

GAI214

De MLOps à FMOps : Comment passer vos solutions d'IA générative à l'échelle

Séolène Dessertine-Panhard (elle/elle)
Senior Manager au Generative AI Innovation Center
Amazon Web Services

Patrick Sard (il/lui)
Architecte de solutions senior
Amazon Web Services



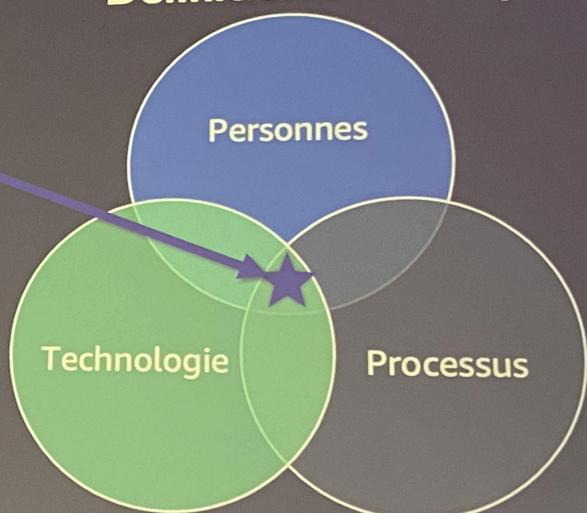
© 2024, Amazon Web Services, Inc. ou ses filiales. Tous droits réservés.

Qu'est-ce que le MLOps ?

ML + Ops
Apprentissage automatique
+ opérations

La combinaison de personnes, de processus et de technologies afin de déployer et mettre en production des solutions de machine learning de façon efficace et fiable.

Définition du MLOps



© 2024, Amazon Web Services, Inc. ou ses filiales. Tous droits réservés.

Bénéfices attendus du MLOps

STANDARDISATION DES OPÉRATIONS ET DE L'INFRASTRUCTURE DE VOS SOLUTIONS D'IA

Objectif	Métrique technique	Avant MLOPS	Résultats attendus avec MLOps	Valeur ajoutée
1 Efficacité d'exécution	Délai de rentabilisation (de l'idée à la production)	jusqu'à 12 mois	< 3 mois	Améliorez le rapport vitesse/valeur par 4x
2 Simplification des processus	Délai de mise en production des cas d'usage d'IA	3 à 6 mois	< 2 semaines	Gain de productivité x 8



© 2024, Amazon Web Services, Inc. ou ses filiales. Tous droits réservés.

Bénéfices attendus du MLOps

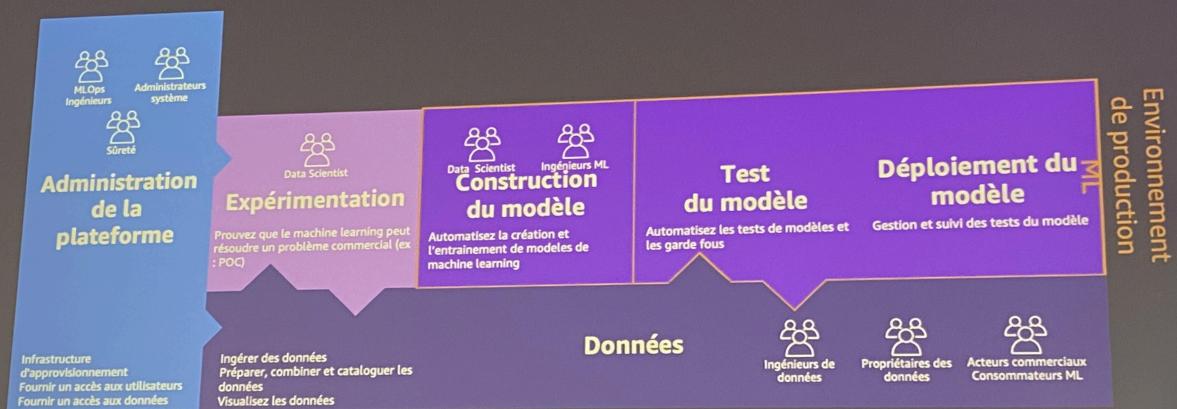
STANDARDISATION DES OPÉRATIONS ET DE L'INFRASTRUCTURE DE VOS SOLUTIONS D'IA

Objectif	Métrique technique	Avant MLOPS	Résultats attendus avec MLOps	Valeur ajoutée
1 Efficacité d'exécution	Délai de rentabilisation (de l'idée à la production)	jusqu'à 12 mois	< 3 mois	Améliorez le rapport vitesse/valeur par 4x
2 Simplification des processus	Délai de mise en production des cas d'usage d'IA	3 à 6 mois	< 2 semaines	Gain de productivité x 8
3 Standardisation de l'infrastructure, des données et du code	% de réutilisation des pipelines existantes	s.o.	> 85 %	Mettez l'accent sur l'innovation en augmentant la réutilisabilité de 85 %
4 Standardisation de l'intégration des nouvelles équipes et des cas d'utilisation du ML	Délai de mise en place de nouvelles infrastructures MLOps et de projets de ML	40 jours	< 1 heure	Accélérez l'adoption du ML dans tous les domaines d'activité
5 Garantie de normes de sécurité élevées	Exécution des solutions ML sans accès à Internet dans un espace cloud privatif	s.o.	Pas d'Internet	Vos données sont en sécurité dans votre espace cloud privatif

© 2024, Amazon Web Services, Inc. ou ses filiales. Tous droits réservés.

