operação falhará se você não tiver a permissão do IAM nas operações da glue: Get* API.

Política de permissões do IAM de JobExecutionRole: A JobExecutionfunção deve ter as declarações de política em sua política de permissões.

```
"Version": "2012-10-17",
"Statement": [
  {
    "Sid": "GlueCatalogAccess",
    "Effect": "Allow",
    "Action": [
      "glue:Get*",
      "glue:Create*",
      "glue:Update*"
    ],
    "Resource": [
      11 * 11
    ]
  },
    "Sid": "LakeFormationAccess",
    "Effect": "Allow",
    "Action": [
      "lakeformation:GetDataAccess"
    ],
    "Resource": [
      11 * 11
```

```
}

},

{
    "Sid": "CreateCertificateAccessForTLS",
    "Effect": "Allow",
    "Action": [
        "emr-containers:CreateCertificate"
],
    "Resource": [
        "*"
]
}
```

Política de confiança do IAM para JobExecutionRole:

```
"Version": "2012-10-17",
"Statement": [
 {
   "Sid": "TrustQueryEngineRoleForSystemDriver",
   "Effect": "Allow",
   "Action": [
      "sts:AssumeRole",
     "sts:TagSession"
   ],
   "Resource": [
      "arn:aws:iam::*:role/QueryExecutionRole"
   ],
   "Condition": {
      "StringLike": {
        "aws:RequestTag/LakeFormationAuthorizedCaller": "EMR on EKS Engine"
     }
   }
 },
   "Sid": "TrustQueryEngineRoleForSystemExecutor",
   "Effect": "Allow",
   "Action": [
```

```
"sts:AssumeRole"
],
    "Resource": [
        "arn:aws:iam::*:role/QueryEngineRole"
]
}
]
}
```

Configure a política de confiança da função de execução de Job para confiar no namespace de usuário do Kubernetes:

```
aws emr-containers update-role-trust-policy \
    --cluster-name eks cluster \
    --namespace eks User namespace \
    --role-name job_execution_role_name
```

Para obter mais informações, consulte <u>Atualizar a política de confiança da função de execução do</u> trabalho.

Etapa 4: Configurar a configuração de segurança

Para executar um trabalho habilitado para Lake Formation, você deve criar uma configuração de segurança.

Certifique-se de que a etiqueta de sessão passada no campo authorizedSessionTagValue possa autorizar o Lake Formation. Defina o valor como aquele configurado em Lake Formation, em<u>Etapa 1:</u> configurar permissões em nível de coluna, linha ou célula com base no Lake Formation.

Etapa 5: criar um cluster virtual

Crie um cluster virtual do Amazon EMR no EKS com uma configuração de segurança.

```
aws emr-containers create-virtual-cluster \
--name my-lf-enabled-vc \
--container-provider '{
    "id": "eks-cluster",
    "type": "EKS",
    "info": {
        "eksInfo": {
            "namespace": "user-namespace"
        }
    }
}' \
--security-configuration-id SecurityConfiguraionId
```

Certifique-se de que o SecurityConfigurationID da etapa anterior seja passado, para que a configuração de autorização do Lake Formation seja aplicada a todos os trabalhos em execução no cluster virtual. Para obter mais informações, consulte Registrar o cluster Amazon EKS no Amazon EMR.

Etapa 6: Enviar um Job no FGAC Enabled VirtualCluster

O processo de envio de trabalhos é o mesmo para trabalhos que não sejam da Lake Formation e da Lake Formation. Para obter mais informações, consulte Enviar um trabalho executado com StartJobRun.

O driver Spark, o executor e os registros de eventos do driver do sistema são armazenados no bucket S3 da conta de AWS serviço para depuração. Recomendamos configurar uma chave KMS gerenciada pelo cliente no Job Run para criptografar todos os registros armazenados no service bucket. AWS Para obter mais informações sobre como habilitar a criptografia de registros, consulte Criptografar o Amazon EMR em registros do EKS.

Considerações e limitações

Observe as seguintes considerações e limitações ao usar o Lake Formation com o Amazon EMR no EKS:

- O Amazon EMR no EKS oferece suporte ao controle de acesso refinado via Lake Formation somente para os formatos de tabela Apache Hive, Apache Iceberg, Apache Hudi e Delta. Os formatos do Apache Hive incluem Parquet, ORC e xSV.
- DynamicResourceAllocationestá ativado por padrão e você não pode desativar
 DynamicResourceAllocation os trabalhos do Lake Formation. Como o valor padrão da spark.dynamicAllocation.maxExecutors configuração do DRA é infinito, configure um valor apropriado com base na sua carga de trabalho.
- As tarefas habilitadas para Lake Formation n\u00e3o oferecem suporte ao uso de EMR personalizado em imagens EKS no driver do sistema e nos executores do sistema.
- Você só pode usar o Lake Formation com trabalhos do Spark.
- O EMR no EKS com Lake Formation suporta apenas uma única sessão do Spark durante todo o trabalho.
- O EMR no EKS com Lake Formation só oferece suporte a consultas de tabelas entre contas compartilhadas por meio de links de recursos.
- As seguintes opções não são compatíveis:
 - Conjuntos de dados distribuídos resilientes (RDD)
 - Streaming do Spark
 - Gravação com as permissões concedidas pelo Lake Formation
 - Controle de acesso para colunas aninhadas
- O EMR no EKS bloqueia funcionalidades que podem prejudicar o isolamento completo do driver do sistema, incluindo as seguintes:
 - UDTs, Hive UDFs e qualquer função definida pelo usuário que envolva classes personalizadas
 - Fontes de dados personalizadas
 - Fornecimento de frascos adicionais para extensão, conector ou comando de metastore do Spark ANALYZE TABLE
- Para impor controles de acesso, EXPLAIN PLAN e operações de DDL, como DESCRIBE TABLE, não expõem informações restritas.
- O Amazon EMR no EKS restringe o acesso aos registros do Spark do driver do sistema em trabalhos habilitados para Lake Formation. Como o driver do sistema é executado com mais

Considerações e limitações 378

acesso, os eventos e logs que o driver do sistema gera podem incluir informações confidenciais. Para evitar que usuários ou códigos não autorizados acessem esses dados confidenciais, o EMR no EKS desativou o acesso aos registros do driver do sistema. Para solucionar problemas, entre em contato com AWS o suporte.

- Se você registrou uma localização de tabela no Lake Formation, o caminho de acesso aos dados passa pelas credenciais armazenadas do Lake Formation, independentemente da permissão do IAM para a função de execução de tarefas do EMR no EKS. Se você configurar incorretamente a função registrada com o local da tabela, os trabalhos enviados que usam a função com a permissão do S3 IAM para o local da tabela falharão.
- Gravar em uma tabela do Lake Formation usa a permissão do IAM em vez das permissões concedidas pelo Lake Formation. Se sua função de execução de tarefas tiver as permissões necessárias do S3, você poderá usá-la para executar operações de gravação.

Observe estas considerações e limitações ao usar o Apache Iceberg:

- Você só pode usar o Apache Iceberg com o catálogo de sessões e não com catálogos nomeados arbitrariamente.
- As tabelas do Iceberg registradas no Lake Formation oferecem suporte apenas às tabelas de metadados history, metadata_log_entries, snapshots, files, manifests e refs. O Amazon EMR oculta as colunas que podem conter dados confidenciais, como partitions, path e summaries. Essa limitação não se aplica às tabelas do Iceberg que não estão registradas no Lake Formation.
- As tabelas que você não registra no Lake Formation oferecem suporte a todos os procedimentos armazenados do Iceberg. Os procedimentos register table e migrate não são compatíveis com nenhuma tabela.
- Recomendamos que você use o Iceberg DataFrameWriter V2 em vez do V1.

Para obter mais informações, consulte Entendendo os conceitos e a terminologia do Amazon EMR no EKS e Habilitar o acesso ao cluster para o Amazon EMR no EKS.

Aviso legal para administradores de dados



Note

Ao conceder acesso aos recursos do Lake Formation a uma função do IAM para o EMR no EKS, você deve garantir que o administrador ou operador do cluster do EMR seja um

Considerações e limitações 379 administrador confiável. Isso é particularmente relevante para os recursos do Lake Formation que são compartilhados entre várias organizações e AWS contas.

Responsabilidades dos administradores do EKS

- O System namespace deve ser protegido. Nenhum usuário, recurso, entidade ou ferramenta teria permissão para ter qualquer permissão RBAC do Kubernetes nos recursos do Kubernetes no namespace. System
- Nenhum usuário, recurso ou entidade, exceto o EMR no serviço EKS, deve ter CREATE acesso ao POD, CONFIG_MAP e SECRET no namespace. User
- Systemdrivers e System executores contêm dados confidenciais. Portanto, os eventos do Spark, os registros do driver do Spark e os registros do executor do Spark no System namespace não devem ser encaminhados para sistemas externos de armazenamento de registros.

Solução de problemas

Registro em log

O EMR no EKS usa perfis de recursos do Spark para dividir a execução do trabalho. O Amazon EMR no EKS usa o perfil do usuário para executar o código que você forneceu, enquanto o perfil do sistema aplica as políticas do Lake Formation. Você pode acessar os registros dos contêineres executados como perfil de usuário configurando a StartJobRun solicitação com MonitoringConfiguration.

Servidor de histórico do Spark

O Spark History Server tem todos os eventos do Spark gerados a partir do perfil do usuário e os eventos editados gerados pelo driver do sistema. Você pode ver todos os contêineres dos drivers do usuário e do sistema na guia Executors. No entanto, os links de logs estão disponíveis somente para o perfil do usuário.

O trabalho falhou com permissões insuficientes do Lake Formation

Certifique-se de que sua função de execução de tarefas tenha as permissões para ser executada SELECT e estar DESCRIBE na tabela que você está acessando.

Solução de problemas 380

Falha na execução do trabalho com RDD

Atualmente, o EMR no EKS não oferece suporte a operações resilientes de conjunto de dados distribuído (RDD) em trabalhos habilitados para Lake Formation.

Não é possível acessar arquivos de dados no Amazon S3

Certifique-se de ter registrado a localização do data lake no Lake Formation.

Exceção de validação de segurança

O EMR no EKS detectou um erro de validação de segurança. Entre em contato com o AWS suporte para obter assistência.

Compartilhando o catálogo de dados e tabelas do AWS Glue entre contas

Você pode compartilhar bancos de dados e tabelas entre contas e ainda usar o Lake Formation.

Para obter mais informações, consulte Compartilhamento de dados entre contas no Lake Formation e Como faço para compartilhar o catálogo de dados e tabelas do AWS Glue entre contas usando o AWS Lake Formation?

Erro de inicialização do Iceberg Job que não define a região AWS

A mensagem é a seguinte:

```
25/02/25 13:33:19 ERROR SparkFGACExceptionSanitizer: Client received error with id =
b921f9e6-f655-491f-b8bd-b2842cdc20c7,
reason = IllegalArgumentException, message = Cannot initialize
LakeFormationAwsClientFactory, please set client.region to a valid aws region
```

Certifique-se de que a configuração do Spark

spark.sql.catalog.catalog_name.client.region esteja definida como uma região válida.

Lançamento de empregos no Iceberg SparkUnsupportedOperationException

A mensagem é a seguinte:

```
25/02/25 13:53:15 ERROR SparkFGACExceptionSanitizer: Client received error with id = 921fef42-0800-448b-bef5-d283d1278ce0,
reason = SparkUnsupportedOperationException, message = Either glue.id or glue.account-id is set with non-default account.
```

Solução de problemas 381

Cross account access with fine-grained access control is only supported with AWS Resource Access Manager.

Certifique-se de que a configuração do Spark

spark.sql.catalog.catalog_name.glue.account-id esteja definida como um ID de conta válido.

Registro em log e monitoramento

Para detectar incidentes, receber alertas quando incidentes ocorrerem e responder a eles, use estas opções com o Amazon EMR no EKS:

- Monitore o Amazon EMR no EKS com AWS CloudTrail <u>AWS CloudTrail</u>fornece um registro das ações realizadas por um usuário, função ou AWS serviço no Amazon EMR no EKS. Ele captura as chamadas do console do Amazon EMR e as chamadas de código para as operações de API do Amazon EMR no EKS como eventos. Dessa forma, você pode determinar a solicitação feita ao Amazon EMR no EKS, o endereço IP do qual a solicitação foi feita, quem fez a solicitação, quando ela foi feita e detalhes adicionais. Para obter mais informações, consulte <u>Registro do Amazon EMR</u> em chamadas de API EKS usando AWS CloudTrail.
- Use CloudWatch Events with Amazon EMR on EKS CloudWatch Events fornece um fluxo quase em tempo real de eventos do sistema que descrevem mudanças nos AWS recursos. CloudWatch Os eventos ficam cientes das mudanças operacionais à medida que elas ocorrem, respondem a elas e tomam medidas corretivas conforme necessário, enviando mensagens para responder ao ambiente, ativando funções, fazendo alterações e capturando informações de estado. Para usar CloudWatch Eventos com o Amazon EMR no EKS, crie uma regra que seja acionada em uma chamada de API do Amazon EMR no EKS via. CloudTrail Para obter mais informações, consulte Monitore trabalhos com a Amazon CloudWatch Events.

Criptografando o Amazon EMR em registros do EKS com armazenamento gerenciado

As seções a seguir mostram como configurar a criptografia para registros.

Habilitar criptografia

Para criptografar registros no armazenamento gerenciado com sua própria chave KMS, use a configuração a seguir ao enviar uma execução de trabalho.

A allowAWSToRetainLogs configuração permite AWS reter os registros do namespace do sistema ao executar um trabalho usando o FGAC nativo. A persistentAppUI configuração permite AWS salvar registros de eventos que são usados para gerar a interface do usuário do Spark. O encryptionKeyArn é usado para especificar o ARN da chave KMS que você deseja usar para criptografar os registros armazenados. AWS

Permissões necessárias para criptografia de logs

O usuário que envia o trabalho ou visualiza a interface do usuário do Spark deve ter permissão kms:Decrypt para realizar as ações e kms:DescribeKey receber kms:GenerateDataKey a chave de criptografia. Essas permissões são usadas para verificar a validade da chave e verificar se o usuário tem as permissões necessárias para ler e gravar registros criptografados com a chave KMS. Se o usuário que envia o trabalho não tiver as permissões de chave necessárias, o Amazon EMR no EKS rejeitará o envio da execução do trabalho.

Exemplo de política do IAM para a função usada para chamar StartJobRun

JSON

Criptografia de logs 383

```
{
    "Action": [
        "kms:DescribeKey",
        "kms:GenerateDataKey"
],
    "Resource": [
        "arn:aws:kms:*:*:key/key-id"
],
    "Effect": "Allow",
    "Sid": "AllowKMSDescribekey"
}
]
```

Você também deve configurar a chave KMS para permitir que persistentappui.elasticmapreduce.amazonaws.com os diretores elasticmapreduce.amazonaws.com de serviço e. kms:GenerateDataKey kms:Decrypt Isso permite que o EMR leia e grave registros criptografados com a chave KMS no armazenamento gerenciado.

Exemplo de política de chaves do KMS

JSON

Criptografia de logs 384

```
"Sid": "AllowKMSDescribekey"
    },
      "Effect": "Allow",
      "Action": [
        "kms:Decrypt",
        "kms:GenerateDataKey"
      ],
      "Resource": [
        11 * 11
      ],
      "Condition": {
        "StringLike": {
          "kms:viaService": "emr-containers.*.amazonaws.com",
          "kms:EncryptionContext:aws:emr-containers:virtualClusterId": "virtual
 cluster id"
        }
      },
      "Sid": "AllowKMSDecryptGenerate"
    },
    {
      "Effect": "Allow",
      "Action": [
        "kms:Decrypt",
        "kms:GenerateDataKey"
      ],
      "Resource": [
        11 * 11
      ],
      "Condition": {
        "StringLike": {
          "kms:EncryptionContext:aws:emr-containers:virtualClusterId": "virtual
 cluster id"
        },
        "ArnLike": {
          "aws:SourceArn": "arn:aws:emr-containers:*:*:/
virtualclusters/virtual_cluster_id"
        }
      },
      "Sid": "AllowKMSDecryptService"
    }
 ]
}
```

Criptografia de logs 385

Como prática recomendada de segurança, recomendamos que você adicione as aws:SourceArn condições kms:viaServicekms:EncryptionContext, e. Essas condições ajudam a garantir que a chave seja usada somente pelo Amazon EMR no EKS e usada somente para registros gerados a partir de trabalhos executados em um cluster virtual específico.

Registro do Amazon EMR em chamadas de API EKS usando AWS CloudTrail

O Amazon EMR no EKS é integrado com AWS CloudTrail, um serviço que fornece um registro das ações realizadas por um usuário, função ou AWS serviço no Amazon EMR no EKS. CloudTrail captura todas as chamadas de API para o Amazon EMR no EKS como eventos. As chamadas capturadas incluem as chamadas do console do Amazon EMR no EKS e as chamadas de código para as operações de API do Amazon EMR no EKS. Se você criar uma trilha, poderá habilitar a entrega contínua de CloudTrail eventos para um bucket do Amazon S3, incluindo eventos para o Amazon EMR no EKS. Se você não configurar uma trilha, ainda poderá ver os eventos mais recentes no CloudTrail console no Histórico de eventos. Usando as informações coletadas por CloudTrail, você pode determinar a solicitação que foi feita ao Amazon EMR no EKS, o endereço IP a partir do qual a solicitação foi feita, quem fez a solicitação, quando ela foi feita e detalhes adicionais.

Para saber mais CloudTrail, consulte o Guia AWS CloudTrail do usuário.

Informações sobre o Amazon EMR sobre EKS em CloudTrail

CloudTrail é ativado em sua AWS conta quando você cria a conta. Quando a atividade ocorre no Amazon EMR no EKS, essa atividade é registrada em um CloudTrail evento junto com outros eventos de AWS serviço no histórico de eventos. Você pode visualizar, pesquisar e baixar eventos recentes em sua AWS conta. Para obter mais informações, consulte <u>Visualização de eventos com histórico de CloudTrail eventos</u>.

Para um registro contínuo dos eventos em sua AWS conta, incluindo eventos do Amazon EMR no EKS, crie uma trilha. Uma trilha permite CloudTrail entregar arquivos de log para um bucket do Amazon S3. Por padrão, quando você cria uma trilha no console, a trilha se aplica a todas as AWS regiões. A trilha registra eventos de todas as regiões na AWS partição e entrega os arquivos de log ao bucket do Amazon S3 que você especificar. Além disso, você pode configurar outros AWS serviços para analisar e agir com base nos dados de eventos coletados nos CloudTrail registros. Para obter mais informações, consulte:

Visão geral da criação de uma trilha

CloudTrail troncos 386

- CloudTrail serviços e integrações suportados
- Configurando notificações do Amazon SNS para CloudTrail
- Recebendo arquivos de CloudTrail log de várias regiões e Recebendo arquivos de CloudTrail log de várias contas

Todas as ações do Amazon EMR no EKS são registradas CloudTrail e documentadas na documentação da API Amazon <u>EMR</u> on EKS. Por exemplo, chamadas para o CreateVirtualCluster StartJobRun e ListJobRuns as ações geram entradas nos arquivos de CloudTrail log.

Cada entrada de log ou evento contém informações sobre quem gerou a solicitação. As informações de identidade ajudam a determinar o seguinte:

- Se a solicitação foi feita com credenciais de usuário root ou AWS Identity and Access Management (IAM).
- Se a solicitação foi feita com credenciais de segurança temporárias de uma função ou de um usuário federado.
- Se a solicitação foi feita por outro AWS serviço.

Para obter mais informações, consulte o elemento Identidade CloudTrail do usuário.

Noções básicas sobre entradas de arquivos de log do Amazon EMR no EKS

Uma trilha é uma configuração que permite a entrega de eventos como arquivos de log para um bucket do Amazon S3 que você especificar. CloudTrail os arquivos de log contêm uma ou mais entradas de log. Um evento representa uma única solicitação de qualquer fonte e inclui informações sobre a ação solicitada, a data e a hora da ação, os parâmetros da solicitação e assim por diante. CloudTrail os arquivos de log não são um rastreamento de pilha ordenado das chamadas públicas de API, portanto, eles não aparecem em nenhuma ordem específica.

O exemplo a seguir mostra uma entrada de CloudTrail registro que demonstra a <u>ListJobRuns</u>ação.

```
{
  "eventVersion": "1.05",
  "userIdentity": {
    "type": "AssumedRole",
    "principalId": "AIDACKCEVSQ6C2EXAMPLE:admin",
```

CloudTrail troncos 387

```
"arn": "arn:aws:sts::012345678910:assumed-role/Admin/admin",
    "accountId": "012345678910",
    "accessKeyId": "AKIAIOSFODNN7EXAMPLE",
    "sessionContext": {
      "sessionIssuer": {
        "type": "Role",
        "principalId": "AIDACKCEVSQ6C2EXAMPLE",
        "arn": "arn:aws:iam::012345678910:role/Admin",
        "accountId": "012345678910",
        "userName": "Admin"
      },
      "webIdFederationData": {},
      "attributes": {
        "mfaAuthenticated": "false",
        "creationDate": "2020-11-04T21:49:36Z"
      }
    }
  },
  "eventTime": "2020-11-04T21:52:58Z",
  "eventSource": "emr-containers.amazonaws.com",
  "eventName": "ListJobRuns",
  "awsRegion": "us-east-1",
  "sourceIPAddress": "203.0.113.1",
  "userAgent": "aws-cli/1.11.167 Python/2.7.10 Darwin/16.7.0 botocore/1.7.25",
  "requestParameters": {
    "virtualClusterId": "1K48XXXXXXHCB"
  },
  "responseElements": null,
  "requestID": "890b8639-e51f-11e7-b038-EXAMPLE",
  "eventID": "874f89fa-70fc-4798-bc00-EXAMPLE",
  "readOnly": true,
  "eventType": "AwsApiCall",
  "recipientAccountId": "012345678910"
}
```

Uso do Amazon S3 Access Grants com o Amazon EMR no EKS

Visão geral do S3 Access Grants para o Amazon EMR no EKS

Com as versões 6.15.0 e superiores do Amazon EMR, o Amazon S3 Access Grants fornece uma solução de controle de acesso escalável que você pode usar para aumentar o acesso aos seus dados do Amazon S3 no Amazon EMR no EKS. Se você tiver uma configuração de permissão

Concessão de Acesso do S3 388

complexa ou grande para os dados do S3, poderá usar a funcionalidade Access Grants para escalar as permissões de dados do S3 para usuários, perfis e aplicações.

Use o S3 Access Grants para aumentar o acesso aos dados do Amazon S3 além das permissões concedidas pelo perfil de runtime ou pelos perfis do IAM que estão anexados às identidades com acesso ao seu cluster do Amazon EMR no EKS.

Para obter mais informações, consulte <u>Managing access with S3 Access Grants for Amazon EMR</u> no Guia de gerenciamento do Amazon EMR e <u>Gerenciar o acesso com o S3 Access Grants</u> no Guia do usuário do Amazon Simple Storage Service.

Esta página descreve os requisitos para executar um trabalho do Spark no Amazon EMR no EKS com a integração do S3 Access Grants. Com o Amazon EMR no EKS, o S3 Access Grants exige uma declaração adicional de política do IAM no perfil de execução do seu trabalho e uma configuração adicional de substituição da API StartJobRun. Para ver as etapas de configuração do S3 Access Grants com outras implantações do Amazon EMR, consulte a seguinte documentação:

- Using S3 Access Grants with Amazon EMR
- Using S3 Access Grants with EMR Serverless

Execução de um cluster do Amazon EMR no EKS com o S3 Access Grants para gerenciamento de dados

Você pode habilitar o S3 Access Grants no Amazon EMR no EKS e executar um trabalho do Spark. Quando sua aplicação solicita dados do S3, o Amazon S3 fornece credenciais temporárias que têm como escopo o bucket, prefixo ou objeto específico.

1. Configure um perfil de execução de trabalhos para o cluster do Amazon EMR no EKS. Inclua as permissões do IAM necessárias para executar os trabalhos do Spark, s3:GetDataAccess e s3:GetAccessGrantsInstanceForPrefix:

```
"Effect": "Allow",
"Action": [
    "s3:GetDataAccess",
    "s3:GetAccessGrantsInstanceForPrefix"
],
"Resource": [ //LIST ALL INSTANCE ARNS THAT THE ROLE IS ALLOWED TO QUERY
    "arn:aws_partition:s3:Region:account-id1:access-grants/default",
```

Executar um cluster 389

```
"arn:aws_partition:s3:Region:account-id2:access-grants/default"
]
}
```

Note

Se você especificar perfis do IAM para a execução de trabalhos que tenham permissões adicionais de acesso direto ao S3, os usuários poderão acessar os dados independentemente das permissões definidas no S3 Access Grants.

2. Envie um trabalho para o cluster do Amazon EMR no EKS com um rótulo de versão 6.15 ou superior do Amazon EMR e a classificação emrfs-site, conforme mostra o exemplo a seguir. Substitua os valores em red text pelos valores apropriados ao seu cenário de uso.

```
{
  "name": "myjob",
  "virtualClusterId": "123456",
  "executionRoleArn": "iam_role_name_for_job_execution",
  "releaseLabel": "emr-7.9.0-latest",
  "jobDriver": {
    "sparkSubmitJobDriver": {
      "entryPoint": "entryPoint_location",
      "entryPointArguments": ["argument1", "argument2"],
       "sparkSubmitParameters": "--class main_class"
    }
  },
  "configurationOverrides": {
    "applicationConfiguration": [
      {
        "classification": "emrfs-site",
        "properties": {
          "fs.s3.s3AccessGrants.enabled": "true",
          "fs.s3.s3AccessGrants.fallbackToIAM": "false"
      }
    ],
  }
}
```

Executar um cluster 390

Considerações sobre o S3 Access Grants com o Amazon EMR no EKS

Para obter informações importantes sobre suporte, compatibilidade e comportamento ao usar o Amazon S3 Access Grants com o Amazon EMR no EKS, consulte Considerações sobre o S3 Access Grants com o Amazon EMR no Guia de gerenciamento do Amazon EMR.

Validação de conformidade para o Amazon EMR no EKS

Auditores terceirizados avaliam a segurança e a conformidade do Amazon EMR no EKS como parte de AWS vários programas de conformidade. Isso inclui SOC, PCI, FedRAMP, HIPAA e outros.

Resiliência no Amazon EMR no EKS

A infraestrutura AWS global é construída em torno de AWS regiões e zonas de disponibilidade. AWS As regiões fornecem várias zonas de disponibilidade fisicamente separadas e isoladas, conectadas a redes de baixa latência, alta taxa de transferência e alta redundância. Com as zonas de disponibilidade, é possível projetar e operar aplicações e bancos de dados que automaticamente executam o failover entre as zonas sem interrupção. As zonas de disponibilidade são altamente disponíveis, tolerantes a falhas e escaláveis que uma ou várias infraestruturas de data center tradicionais.

Para obter mais informações sobre AWS regiões e zonas de disponibilidade, consulte <u>Infraestrutura</u> <u>AWS global</u>.

Além da infraestrutura AWS global, o Amazon EMR no EKS oferece integração com o Amazon S3 por meio do EMRFS para ajudar a suportar suas necessidades de resiliência e backup de dados.

Segurança da infraestrutura no Amazon EMR no EKS

Como um serviço gerenciado, o Amazon EMR é protegido pela segurança de rede AWS global. Para obter informações sobre serviços AWS de segurança e como AWS proteger a infraestrutura, consulte AWS Cloud Security. Para projetar seu AWS ambiente usando as melhores práticas de segurança de infraestrutura, consulte Proteção de infraestrutura no Security Pillar AWS Well-Architected Framework.

Você usa chamadas de API AWS publicadas para acessar o Amazon EMR pela rede. Os clientes devem oferecer compatibilidade com:

Considerações 391

- Transport Layer Security (TLS). Exigimos TLS 1.2 e recomendamos TLS 1.3.
- Conjuntos de criptografia com perfect forward secrecy (PFS) como DHE (Ephemeral Diffie-Hellman) ou ECDHE (Ephemeral Elliptic Curve Diffie-Hellman). A maioria dos sistemas modernos, como Java 7 e versões posteriores, comporta esses modos.

Além disso, as solicitações devem ser assinadas usando um ID da chave de acesso e uma chave de acesso secreta associada a uma entidade principal do IAM. Ou você pode usar o <u>AWS</u>

<u>Security Token Service</u> (AWS STS) para gerar credenciais de segurança temporárias para assinar solicitações.

Análise de configuração e vulnerabilidade

AWS lida com tarefas básicas de segurança, como sistema operacional (SO) convidado e aplicação de patches em bancos de dados, configuração de firewall e recuperação de desastres. Esses procedimentos foram revisados e certificados por terceiros certificados. Para obter mais detalhes, consulte os seguintes recursos da :

- Validação de conformidade para o Amazon EMR no EKS
- · Modelo de responsabilidade compartilhada
- Amazon Web Services: visão geral do processo de segurança (whitepaper)

Conexão com o Amazon EMR no EKS usando um endpoint da VPC de interface

Você pode se conectar diretamente ao Amazon EMR no EKS usando <u>endpoints de interface VPC</u> (AWS PrivateLink) em sua Virtual Private Cloud (VPC) em vez de se conectar pela Internet. Quando você usa uma interface VPC endpoint, a comunicação entre sua VPC e o Amazon EMR no EKS é conduzida inteiramente dentro da rede. AWS Cada VPC endpoint é representado por uma ou mais <u>interfaces de rede elástica (ENIs)</u> com endereços IP privados em suas sub-redes VPC.

A interface VPC endpoint conecta sua VPC diretamente ao Amazon EMR no EKS sem um gateway de internet, dispositivo NAT, conexão VPN ou conexão Direct Connect. AWS As instâncias em sua VPC não precisam de endereços IP públicos para se comunicarem com a API do Amazon EMR no EKS.

Você pode criar uma interface VPC endpoint para se conectar ao Amazon EMR no EKS usando os AWS Management Console comandos or (). AWS Command Line Interface AWS CLI Para obter mais informações, consulte Criação de um endpoint de interface.

Após criar um endpoint da VPC de interface, se você habilitar nomes de host DNS privados para o endpoint, o endpoint padrão do Amazon EMR no EKS será resolvido para seu endpoint da VPC. O endpoint de nome de serviço padrão para o Amazon EMR no EKS estará no formato a seguir.

```
emr-containers.Region.amazonaws.com
```

Se você não habilitar nomes de host DNS privados, a Amazon VPC fornecerá um nome de endpoint DNS que poderá ser usado no formato a seguir.

```
VPC_Endpoint_ID.emr-containers.Region.vpce.amazonaws.com
```

Para obter mais informações, consulte <u>Interface VPC Endpoints (AWS PrivateLink)</u> no Guia do usuário da Amazon VPC. O Amazon EMR no EKS oferece suporte a chamadas para todas as <u>ações</u> de API dentro da sua VPC.

Você pode anexar políticas de endpoint da VPC a um endpoint da VPC para controlar o acesso de entidades principais do IAM. Também é possível associar grupos de segurança a um VPC endpoint para controlar o acesso de entrada e saída com base na origem e no destino do tráfego de rede, como um intervalo de endereços IP. Para obter mais informações, consulte Controlar o acesso a serviços com endpoints da VPC.

Criação de uma política de endpoint da VPC para o Amazon EMR no EKS

É possível criar uma política para endpoints da Amazon VPC para o Amazon EMR no EKS com a finalidade de especificar o seguinte:

- O principal que pode ou não executar ações
- As ações que podem ser executadas
- Os recursos nos quais as ações podem ser executadas

Para obter mais informações, consulte <u>Controlling Access to Services with VPC Endpoints</u> no Guia do usuário da Amazon VPC.

Example Política de VPC Endpoint para negar todo o acesso de uma conta especificada AWS

A política de VPC endpoint a seguir nega à AWS conta 123456789012 todo o acesso aos recursos que usam o endpoint.

```
{
    "Statement": [
        {
             "Action": "*",
             "Effect": "Allow",
             "Resource": "*",
             "Principal": "*"
        },
        {
             "Action": "*",
             "Effect": "Deny",
             "Resource": "*",
             "Principal": {
                 "AWS": [
                     "123456789012"
                 ]
             }
        }
    ]
}
```

Example Política de endpoint da VPC para permitir o acesso à VPC somente a uma entidade principal do IAM (usuário) especificada

A política de VPC endpoint a seguir permite acesso total somente ao usuário *lijuan* do IAM na conta. AWS *123456789012* Todos os outros principais IAM têm acesso negado usando o endpoint.

```
}
]
}
```

Example Política de endpoint da VPC para permitir operações somente leitura do Amazon EMR no EKS

A política de VPC endpoint a seguir permite que somente 123456789012 a AWS conta execute as ações especificadas do Amazon EMR no EKS.

As ações especificadas fornecem o equivalente ao acesso somente leitura para o Amazon EMR no EKS. Todas as outras ações na VPC serão negadas para a conta especificada. Todas as outras contas terão acesso negado. Para obter uma lista de ações do Amazon EMR no EKS, consulte Ações, recursos e chaves de condição para o Amazon EMR no EKS.

```
{
    "Statement": [
        {
            "Action": [
                "emr-containers:DescribeJobRun",
                "emr-containers:DescribeVirtualCluster",
                "emr-containers:ListJobRuns",
                "emr-containers:ListTagsForResource",
                "emr-containers:ListVirtualClusters"
            ],
            "Effect": "Allow",
            "Resource": "*",
            "Principal": {
                "AWS": [
                     "123456789012"
                ]
            }
        }
    ]
}
```

Example Política de endpoint da VPC com negação de acesso a um cluster virtual especificado

A política de VPC endpoint a seguir permite acesso total a todas as contas e entidades principais, mas nega qualquer acesso da AWS conta 123456789012 às ações executadas no cluster virtual com o ID do cluster. A1B2CD34EF5G Outras ações do Amazon EMR no EKS que não oferecem

suporte a permissões em nível de recurso para clusters virtuais ainda são permitidas. Para obter uma lista de ações do Amazon EMR no EKS e seus tipos de recursos correspondentes, consulte <u>Ações, recursos e chaves de condição para o Amazon EMR no EKS</u> no Guia do usuário do AWS Identity and Access Management.

```
{
    "Statement": [
        {
             "Action": "*",
             "Effect": "Allow",
            "Resource": "*",
            "Principal": "*"
        },
        {
            "Action": "*",
            "Effect": "Deny",
             "Resource": "arn:aws:emr-containers:us-west-2:123456789012:/
virtualclusters/A1B2CD34EF5G",
             "Principal": {
                 "AWS": Γ
                     "123456789012"
                 ]
            }
        }
    ]
}
```

Configuração do acesso entre contas para o Amazon EMR no EKS

Você pode configurar o acesso entre contas para o Amazon EMR no EKS. O acesso entre contas permite que os usuários de uma AWS conta executem o Amazon EMR em trabalhos do EKS e acessem os dados subjacentes que pertencem a AWS outra conta.

Pré-requisitos

Para configurar o acesso entre contas para o Amazon EMR no EKS, você concluirá as tarefas enquanto estiver conectado às seguintes AWS contas:

 AccountA- Uma AWS conta em que você criou um cluster virtual do Amazon EMR no EKS registrando o Amazon EMR com um namespace em um cluster EKS.

Acesso entre contas 396

 AccountB- Uma AWS conta que contém um bucket do Amazon S3 ou uma tabela do DynamoDB que você deseja que suas tarefas do Amazon EMR no EKS acessem.

Você deve ter o seguinte em mãos em suas AWS contas antes de configurar o acesso entre contas:

- Um cluster virtual do Amazon EMR no EKS na AccountA em que deseja executar trabalhos.
- Um perfil de execução de trabalho na AccountA que tem as permissões obrigatórias para executar trabalhos no cluster virtual. Para obter mais informações, consulte <u>Criação de um perfil de</u> execução de trabalho e <u>Uso de perfis de execução de trabalho com o Amazon EMR no EKS</u>.

Como acessar um bucket do Amazon S3 ou uma tabela do DynamoDB entre contas

Para configurar o acesso entre contas do Amazon EMR no EKS, conclua as etapas apresentadas a seguir.

- Crie um bucket do Amazon S3, cross-account-bucket, na AccountB. Para mais informações, consulte <u>Criar um bucket</u>. Se desejar ter acesso entre contas para o DynamoDB, você também pode criar uma tabela do DynamoDB na AccountB. Para obter mais informações, consulte <u>Creating a DynamoDB table</u>.
- 2. Crie um perfil do IAM Cross-Account-Role-B na AccountB que possa acessar o crossaccount-bucket.
 - 1. Faça login no console do IAM.
 - 2. Escolha Perfis e crie um novo perfil: Cross-Account-Role-B. Para obter mais informações sobre como criar perfis do IAM, consulte Criação de perfis do IAM no Guia do usuário do IAM.
 - 3. Crie uma política do IAM que especifique as permissões para Cross-Account-Role-B acessar o bucket cross-account-bucket do S3, como demonstra a instrução de política a seguir. Em seguida, anexe a política do IAM ao Cross-Account-Role-B. Para obter mais informações, consulte Creating a New Policy no Guia do usuário do IAM.

JSON

```
{
    "Version": "2012-10-17",
    "Statement": [
        {
```

```
"Effect": "Allow",
    "Action": [
        "s3:*"
],
    "Resource": [
        "arn:aws:s3:::cross-account-bucket",
        "arn:aws:s3:::cross-account-bucket/*"
],
    "Sid": "AllowS3"
}
]
```

Se o acesso ao DynamoDB for necessário, crie uma política do IAM que especifique as permissões de acesso à tabela do DynamoDB entre contas. Em seguida, anexe a política do IAM ao Cross-Account-Role-B. Para obter mais informações, consulte Criar uma tabela do DynamoDB no Guia do usuário do IAM.

A seguir, é apresentada uma política de acesso a uma tabela do DynamoDB, CrossAccountTable.

JSON

3. Edite a relação de confiança para o perfil Cross-Account-Role-B.

- 1. Para configurar a relação de confiança para o perfil, escolha a guia Relações de confiança no console do IAM para o perfil criado na Etapa 2: Cross-Account-Role-B.
- 2. Selecione Editar relação de confiança.
- 3. Adicione o documento de política a seguir, que permite que Job-Execution-Role-A na AccountA assuma esse perfil Cross-Account-Role-B.

JSON

- 4. Conceda que o Job-Execution-Role-A na AccountA tenha a permissão sts assume-role para assumir Cross-Account-Role-B.
 - 1. No console do IAM da AWS contaAccountA, selecioneJob-Execution-Role-A.
 - 2. Adicione a instrução de política a seguir ao Job-Execution-Role-A para permitir a ação AssumeRole no perfil Cross-Account-Role-B.

JSON

```
{
  "Version": "2012-10-17",
  "Statement": [
     {
        "Effect": "Allow",
        "Action": [
            "sts:AssumeRole"
      ],
```

```
"Resource": [
    "arn:aws:iam::*:role/Cross-Account-Role-B"
],
    "Sid": "AllowSTSAssumerole"
}
]
```

 Para obter acesso ao Amazon S3, defina os parâmetros spark-submit (spark conf) apresentados a seguir ao enviar o trabalho para o Amazon EMR no EKS.



Por padrão, o EMRFS usa o perfil de execução do trabalho para acessar o bucket do S3 usando o trabalho. Entretanto, quando customAWSCredentialsProvider é definido como AssumeRoleAWSCredentialsProvider, o EMRFS usa o perfil correspondente que você especifica com ASSUME_ROLE_CREDENTIALS_ROLE_ARN em vez do Job-Execution-Role-A para obter acesso ao Amazon S3.

- --conf spark.hadoop.fs.s3.customAWSCredentialsProvider=com.amazonaws.emr.AssumeRol
- --conf
 spark.kubernetes.driverEnv.ASSUME_ROLE_CREDENTIALS_ROLE_ARN=arn:aws:iam::Ac
 Cross-Account-Role-B \
- --conf
 spark.executorEnv.ASSUME_ROLE_CREDENTIALS_ROLE_ARN=arn:aws:iam::AccountB:ro
 Cross-Account-Role-B \

Note

Você deve definir ASSUME_ROLE_CREDENTIALS_ROLE_ARN para o env do executor e do driver na configuração de trabalho do Spark.

Para obter acesso entre contas do DynamoDB, você deve definir --conf spark.dynamodb.customAWSCredentialsProvider=com.amazonaws.emr.AssumeRoleAWSC1

6. Execute o trabalho do Amazon EMR no EKS com acesso entre contas, como demonstrado pelo exemplo a seguir.

```
aws emr-containers start-job-run \
--virtual-cluster-id 123456 \
--name myjob \
--execution-role-arn execution-role-arn \
--release-label emr-6.2.0-latest \
--job-driver '{"sparkSubmitJobDriver": {"entryPoint": "entryPoint_location",
"entryPointArguments": ["arguments_list"], "sparkSubmitParameters": "--class
<main_class> --conf spark.executor.instances=2 --conf spark.executor.memory=2G
--conf spark.executor.cores=2 --conf spark.driver.cores=1 --conf
spark.hadoop.fs.s3.customAWSCredentialsProvider=com.amazonaws.emr.AssumeRoleAWSCredentials
--conf
spark.kubernetes.driverEnv.ASSUME_ROLE_CREDENTIALS_ROLE_ARN=arn:aws:iam::AccountB:role/
Cross-Account-Role-B --conf
spark.executorEnv.ASSUME_ROLE_CREDENTIALS_ROLE_ARN=arn:aws:iam::AccountB:role/
Cross-Account-Role-B"}} ' \
--configuration-overrides '{"applicationConfiguration": [{"classification":
 "spark-defaults", "properties": {"spark.driver.memory": "2G"}}],
 "monitoringConfiguration": {"cloudWatchMonitoringConfiguration":
{"logGroupName": "log_group_name", "logStreamNamePrefix": "log_stream_prefix"},
 "persistentAppUI":"ENABLED", "s3MonitoringConfiguration": {"logUri": "s3://
my_s3_log_location" }}}'
```

Marcação de recursos do Amazon EMR no EKS

Para ajudar no gerenciamento de recursos do Amazon EMR no EKS, você pode atribuir seus próprios metadados a cada recurso usando etiquetas. Este tópico fornece uma visão geral da função das etiquetas e mostra como você pode criá-las.

Tópicos

- Conceitos Básicos de Tags
- Marcar com tag os recursos do
- Restrições de tag
- Trabalhe com tags usando o AWS CLI e o Amazon EMR na API EKS

Conceitos Básicos de Tags

Uma tag é um rótulo que você atribui a um AWS recurso. Cada tag consiste de uma chave e um valor opcional, que podem ser definidos.

As tags permitem que você categorize seus AWS recursos por atributos como finalidade, proprietário ou ambiente. Caso possua muitos recursos do mesmo tipo, você pode identificar rapidamente um recurso específico com base nas tags atribuídas a ele. Por exemplo, é possível definir um conjunto de etiquetas para os clusters do Amazon EMR no EKS para ajudar a rastrear o proprietário e o nível da pilha de cada cluster. Recomendamos planejar um conjunto consistente de chaves de etiquetas para cada tipo de recurso. Você pode pesquisar e filtrar os recursos de acordo com as tags que adicionar.

Tags não são automaticamente atribuídas aos recursos. Após adicionar uma tag, você pode editar as chaves e os valores das tags ou removê-las de um recurso a qualquer momento. Caso exclua um recurso, todas as respectivas tags também serão excluídas.

As etiquetas não têm significado semântico para o Amazon EMR no EKS e são interpretadas estritamente como uma string de caracteres.

Um valor de tag pode ser uma string vazia, mas não nula. Uma chave de etiqueta não pode ser uma string vazia. Se você adicionar uma etiqueta que tenha a mesma chave de uma etiqueta existente nesse recurso, o novo valor substituirá o antigo.

Conceitos Básicos de Tags 402

Se você usa AWS Identity and Access Management (IAM), você pode controlar quais usuários em sua AWS conta têm permissão para gerenciar tags.

Para obter exemplos de políticas de controle de acesso por etiquetas, consulte <u>Políticas para</u> controle de acesso baseado em etiquetas.

Marcar com tag os recursos do

Você pode etiquetar clusters virtuais e execuções de trabalhos novos ou existentes que estão em estados ativos. Os estados ativos para execuções de trabalhos incluem: PENDING, SUBMITTED, RUNNING e CANCEL_PENDING. Os estados ativos para clusters virtuais incluem: RUNNING, TERMINATING e ARRESTED. Para obter mais informações, consulte <u>Estados de execução de trabalho e Estados de um cluster virtual.</u>

Quando um cluster virtual é encerrado, as etiquetas são limpas e não ficam mais acessíveis.

Se você estiver usando o Amazon EMR na API EKS AWS CLI, o ou um AWS SDK, poderá aplicar tags a novos recursos usando o parâmetro tags na ação relevante da API. Você também pode aplicar etiquetas aos recursos usando a ação da API TagResource.

Você pode usar algumas ações de criação de recursos para especificar etiquetas para um recurso quando ele for criado. Nesse caso, se as etiquetas não puderem ser aplicadas enquanto o recurso estiver sendo criado, ele não será criado. Esse mecanismo garante que os recursos que você pretende etiquetar na criação sejam criados com as etiquetas especificadas ou não sejam criados. Se você etiquetar os recursos no momento da criação, não precisará executar os scripts de marcação personalizados após a criação do recurso.

A tabela a seguir descreve os recursos do Amazon EMR no EKS que podem ser etiquetados.

Recurso	Compatível com tags	Compatível com a propagação de tags	Suporta marcação na criação (Amazon EMR na API EKS e AWS CLI SDK) AWS	API para criação (etiquetas podem ser adicionad as durante a criação)
Cluster virtual	Sim	Não. As etiquetas	Sim	CreateVir tualCluster

Recurso	Compatível com tags	Compatível com a propagação de tags	Suporta marcação na criação (Amazon EMR na API EKS e AWS CLI SDK) AWS	API para criação (etiquetas podem ser adicionad as durante a criação)
		associadas a um cluster virtual não se propagam para execuções de trabalhos enviadas a esse cluster virtual.		
Execuções de trabalhos	Sim	Não	Sim	StartJobRun

Restrições de tag

As restrições básicas a seguir se aplicam a tags:

- Número máximo de tags por recurso 50
- Em todos os recursos, cada chave de tag deve ser exclusiva e possuir apenas um valor.
- Comprimento máximo da chave 128 caracteres Unicode em UTF-8
- Comprimento máximo do valor: 256 caracteres Unicode em UTF-8
- Se seu esquema de marcação for usado em vários AWS serviços e recursos, lembre-se de que outros serviços podem ter restrições quanto aos caracteres permitidos. Em geral, caracteres permitidos incluem letras, números, espaços representáveis em UTF-8 e os caracteres + - = . _ : / @.
- Chaves e valores de tags diferenciam maiúsculas de minúsculas.
- Um valor de tag pode ser uma string vazia, mas não nula. Uma chave de etiqueta não pode ser uma string vazia.

Restrições de tag 404

• Não use aws:, AWS: ou qualquer combinação de letras maiúsculas e minúsculas como prefixo para chaves ou valores. Eles são reservados apenas para AWS uso.

Trabalhe com tags usando o AWS CLI e o Amazon EMR na API EKS

Use os AWS CLI comandos a seguir ou o Amazon EMR nas operações da API EKS para adicionar, atualizar, listar e excluir as tags dos seus recursos.

Tarefa	AWS CLI	Ação API
Adicione ou substitua uma ou mais tags	tag-resource	TagResource
Lista de tags para um recurso	list-tags-for-resource	ListTagsForResource
Exclua uma ou mais tags	untag-resource	UntagResource

Os exemplos a seguir mostram como marcar ou desmarcar recursos usando AWS CLI.

Exemplo 1: etiquetar um cluster virtual existente

O comando a seguir realiza a marcação de um cluster virtual existente.

```
aws emr-containers tag-resource --resource-arn resource_ARN --tags team=devs
```

Exemplo 2: remover a etiqueta de um cluster virtual existente

O comando a seguir exclui uma etiqueta de um cluster virtual existente.

```
aws emr-containers untag-resource --resource-arn resource_ARN --tag-keys tag_key
```

Exemplo 3: listar etiquetas para um recurso

O comando a seguir lista as tags associadas a um recurso existente.

```
aws emr-containers list-tags-for-resource --resource-arn resource_ARN
```

Solução de problemas para o Amazon EMR no EKS

Esta seção descreve como solucionar problemas com o Amazon EMR no EKS. Para obter informações sobre como solucionar problemas gerais com o Amazon EMR, consulte <u>Troubleshoot a cluster</u> no Guia de gerenciamento do Amazon EMR.

Tópicos

- Solução de problemas de trabalhos que usam PersistentVolumeClaims (PVC)
- Solução de problemas de escalabilidade automática vertical do Amazon EMR no EKS
- Solução de problemas do operador do Spark do Amazon EMR no EKS

Solução de problemas de trabalhos que usam PersistentVolumeClaims (PVC)

Se você precisar criar, listar ou excluir PersistentVolumeClaims (PVC) para um trabalho, mas não adicionar permissões de PVC aos contêineres emr-emr da função padrão do Kubernetes, o trabalho falhará quando você o enviar. Sem essas permissões, o perfil emr-containers não pode criar os perfis necessários para o driver do Spark ou para o cliente do Spark. Não é suficiente adicionar permissões ao driver do Spark ou aos perfis de clientes, conforme sugerido pelas mensagens de erro. O perfil primário emr-containers também deve incluir as permissões necessárias. Esta seção explica como adicionar as permissões necessárias ao perfil primário emr-containers.

Verificação

Para verificar se o perfil emr-containers tem ou não as permissões necessárias, defina a variável NAMESPACE com seu próprio valor e execute o seguinte comando:

```
export NAMESPACE=YOUR_VALUE
kubectl describe role emr-containers -n ${NAMESPACE}
```

Além disso, para verificar se os perfis do Spark e dos clientes têm as permissões necessárias, execute o seguinte comando:

```
kubectl describe role emr-containers-role-spark-driver -n ${NAMESPACE}
kubectl describe role emr-containers-role-spark-client -n ${NAMESPACE}
```

Se as permissões não existirem, prossiga com a aplicação de patch, como apresentado a seguir.

Patch

- 1. Se os trabalhos estiverem em execução no momento sem as permissões, interrompa-os.
- 2. Crie um arquivo chamado RBAC_Patch.py da seguinte forma:

```
import os
import subprocess as sp
import tempfile as temp
import json
import argparse
import uuid
def delete_if_exists(dictionary: dict, key: str):
    if dictionary.get(key, None) is not None:
        del dictionary[key]
def doTerminalCmd(cmd):
    with temp.TemporaryFile() as f:
        process = sp.Popen(cmd, stdout=f, stderr=f)
        process.wait()
        f.seek(0)
        msg = f.read().decode()
    return msg
def patchRole(roleName, namespace, extraRules, skipConfirmation=False):
    cmd = f"kubectl get role {roleName} -n {namespace} --output json".split(" ")
    msg = doTerminalCmd(cmd)
    if "(NotFound)" in msg and "Error" in msg:
        print(msg)
        return False
    role = json.loads(msg)
    rules = role["rules"]
    rulesToAssign = extraRules[::]
    passedRules = []
    for rule in rules:
        apiGroups = set(rule["apiGroups"])
        resources = set(rule["resources"])
        verbs = set(rule["verbs"])
        for extraRule in extraRules:
            passes = 0
            apiGroupsExtra = set(extraRule["apiGroups"])
```

```
resourcesExtra = set(extraRule["resources"])
           verbsExtra = set(extraRule["verbs"])
           passes += len(apiGroupsExtra.intersection(apiGroups)) >=
len(apiGroupsExtra)
           passes += len(resourcesExtra.intersection(resources)) >=
len(resourcesExtra)
           passes += len(verbsExtra.intersection(verbs)) >= len(verbsExtra)
           if passes >= 3:
               if extraRule not in passedRules:
                   passedRules.append(extraRule)
                   if extraRule in rulesToAssign:
                       rulesToAssign.remove(extraRule)
               break
   prompt_text = "Apply Changes?"
   if len(rulesToAssign) == 0:
       print(f"The role {roleName} seems to already have the necessary
permissions!")
       prompt_text = "Proceed anyways?"
   for ruleToAssign in rulesToAssign:
       role["rules"].append(ruleToAssign)
   delete_if_exists(role, "creationTimestamp")
   delete_if_exists(role, "resourceVersion")
   delete_if_exists(role, "uid")
   new_role = json.dumps(role, indent=3)
   uid = uuid.uuid4()
   filename = f"Role-{roleName}-New_Permissions-{uid}-TemporaryFile.json"
   try:
       with open(filename, "w+") as f:
           f.write(new_role)
           f.flush()
       prompt = "y"
       if not skipConfirmation:
           prompt = input(
               doTerminalCmd(f"kubectl diff -f {filename}".split(" ")) +
f"\n{prompt_text} y/n: "
           ).lower().strip()
           while prompt != "y" and prompt != "n":
               prompt = input("Please make a valid selection. y/n:
").lower().strip()
       if prompt == "y":
           print(doTerminalCmd(f"kubectl apply -f {filename}".split(" ")))
   except Exception as e:
       print(e)
   os.remove(f"./{filename}")
```

```
if __name__ == '__main__':
    parser = argparse.ArgumentParser()
    parser.add_argument("-n", "--namespace",
                        help="Namespace of the Role. By default its the
 VirtualCluster's namespace",
                        required=True,
                        dest="namespace"
    parser.add_argument("-p", "--no-prompt",
                        help="Applies the patches without asking first",
                        dest="no_prompt",
                        default=False,
                        action="store_true"
    args = parser.parse_args()
    emrRoleRules = [
        {
            "apiGroups": [""],
            "resources": ["persistentvolumeclaims"],
            "verbs": ["list", "create", "delete", "patch"]
         }
    ]
    driverRoleRules = [
        {
            "apiGroups": [""],
            "resources": ["persistentvolumeclaims"],
            "verbs": ["list", "create", "delete", "patch", "deletecollection"]
        },
            "apiGroups": [""],
            "resources": ["services"],
            "verbs": ["get", "list", "describe", "create", "delete", "watch",
 "deletecollection"]
        },
        {
            "apiGroups": [""],
            "resources": ["configmaps", "pods"],
            "verbs": ["deletecollection"]
        }
```

3. Execute o script do Python:

```
python3 RBAC_Patch.py -n ${NAMESPACE}
```

- 4. O comando "kubectl diff" apresentará as diferenças entre as permissões novas e as antigas. Pressione Y para aplicar patches ao perfil.
- 5. Verifique os três perfis com permissões adicionais, como o apresentado:

```
kubectl describe role -n ${NAMESPACE}
```

6. Execute o script do Python:

```
python3 RBAC_Patch.py -n ${NAMESPACE}
```

- 7. Após executar o comando, o comando "kubectl diff" apresentará as diferenças entre as permissões novas e as antigas. Pressione Y para aplicar patches ao perfil.
- 8. Verifique os três perfis com permissões adicionais:

```
kubectl describe role -n ${NAMESPACE}
```

9. Envie o trabalho novamente.

Aplicação manual de patches

Se a permissão requerida pela sua aplicação se aplicar a algo diferente das regras de PVC, você poderá adicionar manualmente as permissões do Kubernetes ao cluster virtual do Amazon EMR, conforme necessário.



Note

O perfil emr-containers é um perfil primário. Isso significa que ele deve fornecer todas as permissões necessárias antes que você possa alterar os perfis subjacentes do driver ou dos clientes.

1. Faça o download das permissões atuais em arquivos YAML ao executar os comandos abaixo:

```
kubectl get role -n ${NAMESPACE} emr-containers -o yaml >> emr-containers-role-
patch.yaml
kubectl get role -n ${NAMESPACE} emr-containers-role-spark-driver -o yaml >> driver-
role-patch.yaml
kubectl get role -n ${NAMESPACE} emr-containers-role-spark-client -o yaml >> client-
role-patch.yaml
```

- 2. Com base na permissão requerida pela aplicação, edite cada arquivo e adicione regras adicionais, como as seguintes:
 - emr-containers-role-patch.yaml

```
- apiGroups:
 resources:
  - persistentvolumeclaims
 verbs:
  - list
  - create
  - delete
  - patch
```

driver-role-patch.yaml

```
apiGroups:
resources:
```

411 Aplicação manual de patches

```
- persistentvolumeclaims
  verbs:
  - list
  - create
  - delete
  - patch
  - deletecollection
- apiGroups:
  _ ""
  resources:
  - services
  verbs:
  - get
  - list
  - describe
  - create
  - delete
  - watch
  - deletecollection
- apiGroups:
  _ ""
  resources:
  - configmaps
  - pods
  verbs:
  - deletecollection
```

· client-role-patch.yaml

```
- apiGroups:
- ""
resources:
- persistentvolumeclaims
verbs:
- list
- create
- delete
- patch
```

- 3. Remova os atributos a seguir em conjunto com os valores deles. Isso é necessário para aplicar a atualização.
 - creationTimestamp
 - resourceVersion

- uid
- 4. Por fim, execute a aplicação de patches:

```
kubectl apply -f emr-containers-role-patch.yaml
kubectl apply -f driver-role-patch.yaml
kubectl apply -f client-role-patch.yaml
```

Solução de problemas de escalabilidade automática vertical do Amazon EMR no EKS

Consulte as seções a seguir se você encontrar problemas ao configurar o operador de escalabilidade automática vertical do Amazon EMR no EKS em um cluster do Amazon EKS com o Operator Lifecycle Manager. Para obter mais informações, incluindo as etapas para concluir a instalação, consulte Uso da escalabilidade automática vertical com trabalhos do Spark no Amazon EMR.