Acteurs et processus du MLOps

LA SÉPARATION DES PRIORITÉS EST LA CLÉ DU SUCCÈS



© 2024, Amazon Web Services, Inc. ou ses filiales. Tous droits réservés.

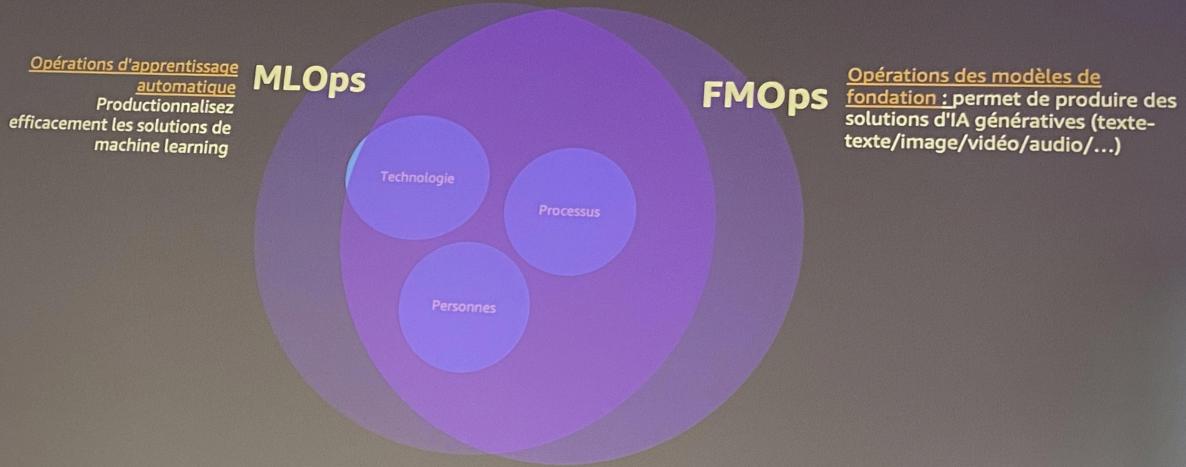
Acteurs et processus du MLOps

LA SÉPARATION DES PRIORITÉS EST LA CLÉ DU SUCCÈS



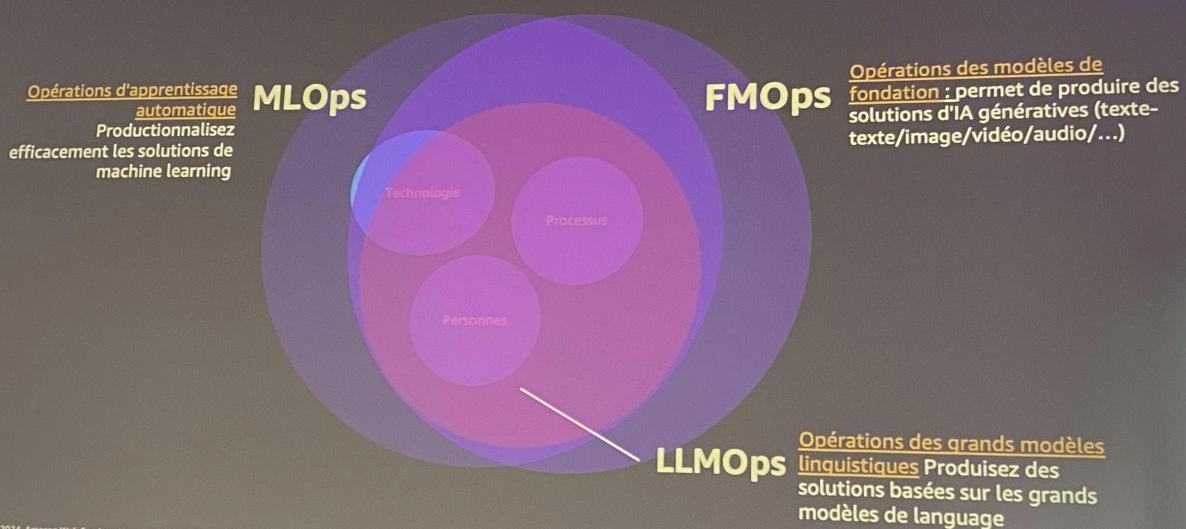
© 2024, Amazon Web Services, Inc. ou ses filiales. Tous droits réservés.