Erro 403 Forbidden

Se você seguiu as etapas em <u>Instalação do Operator Lifecycle Manager (OLM) no cluster do Amazon EKS</u>, executou o comando olm status e ele retornou um erro 403 Forbidden, como o mostrado abaixo, pode ser que você não tenha obtido os tokens de autenticação para o repositório do Amazon ECR para o operador.

Para resolver esse problema, repita a etapa descrita em <u>Instalação do operador de escalabilidade</u> <u>automática vertical do Amazon EMR no EKS</u> para obter os tokens. Em seguida, tente instalar novamente.

```
Error: FATA[0002] Failed to run bundle: pull bundle image: error pulling image IMAGE. error resolving name: unexpected status code [manifests latest]: 403 Forbidden
```

Namespace do Kubernetes não encontrado

Ao <u>configurar o operador de escalabilidade automática vertical do Amazon EMR no EKS</u> em um cluster do Amazon EKS, você poderá receber um erro namespaces not found, como o mostrado aqui:

```
FATA[0020] Failed to run bundle: create catalog: error creating catalog source: namespaces "NAME" not found.
```

Se o namespace especificado não existir, o OLM não instalará o operador de escalabilidade automática vertical. Para resolver esse problema, use o comando apresentado a seguir para criar o namespace. Em seguida, tente instalar novamente.

```
kubectl create namespace NAME
```

Erro ao salvar as credenciais do Docker

Para <u>configurar a escalabilidade automática vertical</u>, você deve autenticar e buscar as imagens do Docker relacionadas à escalabilidade automática vertical do Amazon EMR no EKS. Ao fazer isso, você poderá receber um erro como o seguinte se o Docker não estiver em execução:

```
aws ecr get-login-password \
    --region $REGION | docker login \
    --username AWS \
    --password-stdin $ACCOUNT_ID.dkr.ecr.$REGION.amazonaws.com

Error saving credentials: error storing credentials - err: exit status 1
out: 'Post "http://ipc/registry/credstore-updated": dial unix backend.sock: connect: no such file or directory'
```

Para resolver esse problema, confirme se o Docker está em execução ou abra o Docker Desktop. Em seguida, tente salvar as credenciais novamente.

Solução de problemas do operador do Spark do Amazon EMR no EKS

Consulte as seções a seguir se você encontrar problemas com o operador do Spark do Amazon EMR no EKS. Para obter mais informações, incluindo as etapas para concluir a instalação, consulte Execução de trabalhos do Spark com o operador do Spark.

Erro na instalação do chart do Helm

Se você seguiu as etapas em <u>Instalação do operador do Spark</u> e ele retornou um erro INSTALLATION FAILED, como o mostrado abaixo, ao tentar instalar ou verificar o chart do Helm, pode ser que você não tenha obtido os tokens de autenticação para o repositório do Amazon ECR para o operador.

Erro de credenciais do Docker 414

Para resolver esse problema, repita a descrita etapa em <u>Instalação do operador do Spark</u> para autenticar seu cliente Helm para o registro do Amazon ECR. Em seguida, tente realizar a etapa de instalação novamente.

```
Error: INSTALLATION FAILED: Kubernetes cluster unreachable: the server has asked for the client to provide credentials
```

UnsupportedFileSystemException: Não FileSystem para o esquema "s3"

Você pode encontrar a seguinte exceção no thread "principal":

```
org.apache.hadoop.fs.UnsupportedFileSystemException: No FileSystem for scheme "s3"
```

Se isso ocorrer, adicione as seguintes exceções à especificação SparkApplication:

```
hadoopConf:
   # EMRFS filesystem
   fs.s3.customAWSCredentialsProvider:
 com.amazonaws.auth.WebIdentityTokenCredentialsProvider
   fs.s3.impl: com.amazon.ws.emr.hadoop.fs.EmrFileSystem
   fs.AbstractFileSystem.s3.impl: org.apache.hadoop.fs.s3.EMRFSDelegate
   fs.s3.buffer.dir: /mnt/s3
   fs.s3.getObject.initialSocketTimeoutMilliseconds: "2000"
   mapreduce.fileoutputcommitter.algorithm.version.emr_internal_use_only.EmrFileSystem:
 "2"
   mapreduce.fileoutputcommitter.cleanup-
failures.ignored.emr_internal_use_only.EmrFileSystem: "true"
 sparkConf:
   # Required for EMR Runtime
   spark.driver.extraClassPath: /usr/lib/hadoop-lzo/lib/*:/usr/lib/hadoop/hadoop-
aws.jar:/usr/share/aws/aws-java-sdk/*:/usr/share/aws/emrfs/conf:/usr/share/aws/
emr/emrfs/lib/*:/usr/share/aws/emr/emrfs/auxlib/*:/usr/share/aws/emr/security/conf:/
usr/share/aws/emr/security/lib/*:/usr/share/aws/hmclient/lib/aws-glue-datacatalog-
spark-client.jar:/usr/share/java/Hive-JSON-Serde/hive-openx-serde.jar:/usr/share/aws/
sagemaker-spark-sdk/lib/sagemaker-spark-sdk.jar:/home/hadoop/extrajars/*
   spark.driver.extraLibraryPath: /usr/lib/hadoop/lib/native:/usr/lib/hadoop-lzo/lib/
native:/docker/usr/lib/hadoop/lib/native:/docker/usr/lib/hadoop-lzo/lib/native
   spark.executor.extraClassPath: /usr/lib/hadoop-lzo/lib/*:/usr/lib/hadoop/hadoop-
aws.jar:/usr/share/aws/aws-java-sdk/*:/usr/share/aws/emrfs/conf:/usr/share/aws/
emr/emrfs/lib/*:/usr/share/aws/emr/emrfs/auxlib/*:/usr/share/aws/emr/security/conf:/
usr/share/aws/emr/security/lib/*:/usr/share/aws/hmclient/lib/aws-glue-datacatalog-
```

spark-client.jar:/usr/share/java/Hive-JSON-Serde/hive-openx-serde.jar:/usr/share/aws/
sagemaker-spark-sdk/lib/sagemaker-spark-sdk.jar:/home/hadoop/extrajars/*
 spark.executor.extraLibraryPath: /usr/lib/hadoop/lib/native:/usr/lib/hadoop-lzo/lib/
native:/docker/usr/lib/hadoop/lib/native

Cotas e endpoints de serviço do Amazon EMR no EKS

A seguir são apresentados os endpoints de serviço e as cotas de serviço para o Amazon EMR no EKS. Para se conectar programaticamente a um AWS serviço, você usa um endpoint. Além dos AWS endpoints padrão, alguns AWS serviços oferecem endpoints FIPS em regiões selecionadas. Para obter mais informações, consulte Endpoints de serviço da AWS. As cotas de serviço, também chamadas de limites, correspondem ao número máximo de recursos ou operações de serviço para sua conta da AWS . Para obter mais informações, consulte as Service Quotas do AWS.

Service endpoints

Região da AWS nome	Código	Endpoint	Protocolo
Leste dos EUA (Norte da Virgínia)	us-east-1	<pre>emr-containers.us- east-1.amazonaws.com</pre>	HTTPS
Leste dos EUA (Ohio)	us-east-2	<pre>emr-containers.us- east-2.amazonaws.com</pre>	HTTPS
Oeste dos EUA (Norte da Califórnia)	us-west-1	emr-containers.us- west-1.amazonaws.com	HTTPS
Oeste dos EUA (Oregon)	us-west-2	<pre>emr-containers.us- west-2.amazonaws.com</pre>	HTTPS
Ásia-Pacífico (Tóquio)	ap-northe ast-1	<pre>emr-containers.ap- northeast-1.amazonaws.com</pre>	HTTPS
Ásia-Pacífico (Seul)	ap-northe ast-2	<pre>emr-containers.ap- northeast-2.amazonaws.com</pre>	HTTPS
Asia Pacific (Osaka)	ap-northe ast-3	<pre>emr-containers.ap- northeast-3.amazonaws.com</pre>	HTTPS
Ásia-Pacífico (Mumbai)	ap-south-1	<pre>emr-containers.ap- south-1.amazonaws.com</pre>	HTTPS

Service endpoints 417

Região da AWS nome	Código	Endpoint	Protocolo
Ásia-Pacífico (Hyderabad)	ap-south-2	<pre>emr-containers.ap- south-2.amazonaws.com</pre>	HTTPS
Ásia-Pacífico (Singapura)	ap-southe ast-1	<pre>emr-containers.ap- southeast-1.amazonaws.com</pre>	HTTPS
Ásia-Pacífico (Sydney)	ap-southe ast-2	<pre>emr-containers.ap- southeast-2.amazonaws.com</pre>	HTTPS
Ásia-Pacífico (Jacarta)	ap-southe ast-3	<pre>emr-containers.ap- southeast-3.amazonaws.com</pre>	HTTPS
Ásia-Pacífico (Hong Kong)	ap-east-1	<pre>emr-containers.ap- east-1.amazonaws.com</pre>	HTTPS
África (Cidade do Cabo)	af-south-1	<pre>emr-containers.af- south-1.amazonaws.com</pre>	HTTPS
Canadá (Central)	ca-central-1	emr-containers.ca-central-1.amazonaws.com	HTTPS
China (Ningxia)	cn-northw est-1	<pre>emr-containers.cn- northwest-1.amazon aws.com.cn</pre>	HTTPS
China (Pequim)	cn-north-1	emr-containers.cn- north-1.amazonaws.com.cn	HTTPS
Europa (Frankfurt)	eu-central-1	emr-containers.eu- central-1.amazonaws.com	HTTPS
Europa (Zurique)	eu-central-2	emr-containers.eu- central-2.amazonaws.com	HTTPS
Europa (Irlanda)	eu-west-1	emr-containers.eu- west-1.amazonaws.com	HTTPS

Service endpoints 418

Região da AWS nome	Código	Endpoint	Protocolo
Europa (Londres)	eu-west-2	<pre>emr-containers.eu- west-2.amazonaws.com</pre>	HTTPS
Europa (Paris)	eu-west-3	emr-containers.eu- west-3.amazonaws.com	HTTPS
Europa (Estocolmo)	eu-north-1	emr-containers.eu- north-1.amazonaws.com	HTTPS
Europa (Milão)	eu-south-1	emr-containers.eu- south-1.amazonaws.com	HTTPS
Europa (Espanha)	eu-south-2	<pre>emr-containers.eu- south-2.amazonaws.com</pre>	HTTPS
Israel (Tel Aviv)	il-central-1	<pre>emr-containers.il- central-1.amazonaws.com</pre>	HTTPS
América do Sul (São Paulo)	sa-east-1	<pre>emr-containers.sa- east-1.amazonaws.com</pre>	HTTPS
Oriente Médio (Emirados Árabes Unidos)	me-central-1	<pre>emr-containers.me- central-1.amazonaws.com</pre>	HTTPS
Oriente Médio (Bahrein)	me-south-1	emr-containers.me- south-1.amazonaws.com	HTTPS
AWS GovCloud (Leste dos EUA)	us-gov-east-1	<pre>emr-containers.us-gov- east-1.amazonaws.com</pre>	HTTPS
AWS GovCloud (Oeste dos EUA)	us-gov-we st-1	<pre>emr-containers.us-gov- west-1.amazonaws.com</pre>	HTTPS

Service endpoints 419

Cotas de serviço

O Amazon EMR no EKS limita as seguintes solicitações de API para cada AWS conta por região. Para obter mais informações sobre como a limitação é aplicada, consulte Limitação de solicitações de API na Amazon API Reference. EC2 Você pode solicitar um aumento nas cotas de limitação de API para sua AWS conta seguindo o guia a seguir.

Ação da API	Capacidade máxima do bucket	Taxa de reabastecimento do bucket (por segundo)
CancelJobRun	25	1
CreateJobTemplate	25	1
CreateManagedEndpoint	25	1
CreateSecurityConfiguration	25	1
CreateVirtualCluster	25	1
DeleteJobTemplate	25	1
DeleteManagedEndpoint	25	1
DeleteVirtualCluster	25	1
DescribeJobRun	100	20
DescribeJobTemplate	25	1
DescribeManagedEndpoint	100	5
DescribeSecurityConfiguration	25	1
DescribeVirtualCluster	100	5
GetManagedEndpoint SessionCredentials	25	1
ListJobRuns	100	5

Cotas de serviço 420

Ação da API	Capacidade máxima do bucket	Taxa de reabastecimento do bucket (por segundo)
ListJobTemplates	25	1
ListManagedEndpoints	25	1
ListSecurityConfigurations	25	1
ListVirtualClusters	100	5
StartJobRun	25	1
Throttle quota for all EMR on EKS API requests	200	20

Quando você cria sua AWS conta, definimos cotas padrão (também chamadas de limites) para seus AWS recursos por região. Se você tentar exceder a cota de um recurso, a solicitação falhará. Se isso acontecer, reduza o uso de recursos ou solicite um aumento de cota.

O console Service Quotas é um local central onde você pode visualizar e gerenciar suas cotas de AWS serviços e solicitar aumentos de cota para muitos dos recursos que você usa. Use as informações de cota fornecidas aqui para gerenciar sua AWS infraestrutura. Planeje a solicitação de aumentos das cotas com antecedência antes que sejam necessários.

Visualize cotas e solicite aumentos de cotas

Você pode visualizar suas cotas de serviço atuais para cada região usando o console Service Quotas. Para obter instruções detalhadas, consulte <u>Visualização de cotas de serviço</u> no Guia do usuário de cotas de serviço.

Você pode solicitar um aumento de cota no AWS Management Console ou usando a AWS CLI. As etapas estão detalhadas em Solicitando um aumento de cota.

Versões do Amazon EMR no EKS

Uma versão do Amazon EMR corresponde a um conjunto de aplicações de código aberto do ecossistema de big data. Cada versão inclui diferentes aplicações, componentes e recursos de big data que você seleciona para que o Amazon EMR no EKS implante e configure ao executar seu trabalho.

A partir das versões 5.32.0 e 6.2.0 do Amazon EMR, é possível implantar o Amazon EMR no EKS. Essa opção de implantação não está disponível em versões de liberação anteriores do Amazon EMR. Você deve especificar uma versão de liberação com suporte ao enviar seu trabalho.

O Amazon EMR no EKS usa o seguinte formato de rótulo de versão: emr-x.x.x-latest ou emrx.x.x-yyyymmdd com uma data de versão específica. Por exemplo, emr-7.9.0-latest ou emr-7.9.0-20210129. Ao usar o sufixo -latest, você garante que sua versão do Amazon EMR sempre inclua as atualizações de segurança mais recentes.



Note

Para uma comparação entre o Amazon EMR no EKS e o Amazon EMR em execução, EC2 consulte o Amazon EMR no site. FAQs AWS

Tópicos

- Amazon EMR nas versões EKS 7.9.0
- Amazon EMR nas versões EKS 7.8.0
- Amazon EMR nas versões EKS 7.7.0
- Amazon EMR nas versões do EKS 7.6.0
- Amazon EMR nas versões do EKS 7.5.0
- Amazon EMR nas versões do EKS 7.4.0
- Versões 7.3.0 do Amazon EMR no EKS
- Versões 7.2.0 do Amazon EMR no EKS
- Versões 7.1.0 do Amazon EMR no EKS
- Versões 7.0.0 do Amazon EMR no EKS
- Versões 6.15.0 do Amazon EMR no EKS

- Versões 6.14.0 do Amazon EMR no EKS
- Versões 6.13.0 do Amazon EMR no EKS
- Versões 6.12.0 do Amazon EMR no EKS
- Versões 6.11.0 do Amazon EMR no EKS
- Versões 6.10.0 do Amazon EMR no EKS
- Versões 6.9.0 do Amazon EMR no EKS
- Versões 6.8.0 do Amazon EMR no EKS
- Versões 6.7.0 do Amazon EMR no EKS
- Versões 6.6.0 do Amazon EMR no EKS
- Versões 6.5.0 do Amazon EMR no EKS
- Versões 6.4.0 do Amazon EMR no EKS
- Versões 6.3.0 do Amazon EMR no EKS
- Versões 6.2.0 do Amazon EMR no EKS
- Versões 5.36.0 do Amazon EMR no EKS
- Versões 5.35.0 do Amazon EMR no EKS
- Versões 5.34.0 do Amazon EMR no EKS
- Versões 5.33.0 do Amazon EMR no EKS
- Versões 5.32.0 do Amazon EMR no EKS

Amazon EMR nas versões EKS 7.9.0

Esta página descreve a funcionalidade nova e atualizada do Amazon EMR que é específica para a implantação do Amazon EMR no EKS. Para obter detalhes sobre o Amazon EMR em execução na Amazon EC2 e sobre a versão do Amazon EMR 7.9.0 em geral, consulte Amazon EMR 7.9.0 no Guia de lançamento do Amazon EMR.

Amazon EMR nas versões do EKS 7.9

As seguintes versões do Amazon EMR 7.9.0 estão disponíveis para o Amazon EMR no EKS. Selecione uma versão específica do EMR-7.9.0-xxxx para ver mais detalhes, como a tag de imagem do contêiner relacionada.

Versões 7.9.0 423

Flink releases

As seguintes versões do Amazon EMR 7.9.0 estão disponíveis para o Amazon EMR no EKS quando você executa aplicativos Flink.

- emr-7.9.0-flink-latest
- emr-7.9.0-flink-20250425

Spark releases

As seguintes versões do Amazon EMR 7.9.0 estão disponíveis para o Amazon EMR no EKS quando você executa aplicativos Spark.

- emr-7.9.0-mais recente
- emr-7.9.0-20250425
- emr-7.9.0-spark-rapids-latest
- emr-7.9.0-spark-rapids-20250425
- emr-7.9.0-java11-latest
- emr-7.9.0-java11-20250425
- emr-7.9.0-java8-latest
- emr-7.9.0-java8-20250425
- emr-7.9.0-spark-rapids-java8-latest
- emr-7.9.0-spark-rapids-java8-20250425
- notebook-spark/emr-7.9.0-latest
- notebook-spark/emr-7.9.0-20250425
- notebook-spark/emr-7.9.0-spark-rapids-latest
- notebook-spark/emr-7.9.0-spark-rapids-20250425
- notebook-spark/emr-7.9.0-java11-latest
- notebook-spark/emr-7.9.0-java11-20250425
- notebook-spark/emr-7.9.0-java8-latest
- notebook-spark/emr-7.9.0-java8-20250425
- notebook-spark/emr-7.9.0-spark-rapids-java8-latest

Versões 424

- notebook-spark/emr-7.9.0-spark-rapids-java8-20250425
- notebook-python/emr-7.9.0-latest
- notebook-python/emr-7.9.0-20250425
- notebook-python/emr-7.9.0-spark-rapids-latest
- notebook-python/emr-7.9.0-spark-rapids-20250425
- notebook-python/emr-7.9.0-java11-latest
- notebook-python/emr-7.9.0-java11-20250425
- notebook-python/emr-7.9.0-java8-latest
- notebook-python/emr-7.9.0-java8-20250425
- notebook-python/emr-7.9.0-spark-rapids-java8-latest
- notebook-python/emr-7.9.0-spark-rapids-java8-20250425
- livy/emr-7.9.0-latest
- livy/emr-7.9.0-20250425
- livy/emr-7.9.0-java11-latest
- livy/emr-7.9.0-java11-20250425
- livy/emr-7.9.0-java8-latest
- livy/emr-7.9.0-java8-20250425

Notas da versão

Notas de lançamento do Amazon EMR no EKS 7.9.0

- Aplicações compatíveis: AWS SDK para Java 2.31.16 and 1.12.782, Apache Spark 3.5.5, Apache Hudi 0.15.0-amzn-6, Apache Iceberg 1.7.1-amzn-2, Delta 3.3.0-amzn-1, Apache Spark RAPIDS 25.02.1-amzn-0, Jupyter Enterprise Gateway 2.6.0, Apache Flink 1.20.0-amzn-3, Flink Operator 1.10.0-amzn-3
- Componentes suportados emr-ddbemr-goodies,emr-s3-select,,emrfs,hadoopclient,hudi,hudi-spark,iceberg,spark-kubernetes.
- Classificações de configuração com suporte

Para uso com StartJobRune CreateManagedEndpoint APIs:

Classificações	Descrições
core-site	Altera os valores no arquivo core-site .xml do Hadoop.
emrfs-site	Alterar as configurações do EMRFS.
spark-metrics	Altera os valores no arquivo metrics.p roperties do Spark.
spark-defaults	Altera os valores no arquivo spark-def aults.conf do Spark.
spark-env	Alterar os valores no ambiente do Spark.
spark-hive-site	Altera os valores no arquivo hive-site .xml do Spark.
spark-log4j2	Altera os valores no arquivo log4j2.pr operties do Spark.
emr-job-submitter	Configuração para o pod de envio de trabalho.

Para uso específico com <u>CreateManagedEndpoint</u> APIs:

Classificações	Descrições
jeg-config	Altera os valores no arquivo jupyter_e nterprise_gateway_config.py do Jupyter Enterprise Gateway.
jupyter-kernel-overrides	Altera o valor da imagem do kernel no arquivo de um kernel do Jupyter especificado.

As classificações de configuração permitem que você personalize aplicações. Elas geralmente correspondem a um arquivo XML de configuração da aplicação, como spark-hive-site.xml. Para obter mais informações, consulte Configure Applications.

Alterações

As seguintes alterações estão incluídas na versão 7.9.0 do Amazon EMR no EKS:

Nenhuma alteração no lançamento.

emr-7.9.0-mais recente

Notas de versão: a versão emr-7.9.0-latest direciona para emr-7.9.0-20250425, no momento.

Regiões: emr-7.9.0-latest está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-7.9.0:latest.

emr-7.9.0-20250425

Notas de lançamento: emr-7.9.0-20250425 foi lançado em fevereiro de 2025. Esta é a versão inicial do Amazon EMR 7.9.0 (Spark).

Regiões: emr-emr-7.9.0-20250425 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-7.9.0-20250425.

emr-7.9.0-flink-latest

Notas de lançamento: emr-7.9.0-flink-latest atualmente aponta para emr-7.9.0-flink-20250425

Regiões: emr-7.9.0-flink-latest está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Alterações 427

Etiqueta da imagem do contêiner: emr-7.9.0-flink:latest.

emr-7.9.0-flink-20250425

Notas de lançamento: 7.9.0-flink-20250425 foi lançado em fevereiro de 2025. Esta é a versão inicial do Amazon EMR 7.9.0 (Flink).

Regiões: emr-7.9.0-flink-20250425 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-7.9.0-flink:20250425.

Amazon EMR nas versões EKS 7.8.0

Esta página descreve a funcionalidade nova e atualizada do Amazon EMR que é específica para a implantação do Amazon EMR no EKS. Para obter detalhes sobre o Amazon EMR em execução na Amazon EC2 e sobre a versão do Amazon EMR 7.8.0 em geral, consulte Amazon EMR 7.8.0 no Guia de lançamento do Amazon EMR.

Amazon EMR nas versões do EKS 7.8

As seguintes versões do Amazon EMR 7.8.0 estão disponíveis para o Amazon EMR no EKS. Selecione uma versão específica do EMR-7.8.0-xxxx para ver mais detalhes, como a tag de imagem do contêiner relacionada.

Flink releases

As seguintes versões do Amazon EMR 7.8.0 estão disponíveis para o Amazon EMR no EKS quando você executa aplicativos Flink.

- emr-7.8.0-flink-latest
- emr-7.8.0-flink-20250228

Spark releases

As seguintes versões do Amazon EMR 7.8.0 estão disponíveis para o Amazon EMR no EKS quando você executa aplicativos Spark.

emr-7.8.0-mais recente

emr-7.9.0-flink-20250425 428

- emr-7.8.0-20250228
- emr-7.8.0-spark-rapids-latest
- emr-7.8.0-spark-rapids-20250228
- emr-7.8.0-java11-latest
- emr-7.8.0-java11-20250228
- emr-7.8.0-java8-latest
- emr-7.8.0-java8-20250228
- emr-7.8.0-spark-rapids-java8-latest
- emr-7.8.0-spark-rapids-java8-20250228
- notebook-spark/emr-7.8.0-latest
- notebook-spark/emr-7.8.0-20250228
- notebook-spark/emr-7.8.0-spark-rapids-latest
- notebook-spark/emr-7.8.0-spark-rapids-20250228
- notebook-spark/emr-7.8.0-java11-latest
- notebook-spark/emr-7.8.0-java11-20250228
- notebook-spark/emr-7.8.0-java8-latest
- notebook-spark/emr-7.8.0-java8-20250228
- notebook-spark/emr-7.8.0-spark-rapids-java8-latest
- notebook-spark/emr-7.8.0-spark-rapids-java8-20250228
- notebook-python/emr-7.8.0-latest
- notebook-python/emr-7.8.0-20250228
- notebook-python/emr-7.8.0-spark-rapids-latest
- notebook-python/emr-7.8.0-spark-rapids-20250228
- notebook-python/emr-7.8.0-java11-latest
- notebook-python/emr-7.8.0-java11-20250228
- notebook-python/emr-7.8.0-java8-latest
- notebook-python/emr-7.8.0-java8-20250228
- notebook-python/emr-7.8.0-spark-rapids-java8-latest
- notebook-python/emr-7.8.0-spark-rapids-java8-20250228
- livy/emr-7.8.0-latest

Versões 429

- livy/emr-7.8.0-20250228
- livy/emr-7.8.0-java11-latest
- livy/emr-7.8.0-java11-20250228
- livy/emr-7.8.0-java8-latest
- livy/emr-7.8.0-java8-20250228

Notas da versão

Notas de lançamento do Amazon EMR no EKS 7.8.0

- Aplicações compatíveis: AWS SDK para Java 2.29.52 and 1.12.780, Apache Spark 3.5.4, Apache Hudi 0.15.0-amzn-5, Apache Iceberg 1.7.1-amzn-1, Delta 3.3.0-amzn-0, Apache Spark RAPIDS 24.12.0-amzn-0, Jupyter Enterprise Gateway 2.6.0, Apache Flink 1.20.0-amzn-2, Flink Operator 1.10.0-amzn-2
- Componentes suportados emr-ddbemr-goodies,emr-s3-select,,emrfs,hadoopclient,hudi,hudi-spark,iceberg,spark-kubernetes.
- · Classificações de configuração com suporte

Para uso com StartJobRune CreateManagedEndpoint APIs:

Classificações	Descrições
core-site	Altera os valores no arquivo core-site .xml do Hadoop.
emrfs-site	Alterar as configurações do EMRFS.
spark-metrics	Altera os valores no arquivo metrics.p roperties do Spark.
spark-defaults	Altera os valores no arquivo spark-def aults.conf do Spark.
spark-env	Alterar os valores no ambiente do Spark.
spark-hive-site	Altera os valores no arquivo hive-site .xml do Spark.

Classificações	Descrições
spark-log4j2	Altera os valores no arquivo log4j2.pr operties do Spark.
emr-job-submitter	Configuração para o pod de envio de trabalho.

Para uso específico com CreateManagedEndpoint APIs:

Classificações	Descrições
jeg-config	Altera os valores no arquivo jupyter_e nterprise_gateway_config.py do Jupyter Enterprise Gateway.
jupyter-kernel-overrides	Altera o valor da imagem do kernel no arquivo de um kernel do Jupyter especificado.

As classificações de configuração permitem que você personalize aplicações. Elas geralmente correspondem a um arquivo XML de configuração da aplicação, como spark-hive-site.xml. Para obter mais informações, consulte <u>Configure Applications</u>.

Alterações

As seguintes alterações estão incluídas na versão 7.8.0 do Amazon EMR no EKS:

- · Recursos nativos do FGAC, incluindo:
 - Suporte do Iceberg para executar trabalhos que executam ações em tabelas de formação sem lagos em um cluster virtual de controle de acesso refinado (FGAC). (Há uma alternativa para o IAM.)
 - Suporte de mesa S3
- Conexão Spark

Alterações 431

emr-7.8.0-mais recente

Notas de versão: a versão emr-7.8.0-latest direciona para emr-7.8.0-20250228, no momento.