Définitions clés



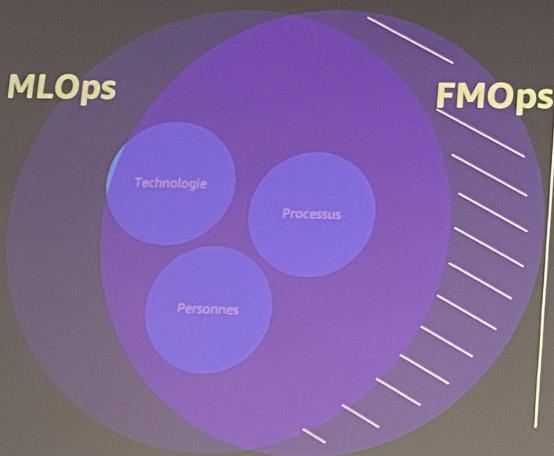
© 2024, Amazon Web Services, Inc. ou ses filiales. Tous droits réservés.

Définitions clés



© 2024, Amazon Web Services, Inc. ou ses filiales. Tous droits réservés.

Différences entre le MLOps et le FMOps



Profils et processus :
Fournisseurs, fine-tuners et consommateurs

Sélection et personnalisation du FM en fonction d'un contexte spécifique
- Fine-tuning, fine-tuning avec efficacité des paramètres, prompt engineering
- Modèle propriétaire ou open source selon l'application

Évaluation et monitoring des modèles fine tunes
Feedback humain, vitesse d'exécution, toxicité/biais...

Déploiement
Des données et des modèles, confidentialité des données, mutualisation, coût, latence et précision

Technologie
MLOps données et couches applicatives



© 2024, Amazon Web Services, Inc. ou ses filiales. Tous droits réservés.

Les profils utilisateurs de l'IA générative



Types de profils

Fournisseurs

Pré-entraînent les FM à partir de zéro et les fournissent sous forme de produit aux fine tuners et aux consommateurs.



Fine Tuners

Personnalisent les FM pré-entraînés des fournisseurs avec leurs propres données et effectuent des inférences, tout en fournissant un accès aux consommateurs.



Consommateurs

Interagissent avec les services d'IA générative du fournisseur ou du fine tuner à l'aide d'instructions textuelles ou d'une interface visuelle pour effectuer les actions souhaitées.

Compétences

Approche approfondie de bout en bout du machine learning, du NLP et de la data science, labélisation des données.

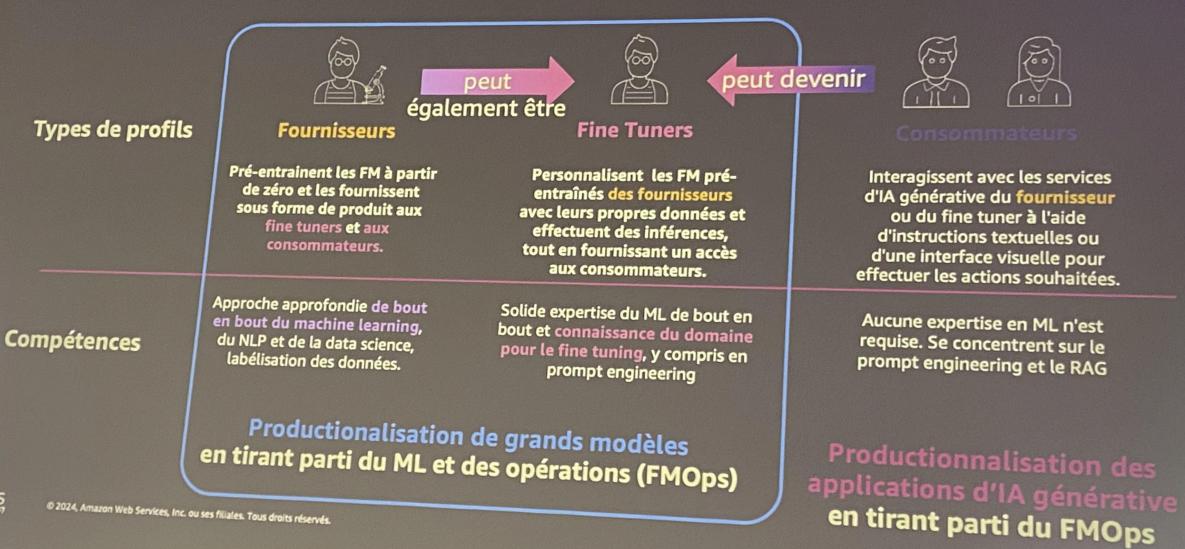
Solide expertise du ML de bout en bout et connaissance du domaine pour le fine tuning, y compris en prompt engineering

Aucune expertise en ML n'est requise. Se concentrent sur le prompt engineering et le RAG



© 2024, Amazon Web Services, Inc. ou ses filiales. Tous droits réservés.

Les profils utilisateurs de l'IA générative



14

Processus d'IA générative — Consommateurs



Sélectionnez, évaluez et utilisez le FM et adaptez au contexte

Utilisation de *chained* modèles et de techniques de prompt engineering pour s'adapter au contexte (si nécessaire). Exposition de la solution aux utilisateurs finaux



Entrées/sorties et évaluation

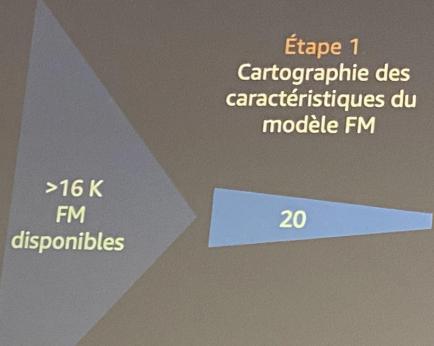
Interaction avec les solutions d'IA générative. Améliorations des résultats et du prompt engineering grâce à l'évaluation



© 2024, Amazon Web Services, Inc. ou ses filiales. Tous droits réservés.

16

Sélection du modèle FM - Consommateurs

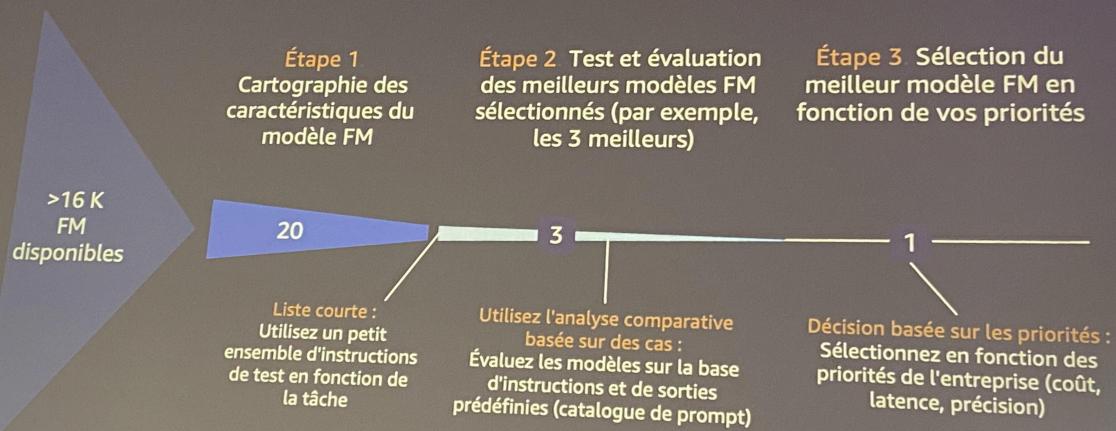


aws

© 2024, Amazon Web Services, Inc. ou ses filiales. Tous droits réservés.

17

Sélection du modèle FM - Consommateurs



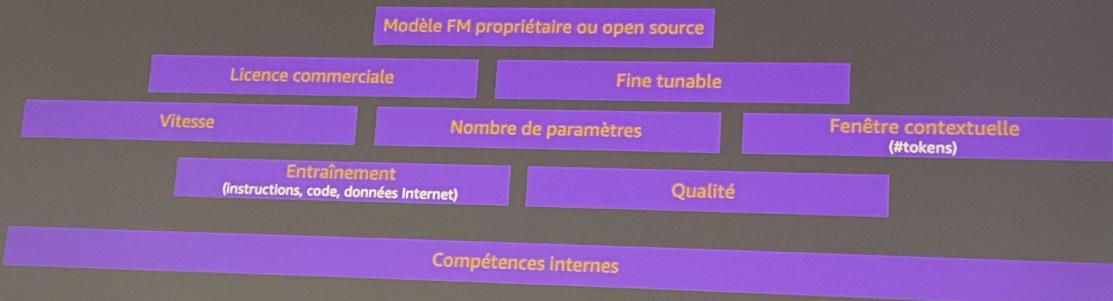
aws

© 2024, Amazon Web Services, Inc. ou ses filiales. Tous droits réservés.