Regiões: emr-7.8.0-latest está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-7.8.0:latest.

emr-7.8.0-20250228

Notas de lançamento: emr-7.8.0-20250228 foi lançado em fevereiro de 2025. Esta é a versão inicial do Amazon EMR 7.8.0 (Spark).

Regiões: emr-emr-7.8.0-20250228 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-7.8.0-20250228.

emr-7.8.0-flink-latest

Notas de lançamento: emr-7.8.0-flink-latest atualmente aponta para emr-7.8.0-flink-20250228

Regiões: emr-7.8.0-flink-latest está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-7.8.0-flink:latest.

emr-7.8.0-flink-20250228

Notas de lançamento: 7.8.0-flink-20250228 foi lançado em fevereiro de 2025. Esta é a versão inicial do Amazon EMR 7.8.0 (Flink).

Regiões: emr-7.8.0-flink-20250228 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-7.8.0-flink:20250228.

emr-7.8.0-mais recente 432

Amazon EMR nas versões EKS 7.7.0

Esta página descreve a funcionalidade nova e atualizada do Amazon EMR que é específica para a implantação do Amazon EMR no EKS. Para obter detalhes sobre o Amazon EMR em execução na Amazon EC2 e sobre a versão do Amazon EMR 7.7.0 em geral, consulte Amazon EMR 7.7.0 no Guia de lançamento do Amazon EMR.

Amazon EMR nas versões do EKS 7.7

As seguintes versões do Amazon EMR 7.7.0 estão disponíveis para o Amazon EMR no EKS. Selecione uma versão específica do EMR-7.7.0-xxxx para ver mais detalhes, como a tag de imagem do contêiner relacionada.

Flink releases

As seguintes versões do Amazon EMR 7.7.0 estão disponíveis para o Amazon EMR no EKS quando você executa aplicativos Flink.

- emr-7.7.0-flink-latest
- emr-7.7.0-flink-20250131

Spark releases

As seguintes versões do Amazon EMR 7.7.0 estão disponíveis para o Amazon EMR no EKS quando você executa aplicativos Spark.

- emr-7.7.0-mais recente
- emr-7.7.0-20250131
- emr-7.7.0-spark-rapids-latest
- emr-7.7.0-spark-rapids-20250131
- emr-7.7.0-java11-latest
- emr-7.7.0-java11-20250131
- emr-7.7.0-java8-latest
- emr-7.7.0-java8-20250131
- emr-7.7.0-spark-rapids-java8-latest
- emr-7.7.0-spark-rapids-java8-20250131

Versões 7.7.0 433

- notebook-spark/emr-7.7.0-latest
- notebook-spark/emr-7.7.0-20250131
- notebook-spark/emr-7.7.0-spark-rapids-latest
- notebook-spark/emr-7.7.0-spark-rapids-20250131
- notebook-spark/emr-7.7.0-java11-latest
- notebook-spark/emr-7.7.0-java11-20250131
- notebook-spark/emr-7.7.0-java8-latest
- notebook-spark/emr-7.7.0-java8-20250131
- notebook-spark/emr-7.7.0-spark-rapids-java8-latest
- notebook-spark/emr-7.7.0-spark-rapids-java8-20250131
- notebook-python/emr-7.7.0-latest
- notebook-python/emr-7.7.0-20250131
- notebook-python/emr-7.7.0-spark-rapids-latest
- notebook-python/emr-7.7.0-spark-rapids-20250131
- notebook-python/emr-7.7.0-java11-latest
- notebook-python/emr-7.7.0-java11-20250131
- notebook-python/emr-7.7.0-java8-latest
- notebook-python/emr-7.7.0-java8-20250131
- notebook-python/emr-7.7.0-spark-rapids-java8-latest
- notebook-python/emr-7.7.0-spark-rapids-java8-20250131
- livy/emr-7.7.0-latest
- livy/emr-7.7.0-20250131
- livy/emr-7.7.0-java11-latest
- livy/emr-7.7.0-java11-20250131
- livy/emr-7.7.0-java8-latest
- livy/emr-7.7.0-java8-20250131

Notas da versão

Notas de lançamento do Amazon EMR no EKS 7.7.0

- Aplicações compatíveis: AWS SDK para Java 2.29.25 and 1.12.779, Apache Spark 3.5.3-amzn-0, Apache Hudi 0.15.0-amzn-3, Apache Iceberg 1.6.1-amzn-2, Delta 3.2.1-amzn-1, Apache Spark RAPIDS 24.10.1-amzn-0, Jupyter Enterprise Gateway 2.6.0, Apache Flink 1.20.0-amzn-0, Flink Operator 1.10.0-amzn-0
- Componentes suportados emr-ddbemr-goodies,emr-s3-select,,emrfs,hadoop-client,hudi,hudi-spark,iceberg,spark-kubernetes.
- Classificações de configuração com suporte

Para uso com StartJobRune CreateManagedEndpoint APIs:

Classificações	Descrições
core-site	Altera os valores no arquivo core-site .xml do Hadoop.
emrfs-site	Alterar as configurações do EMRFS.
spark-metrics	Altera os valores no arquivo metrics.p roperties do Spark.
spark-defaults	Altera os valores no arquivo spark-def aults.conf do Spark.
spark-env	Alterar os valores no ambiente do Spark.
spark-hive-site	Altera os valores no arquivo hive-site .xml do Spark.
spark-log4j2	Altera os valores no arquivo log4j2.pr operties do Spark.
emr-job-submitter	Configuração para o pod de envio de trabalho.

Para uso específico com CreateManagedEndpoint APIs:

Classificações	Descrições
jeg-config	Altera os valores no arquivo jupyter_e nterprise_gateway_config.py do Jupyter Enterprise Gateway.
jupyter-kernel-overrides	Altera o valor da imagem do kernel no arquivo de um kernel do Jupyter especificado.

As classificações de configuração permitem que você personalize aplicações. Elas geralmente correspondem a um arquivo XML de configuração da aplicação, como spark-hive-site.xml. Para obter mais informações, consulte Configure Applications.

Alterações

As seguintes alterações estão incluídas na versão 7.7.0 do Amazon EMR no EKS:

 A versão Iceberg em uso a partir do EMR 7.7.0 não oferece mais suporte ao Java 8. Além disso, o Iceberg está excluído das seguintes imagens do Java 8: emr-7.7.0-java8-latest e. emr-7.7.0-spark-rapids-java8-latest

emr-7.7.0-mais recente

Notas de versão: a versão emr-7.7.0-latest direciona para emr-7.7.0-20250131, no momento.

Regiões: emr-7.7.0-latest está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-7.7.0:latest.

emr-7.7.0-20250131

Notas de lançamento: emr-7.7.0-20250131 foi lançado em fevereiro de 2025. Esta é a versão inicial do Amazon EMR 7.7.0 (Spark).

Alterações 436

Regiões: emr-emr-7.7.0-20250131 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-7.7.0-20250131.

emr-7.7.0-flink-latest

Notas de lançamento: emr-7.7.0-flink-latest atualmente aponta para emr-7.7.0-flink-20250131

Regiões: emr-7.7.0-flink-latest está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-7.7.0-flink:latest.

emr-7.7.0-flink-20250131

Notas de lançamento: 7.7.0-flink-20250131 foi lançado em fevereiro de 2025. Esta é a versão inicial do Amazon EMR 7.7.0 (Flink).

Regiões: emr-7.7.0-flink-20250131 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-7.7.0-flink:20250131.

Amazon EMR nas versões do EKS 7.6.0

Esta página descreve a funcionalidade nova e atualizada do Amazon EMR que é específica para a implantação do Amazon EMR no EKS. Para obter detalhes sobre o Amazon EMR em execução na Amazon EC2 e sobre a versão do Amazon EMR 7.6.0 em geral, consulte Amazon EMR 7.6.0 no Guia de lançamento do Amazon EMR.

Amazon EMR nas versões do EKS 7.6

As seguintes versões do Amazon EMR 7.6.0 estão disponíveis para o Amazon EMR no EKS. Selecione uma versão específica do EMR-7.6.0-xxxx para ver mais detalhes, como a tag de imagem do contêiner relacionada.

emr-7.7.0-flink-latest 437

Flink releases

As seguintes versões do Amazon EMR 7.6.0 estão disponíveis para o Amazon EMR no EKS quando você executa aplicativos Flink.

- emr-7.6.0-flink-latest
- emr-7.6.0-flink-20241213

Spark releases

As seguintes versões do Amazon EMR 7.6.0 estão disponíveis para o Amazon EMR no EKS quando você executa aplicativos Spark.

- emr-7.6.0-mais recente
- emr-7.6.0-20241213
- emr-7.6.0-spark-rapids-latest
- emr-7.6.0-spark-rapids-20241213
- emr-7.6.0-java11-latest
- emr-7.6.0-java11-20241213
- emr-7.6.0-java8-latest
- emr-7.6.0-java8-20241213
- emr-7.6.0-spark-rapids-java8-latest
- emr-7.6.0-spark-rapids-java8-20241213
- notebook-spark/emr-7.6.0-latest
- notebook-spark/emr-7.6.0-20241213
- notebook-spark/emr-7.6.0-spark-rapids-latest
- notebook-spark/emr-7.6.0-spark-rapids-20241213
- notebook-spark/emr-7.6.0-java11-latest
- notebook-spark/emr-7.6.0-java11-20241213
- notebook-spark/emr-7.6.0-java8-latest
- notebook-spark/emr-7.6.0-java8-20241213
- notebook-spark/emr-7.6.0-spark-rapids-java8-latest

Versões 438

- notebook-spark/emr-7.6.0-spark-rapids-java8-20241213
- notebook-python/emr-7.6.0-latest
- notebook-python/emr-7.6.0-20241213
- notebook-python/emr-7.6.0-spark-rapids-latest
- notebook-python/emr-7.6.0-spark-rapids-20241213
- notebook-python/emr-7.6.0-java11-latest
- notebook-python/emr-7.6.0-java11-20241213
- notebook-python/emr-7.6.0-java8-latest
- notebook-python/emr-7.6.0-java8-20241213
- notebook-python/emr-7.6.0-spark-rapids-java8-latest
- notebook-python/emr-7.6.0-spark-rapids-java8-20241213
- livy/emr-7.6.0-latest
- livy/emr-7.6.0-20241213
- livy/emr-7.6.0-java11-latest
- livy/emr-7.6.0-java11-20241213
- livy/emr-7.6.0-java8-latest
- livy/emr-7.6.0-java8-20241213

Notas da versão

Notas de lançamento do Amazon EMR no EKS 7.6.0

- Aplicações compatíveis: AWS SDK para Java 2.29.25 and 1.12.779, Apache Spark 3.5.3-amzn-0, Apache Hudi 0.15.0-amzn-3, Apache Iceberg 1.6.1-amzn-2, Delta 3.2.1-amzn-1, Apache Spark RAPIDS 24.10.1-amzn-0, Jupyter Enterprise Gateway 2.6.0, Apache Flink 1.20.0-amzn-0, Flink Operator 1.10.0-amzn-0
- Componentes com suporte: aws-sagemaker-spark-sdk, emr-ddb, emr-goodies, emr-s3-select, emrfs, hadoop-client, hudi, hudi-spark, iceberg e spark-kubernetes.
- Classificações de configuração com suporte

Para uso com StartJobRune CreateManagedEndpoint APIs:

Classificações	Descrições
core-site	Altera os valores no arquivo core-site .xml do Hadoop.
emrfs-site	Alterar as configurações do EMRFS.
spark-metrics	Altera os valores no arquivo metrics.p roperties do Spark.
spark-defaults	Altera os valores no arquivo spark-def aults.conf do Spark.
spark-env	Alterar os valores no ambiente do Spark.
spark-hive-site	Altera os valores no arquivo hive-site .xml do Spark.
spark-log4j2	Altera os valores no arquivo log4j2.pr operties do Spark.
emr-job-submitter	Configuração para o pod de envio de trabalho.

Para uso específico com <u>CreateManagedEndpoint</u> APIs:

Classificações	Descrições
jeg-config	Altera os valores no arquivo jupyter_e nterprise_gateway_config.py do Jupyter Enterprise Gateway.
jupyter-kernel-overrides	Altera o valor da imagem do kernel no arquivo de um kernel do Jupyter especificado.

As classificações de configuração permitem que você personalize aplicações. Elas geralmente correspondem a um arquivo XML de configuração da aplicação, como spark-hive-site.xml. Para obter mais informações, consulte Configure Applications.

Recursos notáveis

Os seguintes recursos estão incluídos na versão 7.6.0 do Amazon EMR no EKS:

• Suporte de configuração de monitoramento para o operador do Apache Spark — A configuração de monitoramento permite que você configure facilmente o arquivamento de registros do seu aplicativo Spark e dos registros do operador no Amazon S3 ou na Amazon. CloudWatch Você pode escolher um ou ambos. Isso adiciona um sidecar do agente de log aos pods de operador, driver e executor do Spark e, posteriormente, encaminha os registros desses componentes para os coletores configurados. Para obter mais informações, consulte Como usar a configuração de monitoramento para monitorar o operador do Spark Kubernetes e as tarefas do Spark.

Alterações

As seguintes alterações estão incluídas na versão 7.6.0 do Amazon EMR no EKS:

Nenhuma alteração no lançamento.

emr-7.6.0-mais recente

Notas de versão: a versão emr-7.6.0-latest direciona para emr-7.6.0-20241213, no momento.

Regiões: emr-7.6.0-latest está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-7.6.0:latest.

emr-7.6.0-20241213

Notas de lançamento: 7.6.0-20241213 foi lançado em janeiro de 2024. Esta é a versão inicial do Amazon EMR 7.6.0 (Spark).

Atributos 441

Regiões: emr-7.6.0-20241213 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-7.6.0:20241213.

emr-7.6.0-flink-latest

Notas de lançamento: emr-7.6.0-flink-latest atualmente aponta para emr-7.6.0-flink-20241213

Regiões: emr-7.6.0-flink-latest está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-7.6.0-flink:latest.

emr-7.6.0-flink-20241213

Notas de lançamento: 7.6.0-flink-20241213 foi lançado em janeiro de 2024. Esta é a versão inicial do Amazon EMR 7.6.0 (Flink).

Regiões: emr-7.6.0-flink-20241213 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-7.6.0-flink:20241213.

Amazon EMR nas versões do EKS 7.5.0

Esta página descreve a funcionalidade nova e atualizada do Amazon EMR que é específica para a implantação do Amazon EMR no EKS. Para obter detalhes sobre o Amazon EMR em execução na Amazon EC2 e sobre a versão do Amazon EMR 7.5.0 em geral, consulte Amazon EMR <u>7.5.0 no Guia de lançamento do Amazon EMR</u>.

Amazon EMR nas versões do EKS 7.5

As seguintes versões do Amazon EMR 7.5.0 estão disponíveis para o Amazon EMR no EKS. Selecione uma versão específica do EMR-7.5.0-xxxx para ver mais detalhes, como a tag de imagem do contêiner relacionada.

emr-7.6.0-flink-latest 442

Notas da versão

Notas de lançamento do Amazon EMR no EKS 7.5.0

- Aplicações compatíveis: AWS SDK para Java 2.28.8 and 1.12.772, Apache Spark 3.5.2-amzn-1, Apache Hudi 0.15.0-amzn-1, Apache Iceberg 1.6.1-amzn-0, Delta 3.2.0-amzn-1, Apache Spark RAPIDS 24.08.1-amzn-1, Jupyter Enterprise Gateway 2.6.0, Apache Flink 1.19.1-amzn-1, Flink Operator 1.9.0-amzn-0
- Componentes com suporte: aws-sagemaker-spark-sdk, emr-ddb, emr-goodies, emr-s3-select, emrfs, hadoop-client, hudi, hudi-spark, iceberg e spark-kubernetes.

Amazon EMR nas versões do EKS 7.4.0

Esta página descreve a funcionalidade nova e atualizada do Amazon EMR que é específica para a implantação do Amazon EMR no EKS. Para obter detalhes sobre o Amazon EMR em execução na Amazon EC2 e sobre a versão do Amazon EMR 7.4.0 em geral, consulte Amazon EMR <u>7.4.0 no Guia de lançamento do Amazon EMR</u>.

Amazon EMR nas versões do EKS 7.4

As seguintes versões do Amazon EMR 7.4.0 estão disponíveis para o Amazon EMR no EKS. Selecione uma versão específica do EMR-7.4.0-xxxx para ver mais detalhes, como a tag de imagem do contêiner relacionada.

Notas da versão

Notas de lançamento do Amazon EMR no EKS 7.4.0

- Aplicações compatíveis: AWS SDK para Java 2.25.70 and 1.12.772, Apache Spark 3.5.2-amzn-0, Apache Hudi 0.15.0-amzn-1, Apache Iceberg 1.6.1-amzn-0, Delta 3.2.0-amzn-1, Apache Spark RAPIDS 24.08.1-amzn-0, Jupyter Enterprise Gateway 2.6.0, Apache Flink 1.19.1-amzn-0, Flink Operator 1.9.0-amzn-1
- Componentes com suporte: aws-sagemaker-spark-sdk, emr-ddb, emr-goodies, emr-s3-select, emrfs, hadoop-client, hudi, hudi-spark, iceberg e spark-kubernetes.

Versões 7.3.0 do Amazon EMR no EKS

Esta página descreve a funcionalidade nova e atualizada do Amazon EMR que é específica para a implantação do Amazon EMR no EKS. Para obter detalhes sobre o Amazon EMR em execução na Amazon EC2 e sobre a versão do Amazon EMR 7.3.0 em geral, consulte Amazon EMR <u>7.3.0 no</u> Guia de lançamento do Amazon EMR.

Versões 7.3 do Amazon EMR no EKS

As versões 7.3.0 do Amazon EMR a seguir estão disponíveis para o Amazon EMR no EKS. Selecione uma versão específica do emr-7.3.0-XXXX para exibir mais detalhes, como a tag de imagem do contêiner relacionada.

Flink releases

As versões 7.3.0 do Amazon EMR a seguir estão disponíveis para o Amazon EMR no EKS ao executar aplicações do Flink.

- emr-7.3.0-flink-latest
- emr-7.3.0-flink-29240920

Spark releases

As versões 7.3.0 do Amazon EMR a seguir estão disponíveis para o Amazon EMR no EKS ao executar aplicações do Spark.

- emr-7.3.0-latest
- emr-7.3.0-20240920
- emr-7.3.0-spark-rapids-latest
- emr-7.3.0-spark-rapids-29240920
- emr-7.3.0-java11-latest
- emr-7.3.0-java11-29240920
- emr-7.3.0-java8-latest
- emr-7.3.0-java8-29240920
- emr-7.3.0-spark-rapids-java8-latest
- emr-7.3.0-spark-rapids-java8-29240920

Versões 7.3.0 444

- notebook-spark/emr-7.3.0-latest
- notebook-spark/emr-7.3.0-20240920
- notebook-spark/emr-7.3.0-spark-rapids-latest
- notebook-spark/emr-7.3.0-spark-rapids-29240920
- notebook-spark/emr-7.3.0-java11-latest
- notebook-spark/emr-7.3.0-java11-29240920
- notebook-spark/emr-7.3.0-java8-latest
- notebook-spark/emr-7.3.0-java8-29240920
- notebook-spark/emr-7.3.0-spark-rapids-java8-latest
- notebook-spark/emr-7.3.0-spark-rapids-java8-29240920
- notebook-python/emr-7.3.0-latest
- notebook-python/emr-7.3.0-20240920
- notebook-python/emr-7.3.0-spark-rapids-latest
- notebook-python/emr-7.3.0-spark-rapids-29240920
- notebook-python/emr-7.3.0-java11-latest
- notebook-python/emr-7.3.0-java11-29240920
- notebook-python/emr-7.3.0-java8-latest
- notebook-python/emr-7.3.0-java8-29240920
- notebook-python/emr-7.3.0-spark-rapids-java8-latest
- notebook-python/emr-7.3.0-spark-rapids-java8-29240920
- livy/emr-7.3.0-latest
- livy/emr-7.3.0-20240920
- livy/emr-7.3.0-java11-latest
- livy/emr-7.3.0-java11-29240920
- livy/emr-7.3.0-java8-latest
- livy/emr-7.3.0-java8-29240920

Notas da versão

Notas da versão 7.3.0 do Amazon EMR no EKS

- Aplicações compatíveis: AWS SDK para Java 2.25.70 and 1.12.747, Apache Spark 3.5.1-amzn-1, Apache Hudi 0.15.0-amzn-0, Apache Iceberg 1.5.2-amzn-0, Delta 3.2.0-amzn-0, Apache Spark RAPIDS 24.06.1-amzn-0, Jupyter Enterprise Gateway 2.6.0, Apache Flink 1.18.1-amzn-2, Flink Operator 1.9.0-amzn-0
- Componentes com suporte: aws-sagemaker-spark-sdk, emr-ddb, emr-goodies, emr-s3-select, emrfs, hadoop-client, hudi, hudi-spark, iceberg e spark-kubernetes.
- Classificações de configuração com suporte

Para uso com StartJobRune CreateManagedEndpoint APIs:

Classificações	Descrições
core-site	Altera os valores no arquivo core-site .xml do Hadoop.
emrfs-site	Alterar as configurações do EMRFS.
spark-metrics	Altera os valores no arquivo metrics.p roperties do Spark.
spark-defaults	Altera os valores no arquivo spark-def aults.conf do Spark.
spark-env	Alterar os valores no ambiente do Spark.
spark-hive-site	Altera os valores no arquivo hive-site .xml do Spark.
spark-log4j2	Altera os valores no arquivo log4j2.pr operties do Spark.
emr-job-submitter	Configuração para o pod de envio de trabalho.

Para uso específico com CreateManagedEndpoint APIs:

Classificações	Descrições
jeg-config	Altera os valores no arquivo jupyter_e nterprise_gateway_config.py do Jupyter Enterprise Gateway.
jupyter-kernel-overrides	Altera o valor da imagem do kernel no arquivo de um kernel do Jupyter especificado.

As classificações de configuração permitem que você personalize aplicações. Elas geralmente correspondem a um arquivo XML de configuração da aplicação, como spark-hive-site.xml. Para obter mais informações, consulte Configure Applications.

Recursos notáveis

Os recursos a seguir estão incluídos na versão 7.3.0 do Amazon EMR no EKS.

- Atualizações de aplicações: o Amazon EMR no EKS agora inclui o <u>Flink Operator 1.9.0</u>. Além de outros recursos, o Kubernetes do Flink agora permite definir cotas de CPU e memória para o escalador automático.
- Suporte do Apache Iceberg para o Apache Flink: o Apache Iceberg é um formato de tabelas analíticas enormes de código aberto e alta performance. Desde o Amazon EMR 7.3.0, você pode usar as tabelas do Apache Iceberg ao executar o Apache Flink no Amazon EMR no EKS. Para obter mais informações, consulte o artigo <u>Using Apache Iceberg with Amazon EMR on EKS</u> do Amazon EMR no EKS.
- Suporte do Delta Lake para Apache Flink: o Delta Lake é uma estrutura de camada de armazenamento para arquiteturas de lakehouse comumente desenvolvidas no Amazon S3. Com as versões 7.3.0 e superiores do Amazon EMR, você pode usar as tabelas do Delta ao executar o Apache Flink no Amazon EMR no EKS. Para obter mais informações, consulte <u>Using Delta Lake</u> <u>with Amazon EMR on EKS</u>.

Alterações

As alterações a seguir estão incluídas na versão 7.3.0 do Amazon EMR no EKS.

Atributos 447

 Com o Amazon EMR no EKS 7.3.0 e superior, o Apache Flink agora usa o runtime do Java 17 por padrão.

emr-7.3.0-latest

Notas de versão: a versão emr-7.3.0-latest direciona para emr-7.3.0-20240920, no momento.

Regiões: emr-7.3.0-latest está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-7.3.0:latest.

emr-7.3.0-20240920

Notas da versão: 7.3.0-20240920 foi lançada em dezembro de 2023. Esta é a versão inicial do Amazon EMR 7.3.0 (Spark).

Regiões: emr-7.3.0-20240920 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-7.3.0:29240920.

emr-7.3.0-flink-latest

Notas de versão: a versão emr-7.3.0-flink-latest direciona para emr-7.3.0-flink-29240920, no momento.

Regiões: emr-7.3.0-flink-latest está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-7.3.0-flink:latest.

emr-7.3.0-flink-29240920

Notas da versão: 7.3.0-flink-29240920 foi lançada em dezembro de 2023. Esta é a versão inicial do Amazon EMR 7.3.0 (Flink).

Regiões: emr-7.3.0-flink-29240920 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

emr-7.3.0-latest 448

Etiqueta da imagem do contêiner: emr-7.3.0-flink:29240920.

Versões 7.2.0 do Amazon EMR no EKS

Esta página descreve a funcionalidade nova e atualizada do Amazon EMR que é específica para a implantação do Amazon EMR no EKS. Para obter detalhes sobre o Amazon EMR em execução na Amazon EC2 e sobre a versão do Amazon EMR 7.2.0 em geral, consulte Amazon EMR <u>7.2.0 no Guia de lançamento do Amazon EMR</u>.

Versões 7.2 do Amazon EMR no EKS

As versões 7.2.0 do Amazon EMR a seguir estão disponíveis para o Amazon EMR no EKS. Selecione uma versão específica do emr-7.2.0-XXXX para exibir mais detalhes, como a tag de imagem do contêiner relacionada.

Flink releases

As versões 7.2.0 do Amazon EMR a seguir estão disponíveis para o Amazon EMR no EKS ao executar aplicações do Flink.

- emr-7.2.0-flink-latest
- emr-7.2.0-flink-20240610

Spark releases

As versões 7.2.0 do Amazon EMR a seguir estão disponíveis para o Amazon EMR no EKS ao executar aplicações do Spark.

- emr-7.2.0-latest
- emr-7.2.0-20240610
- emr-7.2.0-spark-rapids-latest
- emr-7.2.0-spark-rapids-20240610
- emr-7.2.0-java11-latest
- emr-7.2.0-java11-20240610
- emr-7.2.0-java8-latest
- emr-7.2.0-java8-20240610
- emr-7.2.0-spark-rapids-java8-latest

Versões 7.2.0 449

- emr-7.2.0-spark-rapids-java8-20240610
- notebook-spark/emr-7.2.0-latest
- notebook-spark/emr-7.2.0-20240610
- notebook-spark/emr-7.2.0-spark-rapids-latest
- notebook-spark/emr-7.2.0-spark-rapids-20240610
- notebook-spark/emr-7.2.0-java11-latest
- notebook-spark/emr-7.2.0-java11-20240610
- notebook-spark/emr-7.2.0-java8-latest
- notebook-spark/emr-7.2.0-java8-20240610
- notebook-spark/emr-7.2.0-spark-rapids-java8-latest
- notebook-spark/emr-7.2.0-spark-rapids-java8-20240610
- notebook-python/emr-7.2.0-latest
- notebook-python/emr-7.2.0-20240610
- notebook-python/emr-7.2.0-spark-rapids-latest
- notebook-python/emr-7.2.0-spark-rapids-20240610
- notebook-python/emr-7.2.0-java11-latest
- notebook-python/emr-7.2.0-java11-20240610
- notebook-python/emr-7.2.0-java8-latest
- notebook-python/emr-7.2.0-java8-20240610
- notebook-python/emr-7.2.0-spark-rapids-java8-latest
- notebook-python/emr-7.2.0-spark-rapids-java8-20240610
- livy/emr-7.2.0-latest
- livy/emr-7.2.0-20240610
- livy/emr-7.2.0-java11-latest
- livy/emr-7.2.0-java11-20240610
- livy/emr-7.2.0-java8-latest
- livy/emr-7.2.0-java8-20240610

Notas da versão

Notas da versão 7.2.0 do Amazon EMR no EKS

- Aplicações compatíveis: AWS SDK para Java 2.23.18 and 1.12.705, Apache Spark 3.5.1-amzn-1, Apache Hudi 0.14.1-amzn-0, Apache Iceberg 1.5.0-amzn-0, Delta 3.1.0, Apache Spark RAPIDS 24.02.0-amzn-1, Jupyter Enterprise Gateway 2.6.0, Apache Flink 1.18.1-amzn-0, Flink Operator 1.8.0-amzn-1
- Componentes com suporte: aws-sagemaker-spark-sdk, emr-ddb, emr-goodies, emr-s3-select, emrfs, hadoop-client, hudi, hudi-spark, iceberg e spark-kubernetes.
- Classificações de configuração com suporte

Para uso com StartJobRune CreateManagedEndpoint APIs:

Classificações	Descrições
core-site	Altera os valores no arquivo core-site .xml do Hadoop.
emrfs-site	Alterar as configurações do EMRFS.
spark-metrics	Altera os valores no arquivo metrics.p roperties do Spark.
spark-defaults	Altera os valores no arquivo spark-def aults.conf do Spark.
spark-env	Alterar os valores no ambiente do Spark.
spark-hive-site	Altera os valores no arquivo hive-site .xml do Spark.
spark-log4j2	Altera os valores no arquivo log4j2.pr operties do Spark.
emr-job-submitter	Configuração para o pod de envio de trabalho.