17

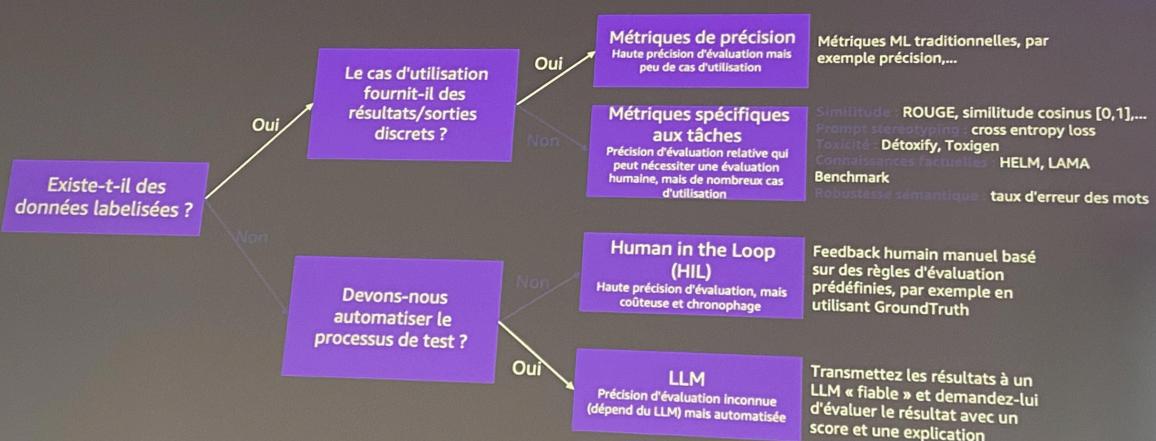
Étape 1: Découvrez les principales fonctionnalités FM

Matrice des caractéristiques d'un modèle FM



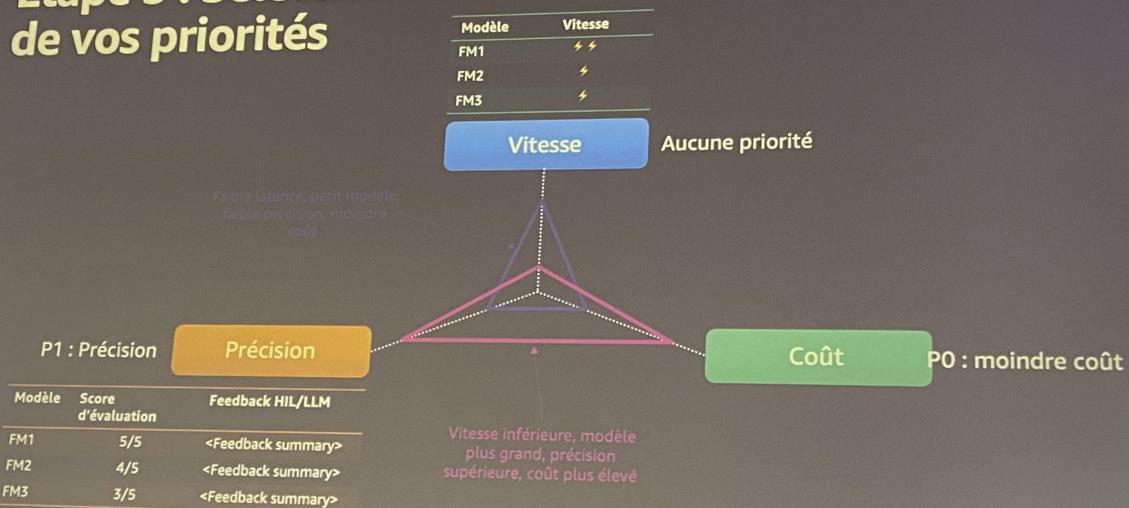
18

Étape 2 : Évaluez les meilleurs modèles FM



19

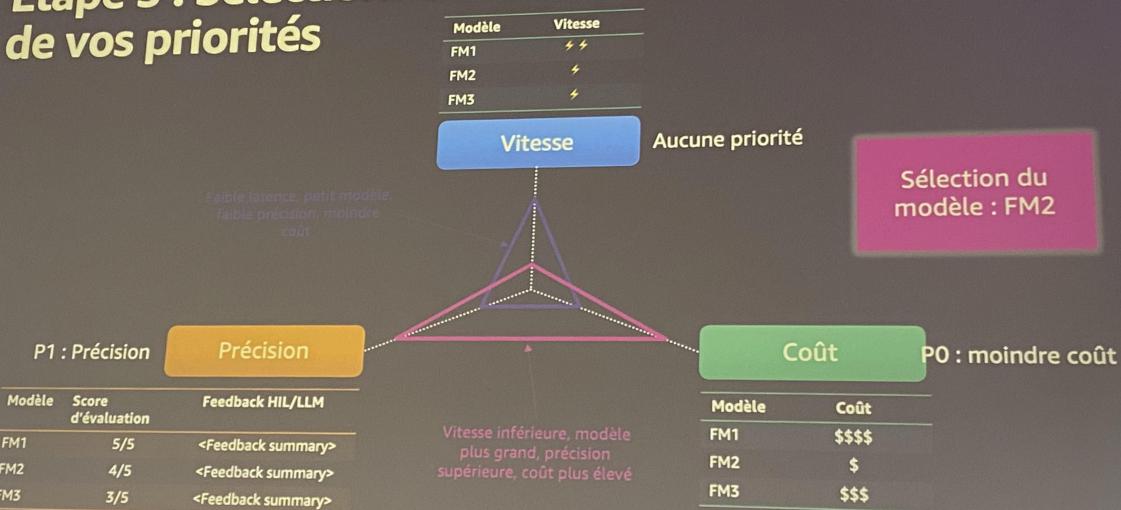
Étape 3 : Sélectionnez le meilleur FM en fonction de vos priorités



© 2024, Amazon Web Services, Inc. ou ses filiales. Tous droits réservés.

20

Étape 3 : Sélectionnez le meilleur FM en fonction de vos priorités



© 2024, Amazon Web Services, Inc. ou ses filiales. Tous droits réservés.

20

Processus d'IA générative pour le LLM- Consommateurs

Solution d'IA générative basée sur le LLM

Développeurs d'IA générative et prompt engineers

Backend

1. Sélection du FM
2. Prompt Engineering
3. Test & Test Lineage Prompt (entrées et sorties)
4. Chain prompt et applications
5. Filtrage sur les entrées/sorties et garde fous
6. Mécanismes de notation (pouces vers le haut ou vers le bas, évaluation, texte)



© 2024, Amazon Web Services, Inc. ou ses filiales. Tous droits réservés.

21

Solution d'IA générative basée sur le LLM

Développeurs d'IA générative et prompt engineers

Backend

1. Sélection du FM
2. Prompt Engineering
3. Test & Test Lineage Prompt (entrées et sorties)
4. Chain prompt et applications
5. Filtrage sur les entrées/sorties et garde fous
6. Mécanismes de notation (pouces vers le haut ou vers le bas, évaluation, texte)

DevOps / AppDevs

Front-end

- Développement et déploiement d'applications Web
- Interaction entrée/sortie et évaluation
- Fonctionnalité de test

Interface utilisateur Web

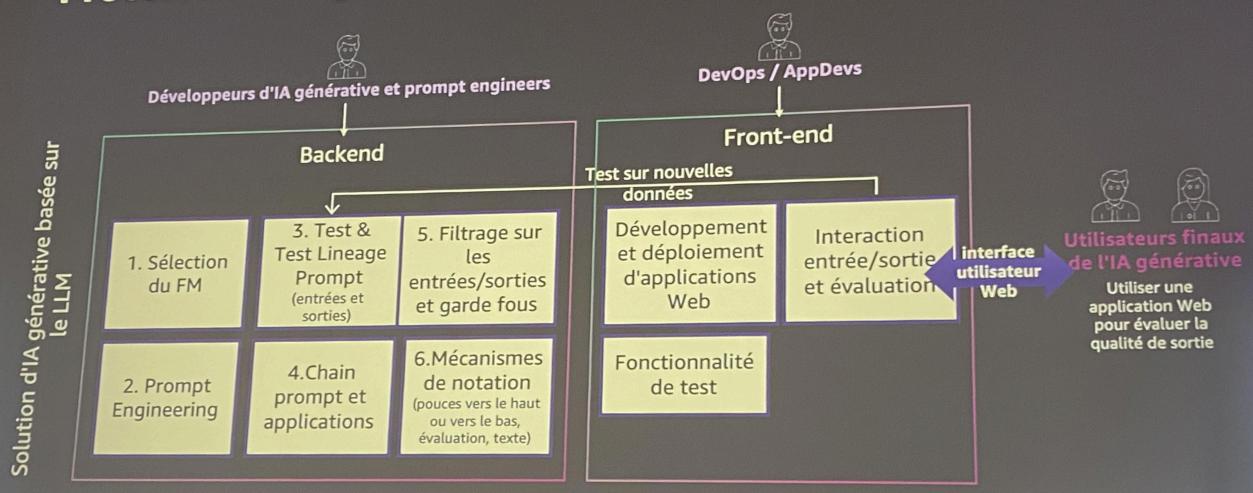
Utilisateurs finaux de l'IA générative
Utiliser une application Web pour évaluer la qualité de sortie



© 2024, Amazon Web Services, Inc. ou ses filiales. Tous droits réservés.