Para uso específico com CreateManagedEndpoint APIs:

Classificações	Descrições
jeg-config	Altera os valores no arquivo jupyter_e nterprise_gateway_config.py do Jupyter Enterprise Gateway.
jupyter-kernel-overrides	Altera o valor da imagem do kernel no arquivo de um kernel do Jupyter especificado.

As classificações de configuração permitem que você personalize aplicações. Elas geralmente correspondem a um arquivo XML de configuração da aplicação, como spark-hive-site.xml. Para obter mais informações, consulte Configure Applications.

Recursos notáveis

Os recursos a seguir estão incluídos na versão 7.2.0 do Amazon EMR no EKS.

- Atualizações da aplicação: as atualizações da aplicação do Amazon EMR no EKS 7.2.0 incluem Spark 3.5.1, Flink 1.18.1 e operador do Flink 1.8.0.
- Escalador automático para atualizações do Flink: a versão 7.2.0 usa a configuração de código aberto job.autoscaler.restart.time-tracking.enabled para permitir a estimativa do tempo de redimensionamento, para que você não precise mais atribuir valores empíricos manualmente ao reiniciar. Se você executar a versão 7.1.0 ou inferior, ainda poderá usar o ajuste de escala automático do Amazon EMR.
- Integração do Apache Hudi com Apache Flink no Amazon EMR no EKS: essa versão adiciona uma integração entre o Apache Hudi e o Apache Flink, para que você possa usar o operador do Kubernetes do Flink ao executar trabalhos do Hudi. O Hudi permite usar operações em nível de registro que podem simplificar o gerenciamento de dados e o desenvolvimento do pipeline de dados.
- Integração do Amazon S3 Express One Zone com o Amazon EMR no EKS: com a versão 7.2.0 e superior, você pode fazer upload de dados no S3 Express One Zone com o Amazon EMR no EKS.
 O S3 Express One Zone é uma classe de armazenamento de zona única e alta performance do Amazon S3 que oferece acesso consistente a dados de milissegundos de um dígito para a maioria das aplicações sensíveis à latência. Na hora da execução, o S3 Express One Zone oferece o

Atributos 452

armazenamento de objetos na nuvem com a menor latência e a maior performance do Amazon S3.

Suporte para configurações padrão no operador do Spark: o operador do Spark no Amazon EKS
agora oferece suporte às mesmas configurações padrão do modelo de início de execução de
trabalhos no Amazon EMR no EKS para a versão 7.2.0 e posteriores. Isso significa que recursos
como Amazon S3 e EMRFS não exigem mais configurações manuais no arquivo YAML.

emr-7.2.0-latest

Notas de versão: a versão emr-7.2.0-latest direciona para emr-7.2.0-20240610, no momento.

Regiões: emr-7.2.0-latest está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-7.2.0:latest.

emr-7.2.0-20240610

Notas da versão: 7.2.0-20240610 foi lançada em dezembro de 2023. Esta é a versão inicial do Amazon EMR 7.2.0 (Spark).

Regiões: emr-7.2.0-20240610 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-7.2.0:20240610.

emr-7.2.0-flink-latest

Notas de versão: a versão emr-7.2.0-flink-latest direciona para emr-7.2.0-flink-20240610, no momento.

Regiões: emr-7.2.0-flink-latest está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-7.2.0-flink:latest.

emr-7.2.0-flink-20240610

Notas da versão: 7.2.0-flink-20240610 foi lançada em dezembro de 2023. Esta é a versão inicial do Amazon EMR 7.2.0 (Flink).

emr-7.2.0-latest 453

Regiões: emr-7.2.0-flink-20240610 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-7.2.0-flink:20240610.

Versões 7.1.0 do Amazon EMR no EKS

Esta página descreve a funcionalidade nova e atualizada do Amazon EMR que é específica para a implantação do Amazon EMR no EKS. Para obter detalhes sobre o Amazon EMR em execução na Amazon EC2 e sobre a versão do Amazon EMR 7.1.0 em geral, consulte Amazon EMR <u>7.1.0 no</u> Guia de lançamento do Amazon EMR.

Versões 7.1 do Amazon EMR no EKS

As versões 7.1.0 do Amazon EMR a seguir estão disponíveis para o Amazon EMR no EKS. Selecione uma versão específica do emr-7.1.0-XXXX para exibir mais detalhes, como a tag de imagem do contêiner relacionada.

Flink releases

As versões 7.1.0 do Amazon EMR a seguir estão disponíveis para o Amazon EMR no EKS ao executar aplicações do Flink.

- emr-7.1.0-flink-latest
- emr-7.1.0-flink-20240321

Spark releases

As versões 7.1.0 do Amazon EMR a seguir estão disponíveis para o Amazon EMR no EKS ao executar aplicações do Spark.

- emr-7.1.0-latest
- emr-7.1.0-20240321
- emr-7.1.0-spark-rapids-latest
- emr-7.1.0-spark-rapids-20240321
- emr-7.1.0-java11-latest
- emr-7.1.0-java11-20240321

Versões 7.1.0 454

- emr-7.1.0-java8-latest
- emr-7.1.0-java8-20240321
- emr-7.1.0-spark-rapids-java8-latest
- emr-7.1.0-spark-rapids-java8-20240321
- notebook-spark/emr-7.1.0-latest
- notebook-spark/emr-7.1.0-20240321
- notebook-spark/emr-7.1.0-spark-rapids-latest
- notebook-spark/emr-7.1.0-spark-rapids-20240321
- notebook-spark/emr-7.1.0-java11-latest
- notebook-spark/emr-7.1.0-java11-20240321
- notebook-spark/emr-7.1.0-java8-latest
- notebook-spark/emr-7.1.0-java8-20240321
- notebook-spark/emr-7.1.0-spark-rapids-java8-latest
- notebook-spark/emr-7.1.0-spark-rapids-java8-20240321
- notebook-python/emr-7.1.0-latest
- notebook-python/emr-7.1.0-20240321
- notebook-python/emr-7.1.0-spark-rapids-latest
- notebook-python/emr-7.1.0-spark-rapids-20240321
- notebook-python/emr-7.1.0-java11-latest
- notebook-python/emr-7.1.0-java11-20240321
- notebook-python/emr-7.1.0-java8-latest
- notebook-python/emr-7.1.0-java8-20240321
- notebook-python/emr-7.1.0-spark-rapids-java8-latest
- notebook-python/emr-7.1.0-spark-rapids-java8-20240321
- livy/emr-7.1.0-latest
- livy/emr-7.1.0-20240321
- livy/emr-7.1.0-java11-latest
- livy/emr-7.1.0-java11-20240321
- livy/emr-7.1.0-java8-latest
- livy/emr-7.1.0-java8-20240321

Versões 455

Notas da versão

Notas da versão 7.1.0 do Amazon EMR no EKS

- Aplicações compatíveis: AWS SDK para Java 2.23.18 and 1.12.656, Apache Spark 3.5.0-amzn-1, Apache Hudi 0.14.1-amzn-0, Apache Iceberg 1.4.3-amzn-0, Delta 3.0.0, Apache Spark RAPIDS 23.10.0-amzn-1, Jupyter Enterprise Gateway 2.6.0, Apache Flink 1.18.1-amzn-0, Flink Operator 1.6.1-amzn-1
- Componentes com suporte: aws-sagemaker-spark-sdk, emr-ddb, emr-goodies, emr-s3-select, emrfs, hadoop-client, hudi, hudi-spark, iceberg e spark-kubernetes.
- Classificações de configuração com suporte

Para uso com StartJobRune CreateManagedEndpoint APIs:

Classificações	Descrições
core-site	Altera os valores no arquivo core-site .xml do Hadoop.
emrfs-site	Alterar as configurações do EMRFS.
spark-metrics	Altera os valores no arquivo metrics.p roperties do Spark.
spark-defaults	Altera os valores no arquivo spark-def aults.conf do Spark.
spark-env	Alterar os valores no ambiente do Spark.
spark-hive-site	Altera os valores no arquivo hive-site .xml do Spark.
spark-log4j2	Altera os valores no arquivo log4j2.pr operties do Spark.
emr-job-submitter	Configuração para o pod de envio de trabalho.

Para uso específico com CreateManagedEndpoint APIs:

Classificações	Descrições
jeg-config	Altera os valores no arquivo jupyter_e nterprise_gateway_config.py do Jupyter Enterprise Gateway.
jupyter-kernel-overrides	Altera o valor da imagem do kernel no arquivo de um kernel do Jupyter especificado.

As classificações de configuração permitem que você personalize aplicações. Elas geralmente correspondem a um arquivo XML de configuração da aplicação, como spark-hive-site.xml. Para obter mais informações, consulte Configure Applications.

Recursos notáveis

Os recursos a seguir estão incluídos na versão 7.1.0 do Amazon EMR no EKS.

• Suporte do Apache Livy para Amazon EMR no EKS: com o Amazon EMR no EKS nas versões 7.1.0 e superiores, você pode usar o Apache Livy em um cluster do Amazon EKS para criar uma interface REST do Apache Livy e enviar trabalhos do Spark ou trechos de código do Spark. Isso permite recuperar os resultados de forma síncrona e assíncrona, enquanto ainda aproveita os benefícios do Amazon EMR no EKS, como o runtime do Spark otimizado para Amazon EMR, endpoints do Livy habilitados por SSL e uma experiência de configuração programática.

emr-7.1.0-latest

Notas de versão: a versão emr-7.1.0-latest direciona para emr-7.1.0-20240321, no momento.

Regiões: emr-7.1.0-latest está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte <u>Endpoints de serviço do Amazon EMR no EKS</u>.

Etiqueta da imagem do contêiner: emr-7.1.0:latest.

Atributos 457

emr-7.1.0-20240321

Notas da versão: 7.1.0-20240321 foi lançada em dezembro de 2023. Esta é a versão inicial do Amazon EMR 7.1.0 (Spark).

Regiões: emr-7.1.0-20240321 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-7.1.0:20240321.

emr-7.1.0-flink-latest

Notas de versão: a versão emr-7.1.0-flink-latest direciona para emr-7.1.0-flink-20240321, no momento.

Regiões: emr-7.1.0-flink-latest está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-7.1.0-flink:latest.

emr-7.1.0-flink-20240321

Notas da versão: 7.1.0-flink-20240321 foi lançada em dezembro de 2023. Esta é a versão inicial do Amazon EMR 7.1.0 (Flink).

Regiões: emr-7.1.0-flink-20240321 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-7.1.0-flink:20240321.

Versões 7.0.0 do Amazon EMR no EKS

Esta página descreve a funcionalidade nova e atualizada do Amazon EMR que é específica para a implantação do Amazon EMR no EKS. Para obter detalhes sobre o Amazon EMR em execução na Amazon EC2 e sobre a versão do Amazon EMR 7.0.0 em geral, consulte Amazon EMR <u>7.0.0 no</u> Guia de lançamento do Amazon EMR.

emr-7.1.0-20240321 458

Versões 7.0 do Amazon EMR no EKS

As versões 7.0.0 do Amazon EMR apresentadas a seguir estão disponíveis para o Amazon EMR no EKS. Selecione uma versão específica do emr-7.0.0-XXXX para visualizar mais detalhes, como a tag de imagem do contêiner relacionada.

Flink releases

As versões 7.0.0 do Amazon EMR apresentadas a seguir estão disponíveis para o Amazon EMR no EKS ao executar aplicações do Flink.

- emr-7.0.0-flink-latest
- emr-7.0.0-flink-2024321
- emr-7.0.0-flink-20231211

Spark releases

As versões 7.0.0 do Amazon EMR apresentadas a seguir estão disponíveis para o Amazon EMR no EKS ao executar aplicações do Spark.

- emr-7.0.0-latest
- emr-7.0.0-20231211
- emr-7.0.0-spark-rapids-latest
- emr-7.0.0-spark-rapids-20231211
- emr-7.0.0-java11-latest
- emr-7.0.0-java11-20231211
- emr-7.0.0-java8-latest
- emr-7.0.0-java8-20231211
- emr-7.0.0-spark-rapids-java8-latest
- emr-7.0.0-spark-rapids-java8-20231211
- notebook-spark/emr-7.0.0-latest
- notebook-spark/emr-7.0.0-20231211
- notebook-spark/emr-7.0.0-spark-rapids-latest
- notebook-spark/emr-7.0.0-spark-rapids-20231211
- notebook-spark/emr-7.0.0-java11-latest

Versões 459

- notebook-spark/emr-7.0.0-java11-20231211
- notebook-spark/emr-7.0.0-java8-latest
- notebook-spark/emr-7.0.0-java8-20231211
- notebook-spark/emr-7.0.0-spark-rapids-java8-latest
- notebook-spark/emr-7.0.0-spark-rapids-java8-20231211
- notebook-python/emr-7.0.0-latest
- notebook-python/emr-7.0.0-20231211
- notebook-python/emr-7.0.0-spark-rapids-latest
- notebook-python/emr-7.0.0-spark-rapids-20231211
- notebook-python/emr-7.0.0-java11-latest
- notebook-python/emr-7.0.0-java11-20231211
- notebook-python/emr-7.0.0-java8-latest
- notebook-python/emr-7.0.0-java8-20231211
- notebook-python/emr-7.0.0-spark-rapids-java8-latest
- notebook-python/emr-7.0.0-spark-rapids-java8-20231211

Notas da versão

Notas da versão 7.0.0 do Amazon EMR no EKS

- Aplicações compatíveis: AWS SDK para Java 2.20.160-amzn-0 and 1.12.595, Apache Spark 3.5.0-amzn-0, Apache Flink 1.18.0-amzn-0, Flink Operator 1.6.1, Apache Hudi 0.14.0-amzn-1, Apache Iceberg 1.4.2-amzn-0, Delta 3.0.0, Apache Spark RAPIDS 23.10.0-amzn-0, Jupyter Enterprise Gateway 2.6.0
- Componentes com suporte: aws-sagemaker-spark-sdk, emr-ddb, emr-goodies, emr-s3-select, emrfs, hadoop-client, hudi, hudi-spark, iceberg e spark-kubernetes.
- Classificações de configuração com suporte

Para uso com StartJobRune CreateManagedEndpoint APIs:

Classificações	Descrições
core-site	Altera os valores no arquivo core-site .xml do Hadoop.
emrfs-site	Alterar as configurações do EMRFS.
spark-metrics	Altera os valores no arquivo metrics.p roperties do Spark.
spark-defaults	Altera os valores no arquivo spark-def aults.conf do Spark.
spark-env	Alterar os valores no ambiente do Spark.
spark-hive-site	Altera os valores no arquivo hive-site .xml do Spark.
spark-log4j	Altera os valores no arquivo log4j2.pr operties do Spark.
emr-job-submitter	Configuração para o pod de envio de trabalho.

Para uso específico com <u>CreateManagedEndpoint</u> APIs:

Classificações	Descrições
jeg-config	Altera os valores no arquivo jupyter_e nterprise_gateway_config.py do Jupyter Enterprise Gateway.
jupyter-kernel-overrides	Altera o valor da imagem do kernel no arquivo de um kernel do Jupyter especificado.

As classificações de configuração permitem que você personalize aplicações. Elas geralmente correspondem a um arquivo XML de configuração da aplicação, como spark-hive-site.xml. Para obter mais informações, consulte Configure Applications.

Recursos notáveis

Os atributos apresentados a seguir estão inclusos na versão 7.0 do Amazon EMR no EKS.

- Atualizações da aplicação: as atualizações da aplicação Amazon EMR no EKS 7.0.0 incluem Spark 3.5, Flink 1.18 e Flink Operator 1.6.1.
- Ajuste automático de parâmetros do escalador automático do Flink: os parâmetros padrão que o escalador automático do Flink usa nos cálculos de ajuste de escala podem não ser o valor ideal para um determinado trabalho. A versão 7.0.0 do Amazon EMR no EKS usa tendências do histórico de métricas específicas capturadas para calcular o parâmetro ideal personalizado para o trabalho.

Alterações

A alterações a seguir estão inclusas na versão 7.0 do Amazon EMR no EKS.

- Amazon Linux 2023: com as versões 7.0.0 e superiores do Amazon EMR no EKS, todas as imagens de contêiner são baseadas no Amazon Linux 2023.
- O Spark usa o Java 17 como runtime padrão: o Spark da versão 7.0.0 do Amazon EMR no EKS usa o Java 17 como runtime padrão. Se precisar, você pode alternar para usar o Java 8 ou o Java 11 com o rótulo de versão correspondente, conforme fornecido na lista Versões 7.0 do Amazon
 EMR no EKS.

emr-7.0.0-latest

Notas de versão: a versão emr-7.0.0-latest direciona para emr-7.0.0-2024321, no momento.

Regiões: emr-7.0.0-latest está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-7.0.0:latest.

Atributos 462

emr-7.0.0-2024321

Notas de versão: 7.0.0-2024321 foi lançada em 11 de março de 2024. Em comparação com a versão anterior, esta versão foi atualizada com os pacotes do Amazon Linux atualizados recentemente e com as correções críticas.

Regiões: emr-7.0.0-2024321 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-7.0.0:2024321.

emr-7.0.0-20231211

Notas da versão: 7.0.0-20231211 foi lançada em dezembro de 2023. Esta é a versão inicial do Amazon EMR 7.0.0 (Spark).

Regiões: emr-7.0.0-20231211 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-7.0.0:20231211.

emr-7.0.0-flink-latest

Notas de versão: a versão emr-7.0.0-flink-latest direciona para emr-7.0.0-flink-2024321, no momento.

Regiões: emr-7.0.0-flink-latest está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-7.0.0-flink:latest.

emr-7.0.0-flink-2024321

Notas de versão: 7.0.0-flink-2024321 foi lançada em 11 de março de 2024. Em comparação com a versão anterior, esta versão foi atualizada com os pacotes do Amazon Linux atualizados recentemente e com as correções críticas.

Regiões: emr-7.0.0-flink-2024321 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

emr-7.0.0-2024321 463

Etiqueta da imagem do contêiner: emr-7.0.0-flink:2024321.

emr-7.0.0-flink-20231211

Notas da versão: 7.0.0-flink-20231211 foi lançada em dezembro de 2023. Esta é a versão inicial do Amazon EMR 7.0.0 (Flink).

Regiões: emr-7.0.0-flink-20231211 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-7.0.0-flink:20231211.

Versões 6.15.0 do Amazon EMR no EKS

Esta página descreve a funcionalidade nova e atualizada do Amazon EMR que é específica para a implantação do Amazon EMR no EKS. Para obter detalhes sobre o Amazon EMR em execução na Amazon EC2 e sobre a versão 6.15.0 do Amazon EMR em geral, consulte Amazon EMR 6.15.0 no Guia de lançamento do Amazon EMR.

Versões 6.15 do Amazon EMR no EKS

As versões 6.15.0 do Amazon EMR apresentadas a seguir estão disponíveis para o Amazon EMR no EKS. Selecione uma versão específica do emr-6.15.0-XXXX para visualizar mais detalhes, como a tag de imagem do contêiner relacionada.

Flink releases

As versões 6.15.0 do Amazon EMR apresentadas a seguir estão disponíveis para o Amazon EMR no EKS ao executar aplicações do Flink.

- emr-6.15.0-flink-latest
- emr-6.15.0-flink-20240105
- emr-6.15.0-flink-20231109

Spark releases

As versões 6.15.0 do Amazon EMR apresentadas a seguir estão disponíveis para o Amazon EMR no EKS ao executar aplicações do Spark.

emr-7.0.0-flink-20231211 464

- emr-6.15.0-latest
- emr-6.15.0-20231109
- emr-6.15.0-spark-rapids-latest
- emr-6.15.0-spark-rapids-20231109
- emr-6.15.0-java11-latest
- emr-6.15.0-java11-20231109
- emr-6.15.0-java17-latest
- emr-6.15.0-java17-20231109
- emr-6.15.0-java17-al2023-latest
- emr-6.15.0-java17-al2023-20231109
- emr-6.15.0-spark-rapids-java17-latest
- emr-6.15.0-spark-rapids-java17-20231109
- emr-6.15.0-spark-rapids-java17-al2023-latest
- emr-6.15.0-spark-rapids-java17-al2023-20231109
- notebook-spark/emr-6.15.0-latest
- notebook-spark/emr-6.15.0-20231109
- notebook-spark/emr-6.15.0-spark-rapids-latest
- notebook-spark/emr-6.15.0-spark-rapids-20231109
- notebook-spark/emr-6.15.0-java11-latest
- notebook-spark/emr-6.15.0-java11-20231109
- notebook-spark/emr-6.15.0-java17-latest
- notebook-spark/emr-6.15.0-java17-20231109
- notebook-spark/emr-6.15.0-java17-al2023-latest
- notebook-spark/emr-6.15.0-java17-al2023-20231109
- notebook-python/emr-6.15.0-latest
- notebook-python/emr-6.15.0-20231109
- notebook-python/emr-6.15.0-spark-rapids-latest
- notebook-python/emr-6.15.0-spark-rapids-20231109
- notebook-python/emr-6.15.0-java11-latest

Versões 465

- notebook-python/emr-6.15.0-java11-20231109
- notebook-python/emr-6.15.0-java17-latest
- notebook-python/emr-6.15.0-java17-20231109
- notebook-python/emr-6.15.0-java17-al2023-latest
- notebook-python/emr-6.15.0-java17-al2023-20231109

Notas da versão

Notas da versão 6.15.0 do Amazon EMR no EKS

- Aplicações compatíveis: AWS SDK para Java 1.12.569, Apache Spark 3.4.1-amzn-2, Apache Flink 1.17.1-amzn-1, Apache Hudi 0.14.0-amzn-0, Apache Iceberg 1.4.0-amzn-0, Delta 2.4.0, Apache Spark RAPIDS 23.08.01-amzn-0, Jupyter Enterprise Gateway 2.6.0
- Componentes com suporte: aws-sagemaker-spark-sdk, emr-ddb, emr-goodies, emr-s3-select, emrfs, hadoop-client, hudi, hudi-spark, iceberg e spark-kubernetes.
- Classificações de configuração com suporte

Para uso com StartJobRune CreateManagedEndpoint APIs:

Classificações	Descrições
core-site	Altera os valores no arquivo core-site .xml do Hadoop.
emrfs-site	Alterar as configurações do EMRFS.
spark-metrics	Altera os valores no arquivo metrics.p roperties do Spark.
spark-defaults	Altera os valores no arquivo spark-def aults.conf do Spark.
spark-env	Alterar os valores no ambiente do Spark.
spark-hive-site	Altera os valores no arquivo hive-site .xml do Spark.

Classificações	Descrições
spark-log4j	Altera os valores no arquivo log4j2.pr operties do Spark.
emr-job-submitter	Configuração para o pod de envio de trabalho.

Para uso específico com CreateManagedEndpoint APIs:

Classificações	Descrições
jeg-config	Altera os valores no arquivo jupyter_e nterprise_gateway_config.py do Jupyter Enterprise Gateway.
jupyter-kernel-overrides	Altera o valor da imagem do kernel no arquivo de um kernel do Jupyter especificado.

As classificações de configuração permitem que você personalize aplicações. Elas geralmente correspondem a um arquivo XML de configuração da aplicação, como spark-hive-site.xml. Para obter mais informações, consulte <u>Configure Applications</u>.

Recursos notáveis

Os recursos apresentados a seguir estão inclusos na versão 6.15 do Amazon EMR no EKS.

• Amazon EMR no EKS com Apache Flink: com a versão 6.15.0 do Amazon EMR no EKS você pode executar aplicações baseadas no Apache Flink com outros tipos de aplicações no mesmo cluster do Amazon EKS. Isso ajuda a melhorar a utilização de recursos e a simplificar o gerenciamento da infraestrutura. Você pode aproveitar as instâncias spot em uma aplicação do Flink com uma desativação tranquila e obter tempos de reinicialização mais rápidos com recuperação refinada e recuperação local de tarefas com o Amazon EBS. Os recursos de acessibilidade e monitoramento incluem a capacidade de iniciar um aplicativo Flink com jars armazenados no Amazon S3, acesso ao AWS Glue Data Catalog, monitoramento da integração com o Amazon S3 e a CloudWatch Amazon e rotação de registros de contêineres.

Atributos 467

emr-6.15.0-latest

Notas de versão: a versão emr-6.15.0-latest direciona para emr-6.15.0-20240105, no momento.

Regiões: emr-6.15.0-latest está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-6.15.0:latest.

emr-6.15.0-20240105

Notas de versão: 6.15.0-20240105 foi lançada em 17 de janeiro de 2024. Em comparação com a versão anterior, esta versão foi atualizada com os pacotes do Amazon Linux atualizados recentemente e com as correções críticas.

Regiões: emr-6.15.0-20240105 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-6.15.0:20240105.

emr-6.15.0-20231109

Notas da versão: 6.15.0-20231109 foi lançada em 17 de novembro de 2023. Esta é a versão inicial do Amazon EMR 6.15.0.

Regiões: emr-6.15.0-20231109 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-6.15.0:20231109.

emr-6.15.0-flink-latest

Notas de versão: a versão emr-6.15.0-flink-latest direciona para emr-6.15.0-flink-20240105, no momento.

Regiões: emr-6.15.0-flink-latest está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-6.15.0-flink:latest.

emr-6.15.0-latest 468

emr-6.15.0-flink-20240105

Notas de versão: 6.15.0-flink-20240105 foi lançada em 17 de janeiro de 2024. Em comparação com a versão anterior, esta versão foi atualizada com os pacotes do Amazon Linux atualizados recentemente e com as correções críticas.

Regiões: emr-6.15.0-flink-20240105 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-6.15.0-flink:20240105.

emr-6.15.0-flink-20231109

Notas da versão: 6.15.0-flink-20231109 foi lançada em 17 de novembro de 2023. Esta é a versão inicial do Amazon EMR 6.15.0.

Regiões: emr-6.15.0-flink-20231109 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-6.15.0-flink:20231109.

Versões 6.14.0 do Amazon EMR no EKS

Esta página descreve a funcionalidade nova e atualizada do Amazon EMR que é específica para a implantação do Amazon EMR no EKS. Para obter detalhes sobre o Amazon EMR em execução na Amazon EC2 e sobre a versão do Amazon EMR 6.14.0 em geral, consulte Amazon EMR 6.14.0 no Guia de lançamento do Amazon EMR.

Versões 6.14 do Amazon EMR no EKS

As versões 6.14.0 do Amazon EMR apresentadas a seguir estão disponíveis para o Amazon EMR no EKS. Selecione uma versão específica do emr-6.14.0-XXXX para visualizar mais detalhes, como a etiqueta de imagem do contêiner relacionada.

- emr-6.14.0-latest
- emr-6.14.0-20231005

emr-6.15.0-flink-20240105 469

- emr-6.14.0-spark-rapids-latest
- emr-6.14.0-spark-rapids-20231005
- emr-6.14.0-java11-latest
- emr-6.14.0-java11-20231005
- emr-6.14.0-java17-latest
- emr-6.14.0-java17-20231005
- emr-6.14.0-java17-al2023-latest
- emr-6.14.0-java17-al2023-20231005
- emr-6.14.0-spark-rapids-java17-latest
- emr-6.14.0-spark-rapids-java17-20231005
- emr-6.14.0-spark-rapids-java17-al2023-latest
- emr-6.14.0-spark-rapids-java17-al2023-20231005
- notebook-spark/emr-6.14.0-latest
- notebook-spark/emr-6.14.0-20231005
- notebook-spark/emr-6.14.0-spark-rapids-latest
- notebook-spark/emr-6.14.0-spark-rapids-20231005
- notebook-spark/emr-6.14.0-java11-latest
- notebook-spark/emr-6.14.0-java11-20231005
- notebook-spark/emr-6.14.0-java17-latest
- notebook-spark/emr-6.14.0-java17-20231005
- notebook-spark/emr-6.14.0-java17-al2023-latest
- notebook-spark/emr-6.14.0-java17-al2023-20231005
- notebook-python/emr-6.14.0-latest
- notebook-python/emr-6.14.0-20231005
- notebook-python/emr-6.14.0-spark-rapids-latest
- notebook-python/emr-6.14.0-spark-rapids-20231005
- notebook-python/emr-6.14.0-java11-latest
- notebook-python/emr-6.14.0-java11-20231005
- notebook-python/emr-6.14.0-java17-latest
- notebook-python/emr-6.14.0-java17-20231005

Versões 470

- notebook-python/emr-6.14.0-java17-al2023-latest
- notebook-python/emr-6.14.0-java17-al2023-20231005

Notas da versão

Notas da versão 6.14.0 do Amazon EMR no EKS

- Aplicativos compatíveis AWS SDK para Java 1.12.543, Apache Spark 3.4.1-amzn-1, Apache Hudi 0.13.1-amzn-2, Apache Iceberg 1.3.0-amzn-0, Delta 2.4.0, Apache Spark RAPIDS 23.06.0-amzn-2, Jupyter Enterprise Gateway 2.7.0
- Componentes com suporte: aws-sagemaker-spark-sdk, emr-ddb, emr-goodies, emr-s3-select, emrfs, hadoop-client, hudi, hudi-spark, iceberg e spark-kubernetes.
- Classificações de configuração com suporte

Para uso com StartJobRune CreateManagedEndpoint APIs:

Classificações	Descrições
core-site	Altera os valores no arquivo core-site .xml do Hadoop.
emrfs-site	Alterar as configurações do EMRFS.
spark-metrics	Altera os valores no arquivo metrics.p roperties do Spark.
spark-defaults	Altera os valores no arquivo spark-def aults.conf do Spark.
spark-env	Alterar os valores no ambiente do Spark.
spark-hive-site	Altera os valores no arquivo hive-site .xml do Spark.
spark-log4j	Altera os valores no arquivo log4j2.pr operties do Spark.
emr-job-submitter	Configuração para o pod de envio de trabalho.

Para uso específico com CreateManagedEndpoint APIs:

Classificações	Descrições
jeg-config	Altera os valores no arquivo jupyter_e nterprise_gateway_config.py do Jupyter Enterprise Gateway.
jupyter-kernel-overrides	Altera o valor da imagem do kernel no arquivo de um kernel do Jupyter especificado.

As classificações de configuração permitem que você personalize aplicações. Elas geralmente correspondem a um arquivo XML de configuração da aplicação, como spark-hive-site.xml. Para obter mais informações, consulte <u>Configure Applications</u>.

Recursos notáveis

Os atributos apresentados a seguir estão inclusos na versão 6.14 do Amazon EMR no EKS.

 Suporte ao <u>Apache Livy</u>: o Amazon EMR no EKS agora oferece suporte ao Apache Livy com spark-submit.

emr-6.14.0-latest

Notas de versão: a versão emr-6.14.0-latest direciona para emr-6.14.0-20231005, no momento.

Regiões: emr-6.14.0-latest está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-6.14.0:latest.

emr-6.14.0-20231005

Notas de versão: 6.14.0-20231005 foi liberada em 17 de outubro de 2023. Esta é a versão inicial do Amazon EMR 6.14.0.

Atributos 472

Regiões: emr-6.14.0-20231005 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-6.14.0:20231005.

Versões 6.13.0 do Amazon EMR no EKS

Esta página descreve a funcionalidade nova e atualizada do Amazon EMR que é específica para a implantação do Amazon EMR no EKS. Para obter detalhes sobre o Amazon EMR em execução na Amazon EC2 e sobre a versão do Amazon EMR 6.13.0 em geral, consulte Amazon EMR 6.13.0 no Guia de lançamento do Amazon EMR.

Versões 6.13 do Amazon EMR no EKS

As versões 6.13.0 do Amazon EMR apresentadas a seguir estão disponíveis para o Amazon EMR no EKS. Selecione uma versão específica do emr-6.13.0-XXXX para visualizar mais detalhes, como a etiqueta de imagem do contêiner relacionada.

- emr-6.13.0-latest
- emr-6.13.0-20230814
- emr-6.13.0-spark-rapids-latest
- emr-6.13.0-spark-rapids-20230814
- emr-6.13.0-java11-latest
- emr-6.13.0-java11-20230814
- emr-6.13.0-java17-latest
- emr-6.13.0-java17-20230814
- emr-6.13.0-java17-al2023-latest
- emr-6.13.0-java17-al2023-20230814
- emr-6.13.0-spark-rapids-java17-latest
- emr-6.13.0-spark-rapids-java17-20230814
- emr-6.13.0-spark-rapids-java17-al2023-latest
- emr-6.13.0-spark-rapids-java17-al2023-20230814
- notebook-spark/emr-6.13.0-latest

Versões 6.13.0 473

- notebook-spark/emr-6.13.0-20230814
- notebook-spark/emr-6.13.0-spark-rapids-latest
- notebook-spark/emr-6.13.0-spark-rapids-20230814
- notebook-spark/emr-6.13.0-java11-latest
- notebook-spark/emr-6.13.0-java11-20230814
- notebook-spark/emr-6.13.0-java17-latest
- notebook-spark/emr-6.13.0-java17-20230814
- notebook-spark/emr-6.13.0-java17-al2023-latest
- notebook-spark/emr-6.13.0-java17-al2023-20230814
- notebook-python/emr-6.13.0-latest
- notebook-python/emr-6.13.0-20230814
- notebook-python/emr-6.13.0-spark-rapids-latest
- notebook-python/emr-6.13.0-spark-rapids-20230814
- notebook-python/emr-6.13.0-java11-latest
- notebook-python/emr-6.13.0-java11-20230814
- notebook-python/emr-6.13.0-java17-latest
- notebook-python/emr-6.13.0-java17-20230814
- notebook-python/emr-6.13.0-java17-al2023-latest
- notebook-python/emr-6.13.0-java17-al2023-20230814

Notas da versão

Notas da versão 6.13.0 do Amazon EMR no EKS

- Aplicativos compatíveis AWS SDK para Java 1.12.513, Apache Spark 3.4.1-amzn-0, Apache Hudi 0.13.1-amzn-0, Apache Iceberg 1.3.0-amzn-0, Delta 2.4.0, Apache Spark RAPIDS 23.06.0-amzn-1, Jupyter Enterprise Gateway 2.6.0.amzn
- Componentes com suporte: aws-sagemaker-spark-sdk, emr-ddb, emr-goodies, emr-s3-select, emrfs, hadoop-client, hudi, hudi-spark, iceberg e spark-kubernetes.
- Classificações de configuração com suporte

Para uso com StartJobRune CreateManagedEndpoint APIs:

Classificações	Descrições
core-site	Altera os valores no arquivo core-site .xml do Hadoop.
emrfs-site	Alterar as configurações do EMRFS.
spark-metrics	Altera os valores no arquivo metrics.p roperties do Spark.
spark-defaults	Altera os valores no arquivo spark-def aults.conf do Spark.
spark-env	Alterar os valores no ambiente do Spark.
spark-hive-site	Altera os valores no arquivo hive-site .xml do Spark.
spark-log4j	Altera os valores no arquivo log4j2.pr operties do Spark.
emr-job-submitter	Configuração para o pod de envio de trabalho.

Para uso específico com <u>CreateManagedEndpoint</u> APIs:

Classificações	Descrições
jeg-config	Altera os valores no arquivo jupyter_e nterprise_gateway_config.py do Jupyter Enterprise Gateway.
jupyter-kernel-overrides	Altera o valor da imagem do kernel no arquivo de um kernel do Jupyter especificado.

As classificações de configuração permitem que você personalize aplicações. Elas geralmente correspondem a um arquivo XML de configuração da aplicação, como spark-hive-site.xml. Para obter mais informações, consulte Configure Applications.

Recursos notáveis

Os recursos apresentados a seguir estão inclusos na versão 6.13 do Amazon EMR no EKS.

- Amazon Linux 2023 Com o Amazon EMR no EKS 6.13 e superior, você pode iniciar o Spark com
 o AL2 023 como sistema operacional junto com o tempo de execução do Java 17. Para fazer isso,
 use o rótulo de versão al2023 em seu nome. Por exemplo: emr-6.13.0-java17-al2023latest. Recomendamos que você valide e execute testes de desempenho antes de mover suas
 cargas de trabalho de produção para AL2 023 e Java 17.
- Amazon EMR no EKS com o Apache Flink (pré-visualização pública): as versões 6.13 do Amazon EMR no EKS e versões superiores oferecem suporte ao Apache Flink, que está disponível em pré-visualização pública. Com este lançamento, você pode executar sua aplicação baseada no Apache Flink em conjunto com outros tipos de aplicações no mesmo cluster do Amazon EKS. Isso ajuda a melhorar a utilização de recursos e a simplificar o gerenciamento da infraestrutura. Se você já executa estruturas de big data no Amazon EKS, agora é possível permitir que o Amazon EMR automatize o provisionamento e o gerenciamento.

emr-6.13.0-latest

Notas de versão: a versão emr-6.13.0-latest direciona para emr-6.13.0-20230814, no momento.

Regiões: emr-6.13.0-latest está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-6.13.0:latest.

emr-6.13.0-20230814

Notas de versão: 6.13.0-20230814 foi liberada em 7 de setembro de 2023. Esta é a versão inicial do Amazon EMR 6.13.0.

Atributos 476

Regiões: emr-6.13.0-20230814 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-6.13.0:20230814.

Versões 6.12.0 do Amazon EMR no EKS

Esta página descreve a funcionalidade nova e atualizada do Amazon EMR que é específica para a implantação do Amazon EMR no EKS. Para obter detalhes sobre o Amazon EMR em execução na Amazon EC2 e sobre a versão 6.12.0 do Amazon EMR em geral, consulte Amazon EMR 6.12.0 no Guia de lançamento do Amazon EMR.

Versões 6.12 do Amazon EMR no EKS

As versões 6.12.0 do Amazon EMR apresentadas a seguir estão disponíveis para o Amazon EMR no EKS. Selecione uma versão específica do emr-6.12.0-XXXX para visualizar mais detalhes, como a etiqueta de imagem do contêiner relacionada.

- emr-6.12.0-latest
- emr-6.12.0-20240321
- emr-6.12.0-20230701
- emr-6.12.0- spark-rapids-latest
- emr-6.12.0-spark-rapids-20230701
- emr-6.12.0-java11-latest
- emr-6.12.0-java11-20230701
- emr-6.12.0-java17-latest
- emr-6.12.0-java17-20230701
- emr-6.12.0- 17-mais recente spark-rapids-java
- emr-6.12.0- 17-20230701 spark-rapids-java
- notebook-spark/emr-6.12.0-latest
- notebook-spark/emr-6.12.0-20230701
- notebook-spark/emr-6.12.0- spark-rapids-latest
- notebook-spark/emr-6.12.0-spark-rapids-20230701
- notebook-python/emr-6.12.0-latest
- notebook-python/emr-6.12.0-20230701

Versões 6.12.0 477

- notebook-python/emr-6.12.0- spark-rapids-latest
- notebook-python/emr-6.12.0-spark-rapids-20230701

Notas da versão

Notas da versão 6.12.0 do Amazon EMR no EKS

- Aplicativos compatíveis AWS SDK para Java 1.12.490, Apache Spark 3.4.0-amzn-0, Apache Hudi 0.13.1-amzn-0, Apache Iceberg 1.3.0-amzn-0, Delta 2.4.0, Apache Spark RAPIDS 23.06.0-amzn-0, Jupyter Enterprise Gateway 2.6.0
- Componentes com suporte: aws-sagemaker-spark-sdk, emr-ddb, emr-goodies, emr-s3-select, emrfs, hadoop-client, hudi, hudi-spark, iceberg e spark-kubernetes.
- Classificações de configuração com suporte

Para uso com StartJobRune CreateManagedEndpoint APIs:

Classificações	Descrições
core-site	Altera os valores no arquivo core-site .xml do Hadoop.
emrfs-site	Alterar as configurações do EMRFS.
spark-metrics	Altera os valores no arquivo metrics.p roperties do Spark.
spark-defaults	Altera os valores no arquivo spark-def aults.conf do Spark.
spark-env	Alterar os valores no ambiente do Spark.
spark-hive-site	Altera os valores no arquivo hive-site .xml do Spark.
spark-log4j	Altera os valores no arquivo log4j2.pr operties do Spark.
emr-job-submitter	Configuração para o pod de envio de trabalho.

Para uso específico com CreateManagedEndpoint APIs:

Classificações	Descrições
jeg-config	Altera os valores no arquivo jupyter_e nterprise_gateway_config.py do Jupyter Enterprise Gateway.
jupyter-kernel-overrides	Altera o valor da imagem do kernel no arquivo de um kernel do Jupyter especificado.

As classificações de configuração permitem que você personalize aplicações. Elas geralmente correspondem a um arquivo XML de configuração da aplicação, como spark-hive-site.xml. Para obter mais informações, consulte Configure Applications.

Recursos notáveis

Os recursos apresentados a seguir estão inclusos na versão 6.12 do Amazon EMR no EKS.

 Java 17: com a versão 6.12 do Amazon EMR no EKS e versões superiores, é possível iniciar o Spark com o runtime do Java 17. Para fazer isso, transmita emr-6.12.0-java17-latest como uma etiqueta de versão. Recomendamos que você valide e execute testes de performance antes de mover as workloads de produção de versões anteriores da imagem do Java para a imagem do Java 17.

emr-6.12.0-latest

Notas de versão: a versão emr-6.12.0-latest direciona para emr-6.12.0-20240321, no momento.

Regiões: emr-6.12.0-latest está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-6.12.0:latest.

Atributos 479

emr-6.12.0-20240321

Notas de versão: 6.12.0-20240321 foi lançada em 11 de março de 2024. Em comparação com a versão anterior, esta versão foi atualizada com os pacotes do Amazon Linux atualizados recentemente e com as correções críticas.

Regiões: emr-6.12.0-20240321 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-6.12.0:20240321.

emr-6.12.0-20230701

Notas de versão: 6.12.0-20230701 foi liberada em 1.º de julho de 2023. Esta é a versão inicial do Amazon EMR 6.12.0.

Regiões: emr-6.12.0-20230701 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-6.12.0:20230701.

Versões 6.11.0 do Amazon EMR no EKS

Esta página descreve a funcionalidade nova e atualizada do Amazon EMR que é específica para a implantação do Amazon EMR no EKS. Para obter detalhes sobre o Amazon EMR em execução na Amazon EC2 e sobre a versão do Amazon EMR 6.11.0 em geral, consulte Amazon EMR 6.11.0 no Guia de lançamento do Amazon EMR.

Versões 6.11 do Amazon EMR no EKS

As versões 6.11.0 do Amazon EMR apresentadas a seguir estão disponíveis para o Amazon EMR no EKS. Selecione uma versão específica do emr-6.11.0-XXXX para visualizar mais detalhes, como a etiqueta de imagem do contêiner relacionada.

- emr-6.11.0-latest
- emr-6.11.0-20230905
- emr-6.11.0-20230509

emr-6.12.0-20240321 480

- emr-6.11.0- spark-rapids-latest
- emr-6.11.0-spark-rapids-20230509
- emr-6.11.0-java11-latest
- emr-6.11.0-java11-20230509
- notebook-spark/emr-6.11.0-latest
- notebook-spark/emr-6.11.0-20230509
- notebook-python/emr-6.11.0-latest
- notebook-python/emr-6.11.0-20230509

Notas da versão

Notas da versão 6.11.0 do Amazon EMR no EKS

- Aplicativos compatíveis AWS SDK para Java 1.12.446, Apache Spark 3.3.2-amzn-0, Apache Hudi 0.13.0-amzn-0, Apache Iceberg 1.2.0-amzn-0, Delta 2.2.0, Apache Spark RAPIDS 23.02.0-amzn-0, Jupyter Enterprise Gateway 2.6.0
- Componentes com suporte: aws-sagemaker-spark-sdk, emr-ddb, emr-goodies, emr-s3-select, emrfs, hadoop-client, hudi, hudi-spark, iceberg e spark-kubernetes.
- Classificações de configuração com suporte

Para uso com StartJobRune CreateManagedEndpoint APIs:

Classificações	Descrições
core-site	Altera os valores no arquivo core-site .xml do Hadoop.
emrfs-site	Alterar as configurações do EMRFS.
spark-metrics	Altera os valores no arquivo metrics.p roperties do Spark.
spark-defaults	Altera os valores no arquivo spark-def aults.conf do Spark.
spark-env	Alterar os valores no ambiente do Spark.

Classificações	Descrições
spark-hive-site	Altera os valores no arquivo hive-site .xml do Spark.
spark-log4j	Altera os valores no arquivo log4j.pro perties do Spark.

Para uso específico com CreateManagedEndpoint APIs:

Classificações	Descrições
jeg-config	Altera os valores no arquivo jupyter_e nterprise_gateway_config.py do Jupyter Enterprise Gateway.
jupyter-kernel-overrides	Altera o valor da imagem do kernel no arquivo de um kernel do Jupyter especificado.

As classificações de configuração permitem que você personalize aplicações. Elas geralmente correspondem a um arquivo XML de configuração da aplicação, como spark-hive-site.xml. Para obter mais informações, consulte <u>Configure Applications</u>.

Recursos notáveis

Os recursos apresentados a seguir estão inclusos na versão 6.11 do Amazon EMR no EKS.

- Imagem base do Amazon EMR no EKS na galeria pública do Amazon ECR: se você usar a
 funcionalidade imagem personalizada, nossa imagem base fornecerá os arquivos em JARs,
 as configurações e as bibliotecas essenciais para a interação com o Amazon EMR no EKS. É
 possível encontrar a imagem base na Galeria pública do Amazon ECR.
- Alternância de log do contêiner do Spark: a verão 6.11 do Amazon EMR no EKS oferece suporte para a alternância de log do contêiner do Spark. É possível habilitar a funcionalidade com containerLogRotationConfiguration na operação MonitoringConfiguration da API StartJobRun. Você pode configurar rotationSize e maxFilestoKeep para especificar o

Atributos 482

número e o tamanho dos arquivos de log que você deseja que o Amazon EMR no EKS mantenha nos pods do driver e do executor do Spark. Para obter mais informações, consulte <u>Uso da</u> alternância de log do contêiner do Spark.

Suporte ao Volcano no operador do Spark e no spark-submit: a versão 6.11 do Amazon EMR
no EKS oferece suporte à execução de trabalhos do Spark com o Volcano como programador
personalizado do Kubernetes no operador do Spark e no spark-submit. É possível usar recursos
como o agendamento coletivo, o gerenciamento de filas, a preempção e a programação de
compartilhamento equitativo para obter alto throughput de programação e capacidade otimizada.
Para obter mais informações, consulte Uso do Volcano como um programador personalizado para
Apache Spark no Amazon EMR no EKS.

emr-6.11.0-latest

Notas de versão: a versão emr-6.11.0-latest direciona para emr-20230905, no momento.

Regiões: emr-6.11.0-latest está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-6.11.0:latest.

emr-6.11.0-20230905

Notas de versão: 6.11.0-20230905 foi lançada em 29 de setembro de 2023. Em comparação com a versão anterior, esta versão foi atualizada com os pacotes do Amazon Linux atualizados recentemente e com as correções críticas.

Regiões: emr-6.11.0-20230509 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-6.11.0:20230509.

emr-6.11.0-20230509

Notas de versão: 6.11.0-20230509 foi liberada em 9 de maio de 2023. Esta é a versão inicial do Amazon EMR 6.11.0.

Regiões: emr-6.11.0-20230509 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-6.11.0:20230509.

emr-6.11.0-latest 483

Versões 6.10.0 do Amazon EMR no EKS

As versões 6.10.0 do Amazon EMR apresentadas a seguir estão disponíveis para o Amazon EMR no EKS. Selecione uma versão específica do emr-6.10.0-XXXX para visualizar mais detalhes, como a etiqueta de imagem do contêiner relacionada.

- emr-6.10.0-latest
- emr-6.10.0-20230905
- emr-6.10.0-20230624
- emr-6.10.0-20230421
- emr-6.10.0-20230403
- emr-6.10.0-20230220
- emr-6.10.0- spark-rapids-latest
- emr-6.10.0-spark-rapids-20230624
- emr-6.10.0-spark-rapids-20230220
- emr-6.10.0-java11-latest
- emr-6.10.0-java11-20230624
- emr-6.10.0-java11-20230220
- notebook-spark/emr-6.10.0-latest
- notebook-spark/emr-6.10.0-20230624
- notebook-spark/emr-6.10.0-20230220
- notebook-python/emr-6.10.0-latest
- notebook-python/emr-6.10.0-20230624
- notebook-python/emr-6.10.0-20230220

Notas da versão 6.10.0 do Amazon EMR

- Aplicativos suportados AWS SDK para Java 1.12.397, Spark 3.3.1-amzn-0, Hudi 0.12.2-amzn-0, Iceberg 1.1.0-amzn-0, Delta 2.2.0.
- Componentes com suporte: aws-sagemaker-spark-sdk, emr-ddb, emr-goodies, emr-s3-select, emrfs, hadoop-client, hudi, hudi-spark, iceberg e spark-kubernetes.
- Classificações de configuração com suporte:

Versões 6.10.0 484

$Para\ uso\ com\ \underline{StartJobRun}e\ \underline{CreateManagedEndpoint}\ APIs:$

Classificações	Descrições
core-site	Altera os valores no arquivo core-site .xml do Hadoop.
emrfs-site	Alterar as configurações do EMRFS.
spark-metrics	Altera os valores no arquivo metrics.p roperties do Spark.
spark-defaults	Altera os valores no arquivo spark-def aults.conf do Spark.
spark-env	Alterar os valores no ambiente do Spark.
spark-hive-site	Altera os valores no arquivo hive-site .xml do Spark.
spark-log4j	Altera os valores no arquivo log4j.pro perties do Spark.

Para uso específico com <u>CreateManagedEndpoint</u> APIs:

Classificações	Descrições
jeg-config	Altera os valores no arquivo jupyter_e nterprise_gateway_config.py do Jupyter Enterprise Gateway.
jupyter-kernel-overrides	Altera o valor da imagem do kernel no arquivo de um kernel do Jupyter especificado.

Versões 6.10.0 485

As classificações de configuração permitem que você personalize aplicações. Elas geralmente correspondem a um arquivo XML de configuração da aplicação, como spark-hive-site.xml. Para obter mais informações, consulte Configure Applications.

Recursos notáveis

- Operador do Spark: com a versão 6.10.0 do Amazon EMR no EKS e versões posteriores, é possível usar o operador do Kubernetes para Apache Spark ou o operador do Spark para implantar e gerenciar aplicações do Spark com o runtime de versão do Amazon EMR em seus próprios clusters do Amazon EKS. Para obter mais informações, consulte Execução de trabalhos do Spark com o operador do Spark.
- Java 11: com a versão 6.10 do Amazon EMR no EKS e versões superiores, é possível iniciar o Spark com o runtime do Java 11. Para fazer isso, transmita emr-6.10.0-java11-latest como uma etiqueta de versão. Recomendamos que você valide e execute testes de performance antes de mover as workloads de produção da imagem do Java 8 para a imagem do Java 11.
- Para a integração do Amazon Redshift para Apache Spark, a versão 6.10.0 do Amazon EMR no EKS remove a dependência de minimal-json.jar e adiciona automaticamente os arquivos em JARs relacionados ao spark-redshift obrigatórios ao caminho de classe do executor para o Spark: spark-redshift.jar, spark-avro.jar e RedshiftJDBC.jar.

Alterações

 O confirmador otimizado para EMRFS S3 está habilitado por padrão para Parquet, ORC e formatos baseados em texto (incluindo CSV e JSON).

emr-6.10.0-latest

Notas de versão: a versão emr-6.10.0-latest direciona para emr-6.10.0-20230905, no momento.

Regiões: emr-6.10.0-latest está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-6.10.0:latest.

emr-6.10.0-latest 486

emr-6.10.0-20230905

Notas de versão: 6.10.0-20230905 foi lançada em 29 de setembro de 2023. Em comparação com a versão anterior, esta versão foi atualizada com pacotes do Amazon Linux atualizados recentemente e correções críticas.

Regiões: emr-6.10.0-20230905 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-6.10.0:20230905.

emr-6.10.0-20230624

Notas de versão: 6.10.0-20230624 foi liberada em 7 de julho de 2023. Em comparação com a versão anterior, esta versão foi atualizada com pacotes do Amazon Linux atualizados recentemente e correções críticas.

Regiões: emr-6.10.0-20230624 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-6.10.0:20230624.

emr-6.10.0-20230421

Notas de versão: 6.10.0-20230421 foi liberada em 28 de abril de 2023. Em comparação com a versão anterior, esta versão foi atualizada com pacotes do Amazon Linux atualizados recentemente e correções críticas.

Regiões: emr-6.10.0-20230421 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-6.10.0:20230421.

emr-6.10.0-20230403

Notas de versão: 6.10.0-20230403 foi liberada em 12 de abril de 2023. Em comparação com a versão anterior, esta versão foi atualizada com pacotes do Amazon Linux atualizados recentemente e correções críticas.

Regiões: emr-6.10.0-20230403 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

emr-6.10.0-20230905 487

Etiqueta da imagem do contêiner: emr-6.10.0:20230403.

emr-6.10.0-20230220

Notas de versão: emr-6.10.0-20230220 foi liberada em 20 de fevereiro de 2023. Esta é a versão inicial do Amazon EMR 6.10.0.

Regiões: emr-6.10.0-20230220 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-6.10.0:20230220.

Versões 6.9.0 do Amazon EMR no EKS

As versões 6.9.0 do Amazon EMR apresentadas a seguir estão disponíveis para o Amazon EMR no EKS. Selecione uma versão específica do emr-6.9.0-XXXX para visualizar mais detalhes, como a etiqueta de imagem do contêiner relacionada.

- emr-6.9.0-latest
- emr-6.9.0-20230905
- emr-6.9.0-20230624
- emr-6.9.0-20221108
- emr-6.9.0- spark-rapids-latest
- emr-6.9.0-spark-rapids-20230624
- emr-6.9.0-spark-rapids-20221108
- notebook-spark/emr-6.9.0-latest
- notebook-spark/emr-6.9.0-20230624
- notebook-spark/emr-6.9.0-20221108
- notebook-python/emr-6.9.0-latest
- notebook-python/emr-6.9.0-20230624
- notebook-python/emr-6.9.0-20221108

Notas da versão 6.9.0 do Amazon EMR

emr-6.10.0-20230220 488

- Aplicativos suportados AWS SDK para Java 1.12.331, Spark 3.3.0-amzn-1, Hudi 0.12.1-amzn-0, Iceberg 0.14.1-amzn-0, Delta 2.1.0.
- Componentes com suporte: aws-sagemaker-spark-sdk, emr-ddb, emr-goodies, emr-s3-select, emrfs, hadoop-client, hudi, hudi-spark, iceberg e spark-kubernetes.
- Classificações de configuração com suporte:

Para uso com StartJobRune CreateManagedEndpoint APIs:

Classificações	Descrições
core-site	Altera os valores no arquivo core-site.xml do Hadoop.
emrfs-site	Alterar as configurações do EMRFS.
spark-metrics	Alterar os valores no arquivo metrics.p roperties do Spark.
spark-defaults	Alterar os valores no arquivo spark-def aults.conf do Spark.
spark-env	Alterar os valores no ambiente do Spark.
spark-hive-site	Altera os valores no arquivo hive-site.xml do Spark.
spark-log4j	Alterar os valores no arquivo log4j.properties do Spark.

Para uso específico com CreateManagedEndpoint APIs:

Classificações	Descrições
jeg-config	Altera os valores no arquivo jupyter_e nterprise_gateway_config.py do Jupyter Enterprise Gateway.

Versões 6.9.0 489

Classificações	Descrições
jupyter-kernel-overrides	Altera o valor da imagem do kernel no arquivo de um kernel do Jupyter especificado.

As classificações de configuração permitem que você personalize aplicações. Elas geralmente correspondem a um arquivo XML de configuração da aplicação, como spark-hive-site.xml. Para obter mais informações, consulte <u>Configure Applications</u>.

Recursos notáveis

- Nvidia RAPIDS Accelerator para Apache Spark Amazon EMR no EKS para acelerar o Spark usando tipos de instância de unidade de processamento EC2 gráfico (GPU). Para usar a imagem do Spark com o RAPIDS Accelerator, especifique o rótulo de lançamento como emr-6.9.0-. sparkrapids-latest Acesse a página de documentação para saber mais.
- Conector Spark-Redshift: a integração do Amazon Redshift para Apache Spark está inclusa nas versões 6.9.0 e posteriores do Amazon EMR. Anteriormente uma ferramenta de código aberto, a integração nativa é um conector do Spark que você pode usar para criar aplicações do Apache Spark que realizam a leitura e a gravação de dados no Amazon Redshift e no Amazon Redshift sem servidor. Para obter mais informações, consulte <u>Uso da integração do Amazon Redshift para</u> Apache Spark no Amazon EMR no EKS.
- Delta Lake: o <u>Delta Lake</u> é um formato de armazenamento de código aberto que possibilita o
 desenvolvimento de data lakes com consistência transacional, definição consistente de conjuntos
 de dados, alterações de evolução de esquema e suporte a mutações de dados. Acesse <u>Uso do</u>
 <u>Delta Lake</u> para saber mais.
- Modificar PySpark parâmetros Os endpoints interativos agora oferecem suporte à modificação dos parâmetros do Spark associados às PySpark sessões no notebook EMR Studio Jupyter. Visite Modificação dos parâmetros PySpark da sessão para saber mais.