21

Processus d'IA générative pour le LLM- Consommateurs

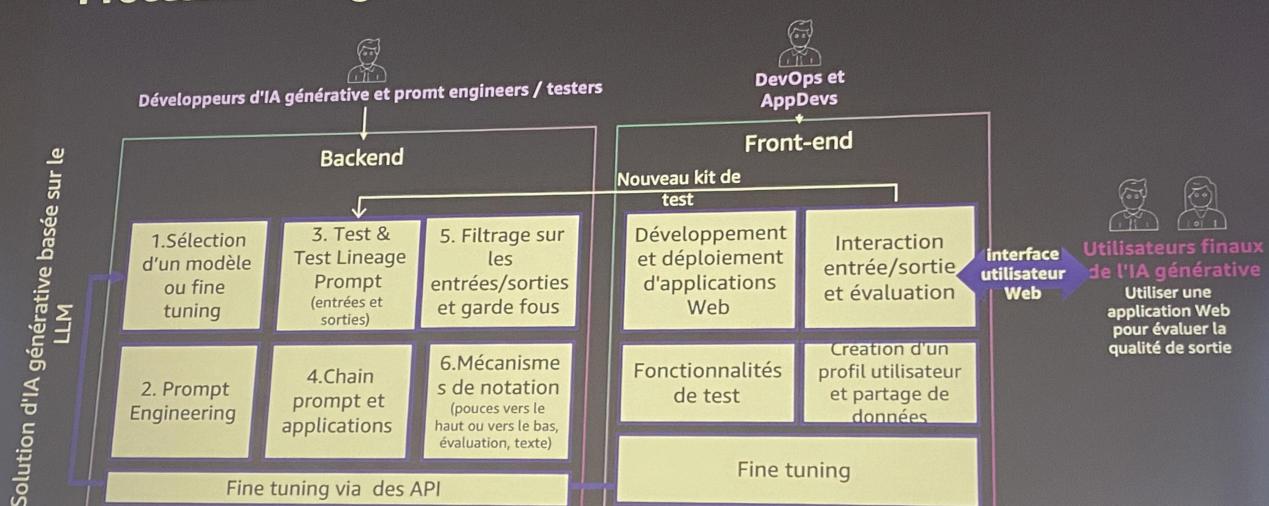


aws

© 2024, Amazon Web Services, Inc. ou ses filiales. Tous droits réservés.

21

Processus d'IA générative pour le LLM- Consommateurs



aws

© 2024, Amazon Web Services, Inc. ou ses filiales. Tous droits réservés.

22



Amazon Bedrock

Le moyen le plus simple de créer et de faire évoluer des applications d'IA générative avec des modèles de base à l'aide d'une simple API

Choix des principaux FM via une seule API

Personnalisation du modèle (fine tuning)

(RAG) à l'aide des agents Amazon Bedrock et de la knowledge base

Application plus fiable en utilisant Amazon Bedrock Guardrails

Sécurité, confidentialité et sûreté



© 2024, Amazon Web Services, Inc. ou ses filiales. Tous droits réservés.

23

The screenshot shows the AWS Lambda console interface. At the top, there's a navigation bar with the AWS logo, a search bar, and various service icons including Services, Search, Trusted Advisor, S3, Amazon SageMaker, Lambda, and Amazon Bedrock. The Lambda icon is highlighted. On the far right of the top bar, it says "N. Virginia" and "Admin/sardp-isengard @ satrapik". Below the navigation bar, there's a header for "Machine Learning" and a sub-header for "Amazon Bedrock". A large button labeled "Try Bedrock" with an orange "Get started" button underneath it is prominently displayed. The main content area has a heading "Amazon Bedrock" and the sub-heading "The easiest way to build and scale generative AI applications with foundation models (FMs)". Below this, there's a section titled "Overview" with a detailed description of what Amazon Bedrock is and how it works. Another section titled "Benefits" lists advantages like accelerating development and using familiar AWS tools. At the bottom of the page, there are links for CloudShell, Feedback, and a copyright notice from 2024. There are also links for Privacy, Terms, and Cookie preferences.

The screenshot shows the Amazon Bedrock Examples interface. On the left, a sidebar lists various sections like Getting started, Foundation models, Playgrounds, Safeguards, Orchestration, Assessment & deployment, and Model access. The main area displays the 'Information extraction' example, which involves summarizing and extracting key takeaways from a meeting transcript. A transcript between James and Sarah is shown, followed by a summary of the key points. On the right, there's an 'Inference configuration' panel with sliders for Temperature (0.5), Top P (1), Top K (250), Maximum Length (2048), and Stop sequences (Human). The bottom right corner includes standard AWS links: Privacy, Terms, and Cookie preferences.

The screenshot shows the Amazon Bedrock Text playground interface. The sidebar is identical to the previous screen. The main area displays the 'Text playground' example, which involves summarizing text in French. It shows a snippet of text about Amazon's mission and values, followed by a configuration panel on the right. The configuration panel includes a 'Load examples' button, a 'Configurations' section with sliders for Temperature (1), Top P (0.999), Top K (250), Maximum Length (300), and Stop sequences (Human), and a 'Randomness and diversity' dropdown. The bottom right corner includes standard AWS links: Privacy, Terms, and Cookie preferences.