Problemas resolvidos

 Ao usar o conector DynamoDB com o Spark nas versões 6.6.0, 6.7.0 e 6.8.0 do Amazon EMR, todas as leituras da tabela retornam um resultado vazio, mesmo que a divisão de entrada faça referência a dados que não estão vazios. A versão 6.9.0 do Amazon EMR corrige esse problema.

Versões 6.9.0 490

 A versão 6.8.0 do Amazon EMR no EKS preenche incorretamente o hash de compilação nos metadados dos arquivos em Parquet gerados usando o <u>Apache Spark</u>. Esse problema pode fazer com que as ferramentas que analisam a string de versão de metadados dos arquivos em Parquet gerados pela versão 6.8.0 do Amazon EMR no EKS apresentem falhas.

Problema conhecido

• Se você usar a integração do Amazon Redshift para Apache Spark e tiver um time, timetz, timestamp ou timestamptz com precisão de microssegundos no formato Parquet, o conector arredondará os valores de tempo para o valor de milissegundo mais próximo. Como solução alternativa, use o parâmetro unload_s3_format do formato de descarregamento de texto.

emr-6.9.0-latest

Notas de versão: a versão emr-6.9.0-latest direciona para emr-6.9.0-20230905, no momento.

Regiões: emr-6.9.0-latest está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-6.9.0:latest.

emr-6.9.0-20230905

Notas da versão: emr-6.9.0-20230905. Em comparação com a versão anterior, esta versão foi atualizada com os pacotes do Amazon Linux atualizados recentemente e com as correções críticas.

Regiões: emr-6.9.0-20230905 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-6.9.0:20230905.

emr-6.9.0-20230624

Notas de versão: emr-6.9.0-20230624 foi liberada em 7 de julho de 2023.

Regiões: emr-6.9.0-20230624 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-6.9.0:20230624.

emr-6.9.0-latest 491

emr-6.9.0-20221108

Notas de versão: emr-6.9.0-20221108 foi liberada em 8 de dezembro de 2022. Esta é a versão inicial do Amazon EMR 6.9.0.

Regiões: emr-6.9.0-20221108 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-6.9.0:20221108.

Versões 6.8.0 do Amazon EMR no EKS

As versões 6.8.0 do Amazon EMR apresentadas a seguir estão disponíveis para o Amazon EMR no EKS. Selecione uma versão específica do emr-6.8.0-XXXX para visualizar mais detalhes, como a etiqueta de imagem do contêiner relacionada.

- emr-6.8.0-latest
- emr-6.8.0-20230905
- emr-6.8.0-20230624
- emr-6.8.0-20221219
- emr-6.8.0-20220802

Notas da versão 6.8.0 do Amazon EMR

- Aplicativos suportados AWS SDK para Java 1.12.170, Spark 3.3.0-amzn-0, Hudi 0.11.1-amzn-0, Iceberg 0.14.0-amzn-0.
- Componentes com suporte: aws-sagemaker-spark-sdk, emr-ddb, emr-goodies, emr-s3-select, emrfs, hadoop-client, hudi, hudi-spark, iceberg e spark-kubernetes.
- Classificações de configuração com suporte:

Classificações	Descrições
core-site	Altera os valores no arquivo core-site.xml do Hadoop.
emrfs-site	Alterar as configurações do EMRFS.

emr-6.9.0-20221108 492

Classificações	Descrições
spark-metrics	Alterar os valores no arquivo metrics.p roperties do Spark.
spark-defaults	Alterar os valores no arquivo spark-def aults.conf do Spark.
spark-env	Alterar os valores no ambiente do Spark.
spark-hive-site	Altera os valores no arquivo hive-site.xml do Spark.
spark-log4j	Alterar os valores no arquivo log4j.properties do Spark.

As classificações de configuração permitem que você personalize aplicações. Elas geralmente correspondem a um arquivo XML de configuração da aplicação, como spark-hive-site.xml. Para obter mais informações, consulte Configure Applications.

Recursos notáveis

- Spark3.3.0: a versão 6.8 Amazon EMR no EKS inclui o Spark 3.3.0, que oferece suporte ao
 uso de rótulos seletores de nós separados para pods de executores e de drivers do Spark.
 Esses novos rótulos permitem que você defina os tipos de nós para os pods de driver e executor
 separadamente na StartJobRun API, sem usar modelos de pod.
 - Propriedade do seletor do nó do driver: spark.kubernetes.driver.node.selector.[labelKey]
 - Propriedade do seletor do nó do executor: spark.kubernetes.executor.node.selector.[labelKey]
- Mensagem aprimorada de falha de trabalho: esta versão apresenta as configurações spark.stage.extraDetailsOnFetchFailures.enabled e spark.stage.extraDetailsOnFetchFailures.maxFailuresToInclude para rastrear falhas de tarefas devido ao código do usuário. Esses detalhes serão usados para aprimorar a mensagem de falha exibida no log do driver quando um estágio é interrompido devido a uma falha na busca aleatória.

Versões 6.8.0 493

Nome da propriedade	Valor padrão	Significado	Desde a versão
spark.sta ge.extraD etailsOnF etchFailu res.enabled	false	Se definida como true, esta proprieda de será usada para aprimorar a mensagem de falha do trabalho exibida no log do driver quando um estágio é interrompido devido a uma falha de busca aleatória. Por padrão, as últimas cinco falhas de tarefas causadas pelo código do usuário são rastreadas e a mensagem de erro de falha é anexada aos logs do driver. Para aumentar o número de falhas de tarefas com exceções de usuário a serem rastreadas, consulte a configura ção spark.sta ge.extraD etailsOnF etchFailu res.maxFa iluresToI nclude.	emr-6.8

Versões 6.8.0 494

Nome da propriedade	Valor padrão	Significado	Desde a versão
spark.sta ge.extraD etailsOnF etchFailu res.maxFa iluresToI nclude	5	Número de falhas de tarefas a serem rastreadas por estágio e por tentativa. Esta propriedade é usada para aprimorar a mensagem de falha do trabalho com exceções de usuário exibidas no log do driver quando um estágio é interromp ido devido a uma falha de busca aleatória. Essa propriedade funciona somente se Config spark.stage. extraDetailsOnFetc hFailures.enabled está definido como verdadeiro.	emr-6.8

Para obter mais informações, consulte a documentação de configuração do Apache Spark.

Problema conhecido

• A versão 6.8.0 do Amazon EMR no EKS preenche incorretamente o hash de compilação nos metadados dos arquivos em Parquet gerados usando o <u>Apache Spark</u>. Esse problema pode fazer com que as ferramentas que analisam a string de versão de metadados dos arquivos em Parquet gerados pela versão 6.8.0 do Amazon EMR no EKS apresentem falhas. Os clientes que analisam a string de versão dos metadados do Parquet e dependem do hash de compilação devem realizar a alteração para uma versão diferente do Amazon EMR e reescrever o arquivo.

Versões 6.8.0 495

Problema resolvido

 Funcionalidade de interrupção do kernel para kernels do PySpark: as workloads interativas em andamento que são acionadas pela execução de células em um caderno podem ser interrompidas usando a funcionalidade Interrupt Kernel. Uma correção foi introduzida para que esta funcionalidade funcione para kernels do PySpark. Isso também está disponível em código aberto em Changes para lidar com interrupções no PySpark Kubernetes Kernel #1115.

emr-6.8.0-latest

Notas de versão: a versão emr-6.8.0-latest direciona para emr-6.8.0-20230624, no momento.

Regiões: emr-6.8.0-latest está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-6.8.0:latest.

emr-6.8.0-20230905

Notas de versão: emr-6.8.0-20230905 foi lançada em 29 de setembro de 2023. Em comparação com a versão anterior, esta versão foi atualizada com os pacotes do Amazon Linux atualizados recentemente e com as correções críticas.

Regiões: emr-6.8.0-20230905 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-6.8.0:20230905.

emr-6.8.0-20230624

Notas de versão: emr-6.8.0-20230624 foi liberada em 7 de julho de 2023. Em comparação com a versão anterior, esta versão foi atualizada com os pacotes do Amazon Linux atualizados recentemente e com as correções críticas.

Regiões: emr-6.8.0-20230624 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-6.8.0:20230624.

emr-6.8.0-latest 496

emr-6.8.0-20221219

Notas de versão: emr-6.8.0-20221219 foi liberada em 19 de janeiro de 2023. Em comparação com a versão anterior, esta versão foi atualizada com os pacotes do Amazon Linux atualizados recentemente e com as correções críticas.

Regiões: emr-6.8.0-20221219 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-6.8.0:20221219.

emr-6.8.0-20220802

Notas de versão: emr-6.8.0-20220802 foi liberada em 27 de setembro de 2022. Esta é a versão inicial do Amazon EMR 6.8.0.

Regiões: emr-6.8.0-20220802 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-6.8.0:20220802.

Versões 6.7.0 do Amazon EMR no EKS

As versões 6.7.0 do Amazon EMR apresentadas a seguir estão disponíveis para o Amazon EMR no EKS. Selecione uma versão específica do emr-6.7.0-XXXX para visualizar mais detalhes, como a etiqueta de imagem do contêiner relacionada.

- emr-6.7.0-latest
- emr-6.7.0-20240321
- emr-6.7.0-20230624
- emr-6.7.0-20221219
- emr-6.7.0-20220630

Notas da versão 6.7.0 do Amazon EMR

 Aplicações com suporte: Spark 3.2.1-amzn-0, Jupyter Enterprise Gateway 2.6, Hudi 0.11-amzn-0 e Iceberg 0.13.1.

emr-6.8.0-20221219 497

- Componentes com suporte: aws-hm-client (conector do Glue), aws-sagemaker-spark-sdk, emr-s3-select, emrfs, emr-ddb e hudi-spark.
- Com a atualização para o JEG 2.6, o gerenciamento do kernel passou a ser assíncrono, o que significa que o JEG não bloqueia transações quando a inicialização do kernel está em andamento.
 Isso melhora muito a experiência do usuário ao fornecer o seguinte:
 - capacidade de executar comandos em cadernos que estão em execução no momento quando outras inicializações de kernel estão em andamento;
 - capacidade de iniciar vários kernels simultaneamente sem afetar os kernels que já estão em execução.
- Classificações de configuração com suporte:

Classificações	Descrições
core-site	Altera os valores no arquivo core-site .xml do Hadoop.
emrfs-site	Alterar as configurações do EMRFS.
spark-metrics	Altera os valores no arquivo metrics.p roperties do Spark.
spark-defaults	Altera os valores no arquivo spark-def aults.conf do Spark.
spark-env	Alterar os valores no ambiente do Spark.
spark-hive-site	Altera os valores no arquivo hive-site .xml do Spark.
spark-log4j	Altera os valores no arquivo log4j.pro perties do Spark.

As classificações de configuração permitem que você personalize aplicações. Elas geralmente correspondem a um arquivo XML de configuração da aplicação, como spark-hive-site.xml. Para obter mais informações, consulte Configuração de aplicações.

Versões 6.7.0 498

Problemas resolvidos

- A versão 6.7 do Amazon EMR no EKS corrige um problema na versão 6.6 relacionado ao uso
 da funcionalidade de modelos de pod do Apache Spark com endpoints interativos. O problema
 estava presente nas versões 6.4, 6.5 e 6.6 do Amazon EMR no EKS. No momento, você pode
 usar modelos de pod para definir como os pods de driver e de executores do Spark são iniciados
 ao usar endpoints interativos para executar analytics interativos.
- Nas versões anteriores do Amazon EMR no EKS, o Jupyter Enterprise Gateway bloqueava transações quando a inicialização do kernel estava em andamento, e isso impedia a execução de sessões de caderno que estavam em execução no momento. Agora, é possível executar comandos em cadernos em execução quando outras inicializações de kernel estiverem em andamento. Você também pode iniciar vários kernels simultaneamente sem o risco de perder a conectividade com os kernels que já estão em execução.

emr-6.7.0-latest

Notas de versão: a versão emr-6.7.0-latest direciona para emr-6.7.0-20240321, no momento.

Regiões: emr-6.7.0-latest está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-6.7.0:latest.

emr-6.7.0-20240321

Notas de versão: emr-6.7.0-20240321 foi lançada em 11 de março de 2024. Em comparação com a versão anterior, esta versão foi atualizada com os pacotes do Amazon Linux atualizados recentemente e com as correções críticas.

Regiões: emr-6.7.0-20240321 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-6.7.0:20240321.

emr-6.7.0-latest 499

emr-6.7.0-20230624

Notas de versão: emr-6.7.0-20230624 foi liberada em 7 de julho de 2023. Em comparação com a versão anterior, esta versão foi atualizada com os pacotes do Amazon Linux atualizados recentemente e com as correções críticas.

Regiões: emr-6.7.0-20230624 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-6.7.0:20230624.

emr-6.7.0-20221219

Notas de versão: emr-6.7.0-20221219 foi liberada em 19 de janeiro de 2023. Em comparação com a versão anterior, esta versão foi atualizada com os pacotes do Amazon Linux atualizados recentemente e com as correções críticas.

Regiões: emr-6.7.0-20221219 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-6.7.0:20221219.

emr-6.7.0-20220630

Notas de versão: emr-6.7.0-20220630 foi liberada em 12 de julho de 2022. Esta é a versão inicial do Amazon EMR 6.7.0.

Regiões: emr-6.7.0-20220630 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-6.7.0:20220630.

Versões 6.6.0 do Amazon EMR no EKS

As versões 6.6.0 do Amazon EMR apresentadas a seguir estão disponíveis para o Amazon EMR no EKS. Selecione uma versão específica do emr-6.6.0-XXXX para visualizar mais detalhes, como a etiqueta de imagem do contêiner relacionada.

- emr-6.6.0-latest
- emr-6.6.0-20240321
- emr-6.6.0-20230624

emr-6.7.0-20230624 500

- emr-6.6.0-20221219
- emr-6.6.0-20220411

Notas da versão 6.6.0 do Amazon EMR

- Aplicações com suporte: Spark 3.2.0-amzn-0, Jupyter Enterprise Gateway (endpoints e prévisualização pública), Hudi 0.10.1-amzn-0 e Iceberg 0.13.1.
- Componentes com suporte: aws-hm-client (conector do Glue), aws-sagemaker-spark-sdk, emr-s3-select, emrfs, emr-ddb e hudi-spark.
- Classificações de configuração com suporte:

Classificações	Descrições
core-site	Altera os valores no arquivo core-site.xml do Hadoop.
emrfs-site	Alterar as configurações do EMRFS.
spark-metrics	Alterar os valores no arquivo metrics.p roperties do Spark.
spark-defaults	Alterar os valores no arquivo spark-def aults.conf do Spark.
spark-env	Alterar os valores no ambiente do Spark.
spark-hive-site	Altera os valores no arquivo hive-site.xml do Spark.
spark-log4j	Alterar os valores no arquivo log4j.properties do Spark.

As classificações de configuração permitem que você personalize aplicações. Eles geralmente correspondem a um arquivo XML de configuração do aplicativo, como spark-hive-site .xml. Para obter mais informações, consulte Configuração de aplicações.

Problema conhecido

Versões 6.6.0 501

A funcionalidade do modelo de pod do Spark com endpoints interativos não funciona nas versões
 6.4, 6.5 e 6.6 do Amazon EMR no EKS.

Problema resolvido

Os logs de endpoint interativos s\(\tilde{a}\) carregados no CloudWatch e no S3.

emr-6.6.0-latest

Notas de versão: a versão emr-6.6.0-latest direciona para emr-6.6.0-20240321, no momento.

Regiões: emr-6.6.0-latest está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-6.6.0:latest.

emr-6.6.0-20240321

Notas de versão: emr-6.6.0-20240321 foi lançada em 11 de março de 2024. Em comparação com a versão anterior, esta versão foi atualizada com os pacotes do Amazon Linux atualizados recentemente e com as correções críticas.

Regiões: emr-6.6.0-20240321 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-6.6.0:20240321.

emr-6.6.0-20230624

Notas de versão: emr-6.6.0-20230624 foi liberada em 27 de janeiro de 2023. Em comparação com a versão anterior, esta versão foi atualizada com os pacotes do Amazon Linux atualizados recentemente e com as correções críticas.

Regiões: emr-6.6.0-20230624 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-6.6.0:20230624.

emr-6.6.0-latest 502

emr-6.6.0-20221219

Notas de versão: emr-6.6.0-20221219 foi liberada em 27 de janeiro de 2023. Em comparação com a versão anterior, esta versão foi atualizada com os pacotes do Amazon Linux atualizados recentemente e com as correções críticas.

Regiões: emr-6.6.0-20221219 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-6.6.0:20221219.

emr-6.6.0-20220411

Notas de versão: emr-6.6.0-20220411 foi liberada em 20 de maio de 2022. Esta é a versão inicial do Amazon EMR 6.6.0.

Regiões: emr-6.6.0-20220411 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-6.6.0:20220411.

Versões 6.5.0 do Amazon EMR no EKS

As versões 6.5.0 do Amazon EMR apresentadas a seguir estão disponíveis para o Amazon EMR no EKS. Selecione uma versão específica do emr-6.5.0-XXXX para visualizar mais detalhes, como a etiqueta de imagem do contêiner relacionada.

- emr-6.5.0-latest
- emr-6.5.0-20240321
- emr-6.5.0-20221219
- emr-6.5.0-20220802
- emr-6.5.0-20211119

Notas da versão 6.5.0 do Amazon EMR

- Aplicações com suporte: Spark 3.1.2-amzn-1 e Jupyter Enterprise Gateway (endpoints e prévisualização pública).
- Componentes com suporte: aws-hm-client (conector do Glue), aws-sagemaker-spark-sdk, emr-s3-select, emrfs, emr-ddb e hudi-spark.

emr-6.6.0-20221219 503

• Classificações de configuração com suporte:

Classificações	Descrições
core-site	Altera os valores no arquivo core-site.xml do Hadoop.
emrfs-site	Alterar as configurações do EMRFS.
spark-metrics	Alterar os valores no arquivo metrics.p roperties do Spark.
spark-defaults	Alterar os valores no arquivo spark-def aults.conf do Spark.
spark-env	Alterar os valores no ambiente do Spark.
spark-hive-site	Altera os valores no arquivo hive-site.xml do Spark.
spark-log4j	Alterar os valores no arquivo log4j.properties do Spark.

As classificações de configuração permitem que você personalize aplicações. Eles geralmente correspondem a um arquivo XML de configuração do aplicativo, como spark-hive-site .xml. Para obter mais informações, consulte Configuração de aplicações.

Problema conhecido

A funcionalidade do modelo de pod do Spark com endpoints interativos não funciona nas versões
 6.4 e 6.5 do Amazon EMR no EKS.

emr-6.5.0-latest

Notas de versão: a versão emr-6.5.0-latest direciona para emr-6.5.0-20240321, no momento.

emr-6.5.0-latest 504

Regiões: emr-6.5.0-latest está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-6.5.0:latest.

emr-6.5.0-20240321

Notas de versão: emr-6.5.0-20240321 foi lançada em 11 de março de 2024. Em comparação com a versão anterior, esta versão foi atualizada com os pacotes do Amazon Linux atualizados recentemente e com as correções críticas.

Regiões: emr-6.5.0-20240321 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-6.5.0:20240321.

emr-6.5.0-20221219

Notas de versão: emr-6.5.0-20221219 foi liberada em 19 de janeiro de 2023. Em comparação com a versão anterior, esta versão foi atualizada com os pacotes do Amazon Linux atualizados recentemente e com as correções críticas.

Regiões: emr-6.5.0-20221219 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-6.5.0:20221219.

emr-6.5.0-20220802

Notas de versão: emr-6.5.0-20220802 foi liberada em 24 de agosto de 2022. Em comparação com a versão anterior, esta versão foi atualizada com os pacotes do Amazon Linux atualizados recentemente.

Regiões: emr-6.5.0-20220802 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-6.5.0:20220802.

emr-6.5.0-20211119

Notas de versão: emr-6.5.0-20211119 foi liberada em 20 de janeiro de 2022. Esta é a versão inicial do Amazon EMR 6.5.0.

emr-6.5.0-20240321 505

Regiões: emr-6.5.0-20211119 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-6.5.0:20211119.

Versões 6.4.0 do Amazon EMR no EKS

As versões 6.4.0 do Amazon EMR apresentadas a seguir estão disponíveis para o Amazon EMR no EKS. Selecione uma versão específica do emr-6.4.0-XXXX para visualizar mais detalhes, como a etiqueta de imagem do contêiner relacionada.

- emr-6.4.0-latest
- emr-6.4.0-20240321
- emr-6.4.0-20221219
- emr-6.4.0-20210830

Notas da versão 6.4.0 do Amazon EMR

- Aplicações com suporte: Spark 3.1.2-amzn-0 e Jupyter Enterprise Gateway (endpoints e prévisualização pública).
- Componentes com suporte: aws-hm-client (conector do Glue), aws-sagemaker-spark-sdk, emr-s3-select, emrfs, emr-ddb e hudi-spark.
- Classificações de configuração com suporte:

Classificações	Descrições
core-site	Altera os valores no arquivo core-site.xml do Hadoop.
emrfs-site	Alterar as configurações do EMRFS.
spark-metrics	Alterar os valores no arquivo metrics.p roperties do Spark.
spark-defaults	Alterar os valores no arquivo spark-def aults.conf do Spark.

Versões 6.4.0 506

Classificações	Descrições
spark-env	Alterar os valores no ambiente do Spark.
spark-hive-site	Altera os valores no arquivo hive-site.xml do Spark.
spark-log4j	Alterar os valores no arquivo log4j.properties do Spark.

As classificações de configuração permitem que você personalize aplicações. Eles geralmente correspondem a um arquivo XML de configuração do aplicativo, como spark-hive-site .xml. Para obter mais informações, consulte Configuração de aplicações.

Problema conhecido

 A funcionalidade do modelo de pod do Spark com endpoints interativos não funciona na versão 6.4 do Amazon EMR no EKS.

emr-6.4.0-latest

Notas de versão: a versão emr-6.4.0-latest direciona para emr-6.4.0-20240321, no momento.

Regiões: emr-6.4.0-latest está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-6.4.0:latest.

emr-6.4.0-20240321

Notas de versão: emr-6.4.0-20240321 foi lançada em 11 de março de 2024. Em comparação com a versão anterior, esta versão foi atualizada com os pacotes do Amazon Linux atualizados recentemente e com as correções críticas.

Regiões: emr-6.4.0-20240321 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-6.4.0:20240321.

emr-6.4.0-latest 507

emr-6.4.0-20221219

Notas de versão: emr-6.4.0-20221219 foi liberada em 27 de janeiro de 2023. Em comparação com a versão anterior, esta versão foi atualizada com os pacotes do Amazon Linux adicionados recentemente.

Regiões: emr-6.4.0-20221219 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-6.4.0:20221219.

emr-6.4.0-20210830

Notas de versão: emr-6.4.0-20210830 foi liberada em 9 de dezembro de 2021. Esta é a versão inicial do Amazon EMR 6.4.0.

Regiões: emr-6.4.0-20210830 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-6.4.0:20210830.

Versões 6.3.0 do Amazon EMR no EKS

As versões 6.3.0 do Amazon EMR apresentadas a seguir estão disponíveis para o Amazon EMR no EKS. Selecione uma versão específica do emr-6.3.0-XXXX para visualizar mais detalhes, como a etiqueta de imagem do contêiner relacionada.

- emr-6.3.0-latest
- emr-6.3.0-20240321
- emr-6.3.0-20220802
- emr-6.3.0-20211008
- emr-6.3.0-20210802
- emr-6.3.0-20210429

Notas da versão 6.3.0 do Amazon EMR

 Novos recursos: começando com o Amazon EMR 6.3.0 na série de versões 6.x, o Amazon EMR no EKS oferece suporte ao recurso de modelo de pod do Spark. Também é possível ativar o

emr-6.4.0-20221219 508

recurso de alternância de log de eventos do Spark para o Amazon EMR no EKS. Para obter mais informações, consulte Uso de modelos de pod e Uso da alternância de log de eventos do Spark.

- Aplicações com suporte: Spark 3.1.1-amzn-0 e Jupyter Enterprise Gateway (endpoints e prévisualização pública).
- Componentes com suporte: aws-hm-client (conector do Glue), aws-sagemaker-spark-sdk, emr-s3-select, emrfs, emr-ddb e hudi-spark.
- Classificações de configuração com suporte:

Classificações	Descrições
core-site	Altera os valores no arquivo core-site.xml do Hadoop.
emrfs-site	Alterar as configurações do EMRFS.
spark-metrics	Alterar os valores no arquivo metrics.p roperties do Spark.
spark-defaults	Alterar os valores no arquivo spark-def aults.conf do Spark.
spark-env	Alterar os valores no ambiente do Spark.
spark-hive-site	Altera os valores no arquivo hive-site.xml do Spark.
spark-log4j	Alterar os valores no arquivo log4j.properties do Spark.

As classificações de configuração permitem que você personalize aplicações. Eles geralmente correspondem a um arquivo XML de configuração do aplicativo, como spark-hive-site .xml. Para obter mais informações, consulte Configuração de aplicações.

emr-6.3.0-latest

Notas de versão: a versão emr-6.3.0-latest direciona para emr-6.3.0-20240321, no momento.

emr-6.3.0-latest 509

Regiões: emr-6.3.0-latest está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-6.3.0:latest.

emr-6.3.0-20240321

Notas de versão: emr-6.3.0-20240321 foi lançada em 11 de março de 2024. Em comparação com a versão anterior, esta versão foi atualizada com os pacotes do Amazon Linux atualizados recentemente e com as correções críticas.

Regiões: emr-6.3.0-20240321 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-6.3.0:20240321.

emr-6.3.0-20220802

Notas de versão: emr-6.3.0-20220802 foi liberada em 27 de setembro de 2022. Em comparação com a versão anterior, esta versão foi atualizada com os pacotes do Amazon Linux atualizados recentemente.

Regiões: emr-6.3.0-20220802 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-6.3.0:20220802.

emr-6.3.0-20211008

Notas de versão: emr-6.3.0-20211008 foi liberada em 9 de dezembro de 2021. Em comparação com a versão anterior, esta versão contém correções de problemas e atualizações de segurança.

Regiões: emr-6.3.0-20211008 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-6.3.0:20211008.

emr-6.3.0-20210802

Notas de versão: emr-6.3.0-20210802 foi liberada em 2 de agosto de 2021. Em comparação com a versão anterior, esta versão contém correções de problemas e atualizações de segurança.

emr-6.3.0-20240321 510

Regiões: emr-6.3.0-20210802 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-6.3.0:20210802.

emr-6.3.0-20210429

Notas de versão: emr-6.3.0-20210429 foi liberada em 29 de abril de 2021. Esta é a versão inicial do Amazon EMR 6.3.0.

Regiões: emr-6.3.0-20210429 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-6.3.0:20210429.

Versões 6.2.0 do Amazon EMR no EKS

As versões 6.2.0 do Amazon EMR apresentadas a seguir estão disponíveis para o Amazon EMR no EKS. Selecione uma versão específica do emr-6.2.0-XXXX para visualizar mais detalhes, como a etiqueta de imagem do contêiner relacionada.

- emr-6.2.0-latest
- emr-6.2.0-20240321
- emr-6.2.0-20220802
- emr-6.2.0-20211008
- emr-6.2.0-20210802
- emr-6.2.0-20210615
- emr-6.2.0-20210129
- emr-6.2.0-20201218
- emr-6.2.0-20201201

Notas da versão 6.2.0 do Amazon EMR

 Aplicações com suporte: Spark 3.0.1-amzn-0 e Jupyter Enterprise Gateway (endpoints e prévisualização pública).

emr-6.3.0-20210429 511

- Componentes com suporte: aws-hm-client (conector do Glue), aws-sagemaker-spark-sdk, emr-s3-select, emrfs, emr-ddb e hudi-spark.
- Classificações de configuração com suporte:

Classificações	Descrições
core-site	Altera os valores no arquivo core-site.xml do Hadoop.
emrfs-site	Alterar as configurações do EMRFS.
spark-metrics	Alterar os valores no arquivo metrics.p roperties do Spark.
spark-defaults	Alterar os valores no arquivo spark-def aults.conf do Spark.
spark-env	Alterar os valores no ambiente do Spark.
spark-hive-site	Altera os valores no arquivo hive-site.xml do Spark.
spark-log4j	Alterar os valores no arquivo log4j.properties do Spark.

As classificações de configuração permitem que você personalize aplicações. Eles geralmente correspondem a um arquivo XML de configuração do aplicativo, como spark-hive-site .xml. Para obter mais informações, consulte <u>Configuração de aplicações</u>.

emr-6.2.0-latest

Notas de versão: a versão emr-6.2.0-latest direciona para emr-6.2.0-20240321, no momento.

Regiões: emr-6.2.0-latest está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-6.2.0:20240321.

emr-6.2.0-latest 512

emr-6.2.0-20240321

Notas de versão: emr-6.2.0-20240321 foi lançada em 11 de março de 2024. Em comparação com a versão anterior, esta versão foi atualizada com os pacotes do Amazon Linux atualizados recentemente e com as correções críticas.

Regiões: emr-6.2.0-20240321 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-6.2.0:20240321.

emr-6.2.0-20220802

Notas de versão: emr-6.2.0-20220802 foi liberada em 27 de setembro de 2022. Em comparação com a versão anterior, esta versão foi atualizada com os pacotes do Amazon Linux atualizados recentemente.

Regiões: emr-6.2.0-20220802 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte <u>Endpoints de serviço do Amazon EMR no EKS</u>.

Etiqueta da imagem do contêiner: emr-6.2.0:20220802.

emr-6.2.0-20211008

Notas de versão: emr-6.2.0-20211008 foi liberada em 9 de dezembro de 2021. Em comparação com a versão anterior, esta versão contém correções de problemas e atualizações de segurança.

Regiões: emr-6.2.0-20211008 está disponível nas seguintes regiões: Leste dos EUA (Norte da Virgínia), Oeste dos EUA (Oregon), Ásia-Pacífico (Tóquio), Europa (Irlanda) e América do Sul (São Paulo).

Etiqueta da imagem do contêiner: emr-6.2.0:20211008.

emr-6.2.0-20210802

Notas de versão: emr-6.2.0-20210802 foi liberada em 2 de agosto de 2021. Em comparação com a versão anterior, esta versão contém correções de problemas e atualizações de segurança.

Regiões: emr-6.2.0-20210802 está disponível nas seguintes regiões: Leste dos EUA (Norte da Virgínia), Oeste dos EUA (Oregon), Ásia-Pacífico (Tóquio), Europa (Irlanda) e América do Sul (São Paulo).

emr-6.2.0-20240321 513

Etiqueta da imagem do contêiner: emr-6.2.0:20210802.

emr-6.2.0-20210615

Notas de versão: emr-6.2.0-20210615 foi liberada em 15 de junho de 2021. Em comparação com a versão anterior, esta versão contém correções de problemas e atualizações de segurança.

Regiões: emr-6.2.0-20210615 está disponível nas seguintes regiões: Leste dos EUA (Norte da Virgínia), Oeste dos EUA (Oregon), Ásia-Pacífico (Tóquio), Europa (Irlanda) e América do Sul (São Paulo).

Etiqueta da imagem do contêiner: emr-6.2.0:20210615.

emr-6.2.0-20210129

Notas de versão: emr-6.2.0-20210129 foi liberada em 29 de janeiro de 2021. Em comparação com emr-6.2.0-20201218, esta versão contém correções de problemas e atualizações de segurança.

Regiões: emr-6.2.0-20210129 está disponível nas seguintes regiões: Leste dos EUA (Norte da Virgínia), Oeste dos EUA (Oregon), Ásia-Pacífico (Tóquio), Europa (Irlanda) e América do Sul (São Paulo).

Etiqueta da imagem do contêiner: emr-6.2.0-20210129.

emr-6.2.0-20201218

Notas de versão: emr-6.2.0-20201218 foi liberada em 18 de dezembro de 2020. Em comparação com emr-6.2.0-20201201, esta versão contém correções de problemas e atualizações de segurança.

Regiões: emr-6.2.0-20201218 está disponível nas seguintes regiões: Leste dos EUA (Norte da Virgínia), Oeste dos EUA (Oregon), Ásia-Pacífico (Tóquio), Europa (Irlanda) e América do Sul (São Paulo).

Etiqueta da imagem do contêiner: emr-6.2.0-20201218.

emr-6.2.0-20201201

Notas de versão: emr-6.2.0-20201201 foi liberada em 1.º de dezembro de 2020. Esta é a versão inicial do Amazon EMR 6.2.0.

emr-6.2.0-20210615 514

Regiões: emr-6.2.0-20201201 está disponível nas seguintes regiões: Leste dos EUA (Norte da Virgínia), Oeste dos EUA (Oregon), Ásia-Pacífico (Tóquio), Europa (Irlanda) e América do Sul (São Paulo).

Etiqueta da imagem do contêiner: emr-6.2.0-20201201.

Versões 5.36.0 do Amazon EMR no EKS

As versões 5.36.0 do Amazon EMR apresentadas a seguir estão disponíveis para o Amazon EMR no EKS. Selecione uma versão específica do emr-5.36.0-XXXX para visualizar mais detalhes, como a etiqueta de imagem do contêiner relacionada.

- emr-5.36.0-latest
- emr-5.36.0-20240321
- emr-5.36.0-20221219
- emr-5.36.0-20220620
- emr-5.36.0-20220525

Notas da versão 5.36.0 do Amazon EMR

- Os problemas de segurança do log4j2 foram corrigidos.
- Aplicações com suporte: Spark 2.4.8-amzn-2, Jupyter Enterprise Gateway (endpoints e prévisualização pública; o kernel do Scala não tem suporte), livy-0.7.1 e fluentd-4.0.0.
- Componentes suportados aws-hm-client, emr-ddb aws-sagemaker-spark-sdk, emr-goodies, emr-kinesis, kerberos-server.
- Classificações de configuração com suporte:

Classificações	Descrições
core-site	Altera os valores no arquivo core-site.xml do Hadoop.
emrfs-site	Alterar as configurações do EMRFS.
spark-metrics	Alterar os valores no arquivo metrics.p roperties do Spark.

Versões 5.36.0 515

Classificações	Descrições
spark-defaults	Alterar os valores no arquivo spark-def aults.conf do Spark.
spark-env	Alterar os valores no ambiente do Spark.
spark-hive-site	Altera os valores no arquivo hive-site.xml do Spark.
spark-log4j	Alterar os valores no arquivo log4j.properties do Spark.

As classificações de configuração permitem que você personalize aplicações. Eles geralmente correspondem a um arquivo XML de configuração do aplicativo, como spark-hive-site .xml. Para obter mais informações, consulte <u>Configuração de aplicações</u>.

emr-5.36.0-latest

Notas de versão: a versão emr-5.36.0-latest direciona para emr-5.36.0-20240321, no momento.

Regiões: emr-5.36.0-latest está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-5.36.0:latest.

emr-5.36.0-20240321

Notas de versão: emr-5.36.0-20240321 foi lançada em 11 de março de 2024. Em comparação com a versão anterior, esta versão foi atualizada com os pacotes do Amazon Linux atualizados recentemente e com as correções críticas.

Regiões: emr-5.36.0-20240321 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-5.36.0:20240321.

emr-5.36.0-latest 516

emr-5.36.0-20221219

Notas de versão: emr-5.36.0-20221219 foi liberada em 27 de janeiro de 2023. Em comparação com a versão anterior, esta versão foi atualizada com os pacotes do Amazon Linux atualizados recentemente.

Regiões: emr-5.36.0-20221219 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-5.36.0:20221219.

emr-5.36.0-20220620

Notas de versão: emr-5.36.0-20220620 foi liberada em 27 de julho de 2022. Em comparação com a versão anterior, esta versão foi atualizada com os pacotes do Amazon Linux atualizados recentemente.

Regiões: emr-5.36.0-20220620 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-5.36.0:20220620.

emr-5.36.0-20220525

Notas de versão: emr-5.36.0-20220525 foi liberada em 16 de junho de 2022. Esta é a versão inicial do Amazon EMR 5.36.0.

Regiões: emr-5.36.0-20220525 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-5.36.0:20220525.

Versões 5.35.0 do Amazon EMR no EKS

As versões 5.35.0 do Amazon EMR apresentadas a seguir estão disponíveis para o Amazon EMR no EKS. Selecione uma versão específica do emr-5.35.0-XXXX para visualizar mais detalhes, como a etiqueta de imagem do contêiner relacionada.

- emr-5.35.0-latest
- emr-5.35.0-20240321
- emr-5.35.0-20221219

emr-5.36.0-20221219 517

- emr-5.35.0-20220802
- emr-5.35.0-20220307

Notas da versão 5.35.0 do Amazon EMR

- Os problemas de segurança do log4j2 foram corrigidos.
- Aplicações com suporte: Spark 2.4.8-amzn-1, Hudi 0.9.0-amzn-2 e Jupyter Enterprise Gateway (endpoints e pré-visualização pública; o kernel do Scala não tem suporte).
- Componentes suportados aws-hm-client (conector Glue), emr-s3-select aws-sagemaker-sparksdk, emrfs, emr-ddb, hudi-spark.
- Classificações de configuração com suporte:

Classificações	Descrições
core-site	Altera os valores no arquivo core-site.xml do Hadoop.
emrfs-site	Alterar as configurações do EMRFS.
spark-metrics	Alterar os valores no arquivo metrics.p roperties do Spark.
spark-defaults	Alterar os valores no arquivo spark-def aults.conf do Spark.
spark-env	Alterar os valores no ambiente do Spark.
spark-hive-site	Altera os valores no arquivo hive-site.xml do Spark.
spark-log4j	Alterar os valores no arquivo log4j.properties do Spark.

As classificações de configuração permitem que você personalize aplicações. Geralmente, eles correspondem a um arquivo XML de configuração do aplicativo, como spark-hive-site .xml. Para obter mais informações, consulte Configuração de aplicações.

Versões 5.35.0 518

emr-5.35.0-latest

Notas de versão: a versão emr-5.35.0-latest direciona para emr-5.35.0-20240321, no momento.

Regiões: emr-5.35.0-latest está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-5.35.0:latest.

emr-5.35.0-20240321

Notas de versão: emr-5.35.0-20240321 foi lançada em 11 de março de 2024. Em comparação com a versão anterior, esta versão foi atualizada com os pacotes do Amazon Linux atualizados recentemente e com as correções críticas.

Regiões: emr-5.35.0-20240321 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-5.35.0:20240321.

emr-5.35.0-20221219

Notas de versão: emr-5.35.0-20221219 foi liberada em 27 de janeiro de 2023. Em comparação com a versão anterior, esta versão foi atualizada com os pacotes do Amazon Linux atualizados recentemente.

Regiões: emr-5.35.0-20221219 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-5.35.0:20221219.

emr-5.35.0-20220802

Notas de versão: emr-5.35.0-20220802 foi liberada em 27 de setembro de 2022. Em comparação com a versão anterior, esta versão foi atualizada com os pacotes do Amazon Linux atualizados recentemente.

Regiões: emr-5.35.0-20220802 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

emr-5.35.0-latest 519

Etiqueta da imagem do contêiner: emr-5.35.0:20220802.

emr-5.35.0-20220307

Notas de versão: emr-5.35.0-20220307 foi liberada em 30 de março de 2022. Em comparação com a versão anterior, esta versão foi atualizada com os pacotes do Amazon Linux atualizados recentemente.

Regiões: emr-5.35.0-20220307 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-5.35.0:20220307.

Versões 5.34.0 do Amazon EMR no EKS

As versões 5.34.0 do Amazon EMR apresentadas a seguir estão disponíveis para o Amazon EMR no EKS. Selecione uma versão específica do emr-5.34.0-XXXX para visualizar mais detalhes, como a etiqueta de imagem do contêiner relacionada.

- emr-5.34.0-latest
- emr-5.34.0-20240321
- emr-5.34.0-20220802

Notas da versão 5.34.0 do Amazon EMR

- Aplicações com suporte: Spark 2.4.8-amzn-0 e Jupyter Enterprise Gateway (endpoints e prévisualização pública; o kernel do Scala não tem suporte).
- Componentes com suporte: aws-hm-client (conector do Glue), aws-sagemaker-spark-sdk, emr-s3-select, emrfs, emr-ddb e hudi-spark.
- Classificações de configuração com suporte:

Classificações	Descrições
core-site	Altera os valores no arquivo core-site.xml do Hadoop.
emrfs-site	Alterar as configurações do EMRFS.

emr-5.35.0-20220307 520

Classificações	Descrições
spark-metrics	Alterar os valores no arquivo metrics.p roperties do Spark.
spark-defaults	Alterar os valores no arquivo spark-def aults.conf do Spark.
spark-env	Alterar os valores no ambiente do Spark.
spark-hive-site	Altera os valores no arquivo hive-site.xml do Spark.
spark-log4j	Alterar os valores no arquivo log4j.properties do Spark.

As classificações de configuração permitem que você personalize aplicações. Geralmente, eles correspondem a um arquivo XML de configuração do aplicativo, como spark-hive-site .xml. Para obter mais informações, consulte Configuração de aplicações.

emr-5.34.0-latest

Notas de versão: a versão emr-5.34.0-latest direciona para emr-5.34.0-20220802, no momento.

Regiões: emr-5.34.0-latest está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-5.34.0:latest.

emr-5.34.0-20240321

Notas de versão: emr-5.34.0-20240321 foi lançada em 11 de março de 2024. Em comparação com a versão anterior, esta versão foi atualizada com os pacotes do Amazon Linux atualizados recentemente e com as correções críticas.

Regiões: emr-5.34.0-20240321 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

emr-5.34.0-latest 521

Etiqueta da imagem do contêiner: emr-5.34.0:20240321.

emr-5.34.0-20220802

Notas de versão: emr-5.34.0-20220802 foi liberada em 24 de agosto de 2022. Em comparação com a versão anterior, esta versão foi atualizada com os pacotes do Amazon Linux atualizados recentemente.

Regiões: emr-5.34.0-20220802 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-5.34.0:20220802.

emr-5.34.0-20211208

Notas de versão: emr-5.34.0-20211208 foi liberada em 20 de janeiro de 2022. Em comparação com a versão anterior, esta versão foi atualizada com os pacotes do Amazon Linux atualizados recentemente.

Regiões: emr-5.34.0-20211208 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-5.34.0:20211208.

Versões 5.33.0 do Amazon EMR no EKS

As versões 5.33.0 do Amazon EMR apresentadas a seguir estão disponíveis para o Amazon EMR no EKS. Selecione uma versão específica do emr-5.33.0-XXXX para visualizar mais detalhes, como a etiqueta de imagem do contêiner relacionada.

- emr-5.33.0-latest
- emr-5.33.0-20240321
- emr-5.33.0-20221219
- emr-5.33.0-20220802
- emr-5.33.0-20211008
- emr-5.33.0-20210802
- emr-5.33.0-20210615

emr-5.34.0-20220802 522

emr-5.33.0-20210323

Notas da versão 5.33.0 do Amazon EMR

- Novo recurso: começando com o Amazon EMR 5.33.0 na série de versões 5.x, o Amazon EMR no EKS oferece suporte ao recurso de modelo de pod do Spark. Para obter mais informações, consulte Uso de modelos de pod.
- Aplicações com suporte: Spark 2.4.7-amzn-1 e Jupyter Enterprise Gateway (endpoints e prévisualização pública; o kernel do Scala não tem suporte).
- Componentes com suporte: aws-hm-client (conector do Glue), aws-sagemaker-spark-sdk, emr-s3-select, emrfs, emr-ddb e hudi-spark.
- Classificações de configuração com suporte:

Classificações	Descrições
core-site	Altera os valores no arquivo core-site.xml do Hadoop.
emrfs-site	Alterar as configurações do EMRFS.
spark-metrics	Alterar os valores no arquivo metrics.p roperties do Spark.
spark-defaults	Alterar os valores no arquivo spark-def aults.conf do Spark.
spark-env	Alterar os valores no ambiente do Spark.
spark-hive-site	Altera os valores no arquivo hive-site.xml do Spark.
spark-log4j	Alterar os valores no arquivo log4j.properties do Spark.

As classificações de configuração permitem que você personalize aplicações. Geralmente, eles correspondem a um arquivo XML de configuração do aplicativo, como spark-hive-site .xml. Para obter mais informações, consulte Configuração de aplicações.

Versões 5.33.0 523

emr-5.33.0-latest

Notas de versão: a versão emr-5.33.0-latest direciona para emr-5.33.0-20240321, no momento.

Regiões: emr-5.33.0-latest está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-5.33.0:latest.

emr-5.33.0-20240321

Notas de versão: emr-5.33.0-20240321 foi lançada em 11 de março de 2024. Em comparação com a versão anterior, esta versão foi atualizada com os pacotes do Amazon Linux atualizados recentemente e com as correções críticas.

Regiões: emr-5.33.0-20240321 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-5.33.0:20240321.

emr-5.33.0-20221219

Notas de versão: emr-5.33.0-20221219 foi liberada em 19 de janeiro de 2023. Em comparação com a versão anterior, esta versão foi atualizada com os pacotes do Amazon Linux atualizados recentemente e com as correções críticas.

Regiões: emr-5.33.0-20221219 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-5.33.0:20221219.

emr-5.33.0-20220802

Notas de versão: emr-5.33.0-20220802 foi liberada em 24 de agosto de 2022. Em comparação com a versão anterior, esta versão foi atualizada com os pacotes do Amazon Linux atualizados recentemente.

Regiões: emr-5.33.0-20220802 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

emr-5.33.0-latest 524

Etiqueta da imagem do contêiner: emr-5.33.0:20220802.

emr-5.33.0-20211008

Notas de versão: emr-5.33.0-20211008 foi liberada em 9 de dezembro de 2021. Em comparação com a versão anterior, esta versão contém correções de problemas e atualizações de segurança.

Regiões: emr-5.33.0-20211008 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-5.33.0:20211008.

emr-5.33.0-20210802

Notas de versão: emr-5.33.0-20210802 foi liberada em 2 de agosto de 2021. Em comparação com a versão anterior, esta versão contém correções de problemas e atualizações de segurança.

Regiões: emr-5.33.0-20210802 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-5.33.0:20210802.

emr-5.33.0-20210615

Notas de versão: emr-5.33.0-20210615 foi liberada em 15 de junho de 2021. Em comparação com a versão anterior, esta versão contém correções de problemas e atualizações de segurança.

Regiões: emr-5.33.0-20210615 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-5.33.0:20210615.

emr-5.33.0-20210323

Notas de versão: emr-5.33.0-20210323 foi liberada em 23 de março de 2021. Esta é a versão inicial do Amazon EMR 5.33.0.

Regiões: emr-5.33.0-20210323 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

emr-5.33.0-20211008 525

Etiqueta da imagem do contêiner: emr-5.33.0-20210323.

Versões 5.32.0 do Amazon EMR no EKS

As versões 5.32.0 do Amazon EMR apresentadas a seguir estão disponíveis para o Amazon EMR no EKS. Selecione uma versão específica do emr-5.32.0-XXXX para visualizar mais detalhes, como a etiqueta de imagem do contêiner relacionada.

- emr-5.32.0-latest
- emr-5.32.0-20240321
- emr-5.32.0-20220802
- emr-5.32.0-20211008
- emr-5.32.0-20210802
- emr-5.32.0-20210615
- emr-5.32.0-20210129
- emr-5.32.0-20201218
- emr-5.32.0-20201201

Notas da versão 5.32.0 do Amazon EMR

- Aplicações com suporte: Spark 2.4.7-amzn-0 e Jupyter Enterprise Gateway (endpoints e prévisualização pública; o kernel do Scala não é compatível).
- Componentes com suporte: aws-hm-client (conector do Glue), aws-sagemaker-spark-sdk, emr-s3-select, emrfs, emr-ddb e hudi-spark.
- Classificações de configuração com suporte:

Classificações	Descrições
core-site	Altera os valores no arquivo core-site.xml do Hadoop.
emrfs-site	Alterar as configurações do EMRFS.
spark-metrics	Alterar os valores no arquivo metrics.p roperties do Spark.

Versões 5.32.0 526

Classificações	Descrições
spark-defaults	Alterar os valores no arquivo spark-def aults.conf do Spark.
spark-env	Alterar os valores no ambiente do Spark.
spark-hive-site	Altera os valores no arquivo hive-site.xml do Spark.
spark-log4j	Alterar os valores no arquivo log4j.properties do Spark.

As classificações de configuração permitem que você personalize aplicações. Geralmente, eles correspondem a um arquivo XML de configuração do aplicativo, como spark-hive-site .xml. Para obter mais informações, consulte <u>Configuração de aplicações</u>.

emr-5.32.0-latest

Notas de versão: a versão emr-5.32.0-latest direciona para emr-5.32.0-20240321, no momento.

Regiões: emr-5.32.0-latest está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-5.32.0:latest.

emr-5.32.0-20240321

Notas de versão: emr-5.32.0-20240321 foi lançada em 11 de março de 2024. Em comparação com a versão anterior, esta versão foi atualizada com os pacotes do Amazon Linux atualizados recentemente e com as correções críticas.

Regiões: emr-5.32.0-20240321 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-5.32.0:20240321.

emr-5.32.0-latest 527

emr-5.32.0-20220802

Notas de versão: emr-5.32.0-20220802 foi liberada em 24 de agosto de 2022. Em comparação com a versão anterior, esta versão foi atualizada com os pacotes do Amazon Linux atualizados recentemente.

Regiões: emr-5.32.0-20220802 está disponível em todas as regiões com suporte do Amazon EMR no EKS. Para obter mais informações, consulte Endpoints de serviço do Amazon EMR no EKS.

Etiqueta da imagem do contêiner: emr-5.32.0:20220802.

emr-5.32.0-20211008

Notas de versão: emr-5.32.0-20211008 foi liberada em 9 de dezembro de 2021. Em comparação com a versão anterior, esta versão contém correções de problemas e atualizações de segurança.

Regiões: emr-5.32.0-20211008 está disponível nas seguintes regiões: Leste dos EUA (Norte da Virgínia), Oeste dos EUA (Oregon), Ásia-Pacífico (Tóquio), Europa (Irlanda) e América do Sul (São Paulo).

Etiqueta da imagem do contêiner: emr-5.32.0:20211008.

emr-5.32.0-20210802

Notas de versão: emr-5.32.0-20210802 foi liberada em 2 de agosto de 2021. Em comparação com a versão anterior, esta versão contém correções de problemas e atualizações de segurança.

Regiões: emr-5.32.0-20210802 está disponível nas seguintes regiões: Leste dos EUA (Norte da Virgínia), Oeste dos EUA (Oregon), Ásia-Pacífico (Tóquio), Europa (Irlanda) e América do Sul (São Paulo).

Etiqueta da imagem do contêiner: emr-5.32.0:20210802.

emr-5.32.0-20210615

Notas de versão: emr-5.32.0-20210615 foi liberada em 15 de junho de 2021. Em comparação com a versão anterior, esta versão contém correções de problemas e atualizações de segurança.

Regiões: emr-5.32.0-20210615 está disponível nas seguintes regiões: Leste dos EUA (Norte da Virgínia), Oeste dos EUA (Oregon), Ásia-Pacífico (Tóquio), Europa (Irlanda) e América do Sul (São Paulo).

emr-5.32.0-20220802 528

Etiqueta da imagem do contêiner: emr-5.32.0:20210615.

emr-5.32.0-20210129

Notas de versão: emr-5.32.0-20210129 foi liberada em 29 de janeiro de 2021. Em comparação com emr-5.32.0-20201218, esta versão contém correções de problemas e atualizações de segurança.

Regiões: emr-5.32.0-20210129 está disponível nas seguintes regiões: Leste dos EUA (Norte da Virgínia), Oeste dos EUA (Oregon), Ásia-Pacífico (Tóquio), Europa (Irlanda) e América do Sul (São Paulo).

Etiqueta da imagem do contêiner: emr-5.32.0-20210129.

emr-5.32.0-20201218

Notas de versão: 5.32.0-20201218 foi liberada em 18 de dezembro de 2020. Em comparação com 5.32.0-20201201, esta versão contém correções de problemas e atualizações de segurança.

Regiões: emr-5.32.0-20201218 está disponível nas seguintes regiões: Leste dos EUA (Norte da Virgínia), Oeste dos EUA (Oregon), Ásia-Pacífico (Tóquio), Europa (Irlanda) e América do Sul (São Paulo).

Etiqueta da imagem do contêiner: emr-5.32.0-20201218.

emr-5.32.0-20201201

Notas de versão: 5.32.0-20201201 foi liberada em 1.º de dezembro de 2020. Esta é a versão inicial do Amazon EMR 5.32.0.

Regiões: 5.32.0-20201201 está disponível nas seguintes regiões: Leste dos EUA (Norte da Virgínia), Oeste dos EUA (Oregon), Ásia-Pacífico (Tóquio), Europa (Irlanda) e América do Sul (São Paulo).

Etiqueta da imagem do contêiner: emr-5.32.0-20201201.

emr-5.32.0-20210129 529

Histórico do documento

A tabela a seguir descreve as alterações importantes na documentação desde a última versão do Amazon EMR no EKS. Para obter mais informações sobre as atualizações desta documentação, você pode se tornar assinante de um feed RSS.

Alteração	Descrição	Data
Atualização de conteúdo	Políticas gerenciadas para o Amazon EMR no EKS— Permissões adicionai s paraAmazonEMRContainer sServiceRolePolicy	3 de fevereiro de 2025
Nova versão	Amazon EMR nas versões do EKS 7.6.0	10 de janeiro de 2025
Nova versão	Amazon EMR nas versões do EKS 7.5.0	21 de novembro de 2024
Nova versão	Amazon EMR nas versões do EKS 7.4.0	13 de novembro de 2024
Nova versão	Versões 7.3.0 do Amazon EMR no EKS	16 de outubro de 2024
Nova versão	Versões 7.2.0 do Amazon EMR no EKS	25 de julho de 2024
Nova versão	Versões 7.1.0 do Amazon EMR no EKS	17 de abril de 2024
Nova versão	Versões 7.0.0 do Amazon EMR no EKS	22 de dezembro de 2023
Nova versão	Versões 6.15.0 do Amazon EMR no EKS	17 de novembro de 2023
Nova versão	Versões 6.14.0 do Amazon EMR no EKS	17 de outubro de 2023
Atualização de conteúdo	Renomeação de "endpoints gerenciad os" para endpoints interativos. Os	29 de setembro de 2023

Alteração	Descrição	Data
	endpoints interativos estão com disponibilidade geral.	
Nova versão	Versões 6.13.0 do Amazon EMR no EKS e documentos de pré-visualização pública para Execução de trabalhos do Flink com o Amazon EMR no EKS.	12 de setembro de 2023
Nova versão	Versões 6.12.0 do Amazon EMR no EKS	21 de julho de 2023
Novo conteúdo	Uso do Volcano como um programad or personalizado para Apache Spark no Amazon EMR no EKS adicionado	13 de junho de 2023
Novo conteúdo	Uso do Volcano como um programad or personalizado para Apache Spark no Amazon EMR no EKS adicionado	13 de junho de 2023
Novo conteúdo	Uso da alternância de log do contêiner do Spark adicionado	12 de junho de 2023
Atualização de conteúdo	Atualização da documentação de imagens personalizadas para a descoberta de informações sobre as imagens base na galeria pública do Amazon ECR.	8 de junho de 2023
Nova versão	Versões 6.11.0 do Amazon EMR no EKS	8 de junho de 2023
Novo conteúdo	Adição de Execução de trabalhos do Spark com o operador do Spark e reorganização das seções de execuções de trabalhos em Execução de trabalhos do Spark com o Amazon EMR no EKS.	5 de junho de 2023

Alteração	Descrição	Data
Novo conteúdo	Adição de duas seções: <u>Uso da</u> escalabilidade automática vertical com trabalhos do Spark no Amazon EMR e Uso de cadernos Jupyter de hospedage m própria.	4 de maio de 2023
Página para o histórico de documentos	Criação de uma página para o histórico de documentos para o Amazon EMR no EKS.	13 de março de 2023
Página para as políticas gerenciadas	Criação de uma página para as políticas gerenciadas para o Amazon EMR no EKS.	13 de março de 2023

As traduções são geradas por tradução automática. Em caso de conflito entre o conteúdo da tradução e da versão original em inglês, a versão em inglês prevalecerá.