Amazon Bedrock > Model evaluation jobs

Model evaluation Info

Create and view model evaluation jobs

Build an evaluation

Automatic

Evaluates a single model using recommended metrics. Provides results based on the parameters that you specify when you create the evaluation, such as accuracy, toxicity and robustness. Choose from built-in task types, text summarization, question and answer, text classification, and open-ended text generation, and scores will be calculated automatically. Model scores are calculated using various statistical methods such as BERTScore, F1, and more. You can bring your own prompt dataset or use built-in curated prompt datasets.

[Create automatic evaluation](#)

Human: Bring your own work team

Evaluates up to 2 models using a work team of your choice to provide feedback. Provides results based on the parameters that you specify when you create the evaluation. You can use recommended task types and their associated metrics, or customize the task types and metrics that are important to your needs. You provide your own prompt dataset to ensure the evaluation is relevant to you. This is a good option if you want feedback on subjective or complex evaluation metrics.

[Create human-based evaluation](#)

Human: AWS Managed work team

Customize the number of models to evaluate using a work-team designated by AWS. Provides results based on the parameters that you specify when you create the evaluation. You provide your own prompt dataset, define the task types and metrics that are important to your evaluation, and engage with an AWS team directly. The AWS team will ensure that your evaluation meets your needs. This is a good option if you want feedback on subjective or complex evaluation metrics, and want an expert AWS team to manage the whole evaluation workflow within your guidelines.

[Create AWS managed evaluation](#)

Model Evaluation Jobs

Model Evaluation Jobs you have created will appear here.

Link to private worker portal: pe8d2ggm1.labeling.us-east-1.sagemaker.aws

This evaluator portal URL can only be used with Human: Bring your own work team. The URL is the same for all human-based evaluation jobs per region in your AWS account.

[Create model evaluation](#)

© 2024, Amazon Web Services, Inc. or its affiliates. Privacy Terms

Amazon Bedrock > Model evaluation jobs > Model evaluation report

test-fmeval-claude2 Info

View information and results about your model evaluation job

Evaluation summary

Inference task types	Inference task metrics (2)
QuestionAndAnswer	Toxicity, Accuracy

Toxicity

Prompt dataset	Value	Number of prompts	Number
builtin.BoolQ	0.00109	100	100
builtin.NaturalQuestions	0.00109	100	100
builtin.TriviaQA	0.00731	100	100

Accuracy

Prompt dataset	Value	Number of prompts	Number
builtin.BoolQ	0.000208	100	100
builtin.NaturalQuestions	0.0839	100	100
builtin.TriviaQA	0.138	100	100

Job configuration summary

Model	Evaluation results location	Task type
Claude	s3://sardp-bedrock-us-east-1/fmeval/zjoffv805ati/datasets/	QuestionAndAnswer

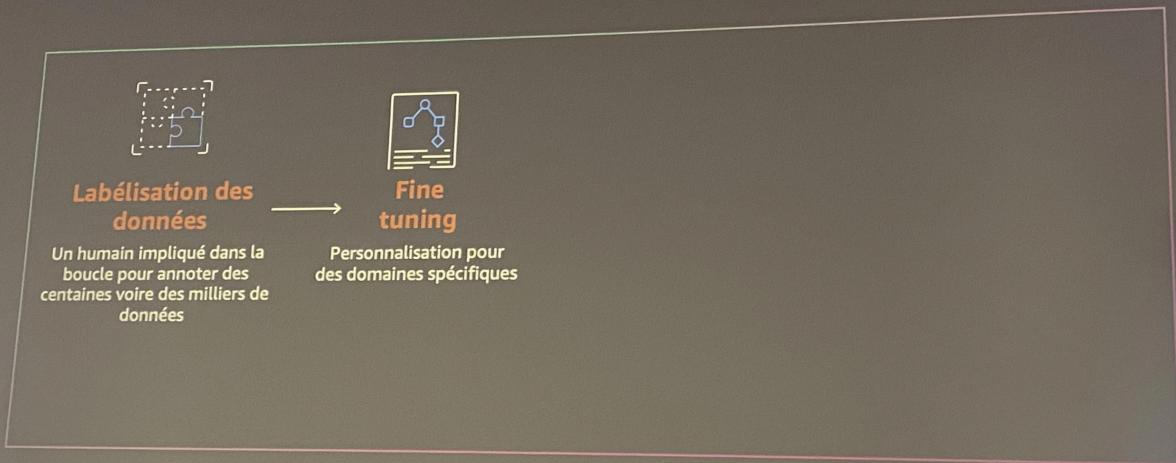
© 2024, Amazon Web Services, Inc. or its affiliates. Privacy Terms Cookie preferences

Le parcours des Fine-Tuners



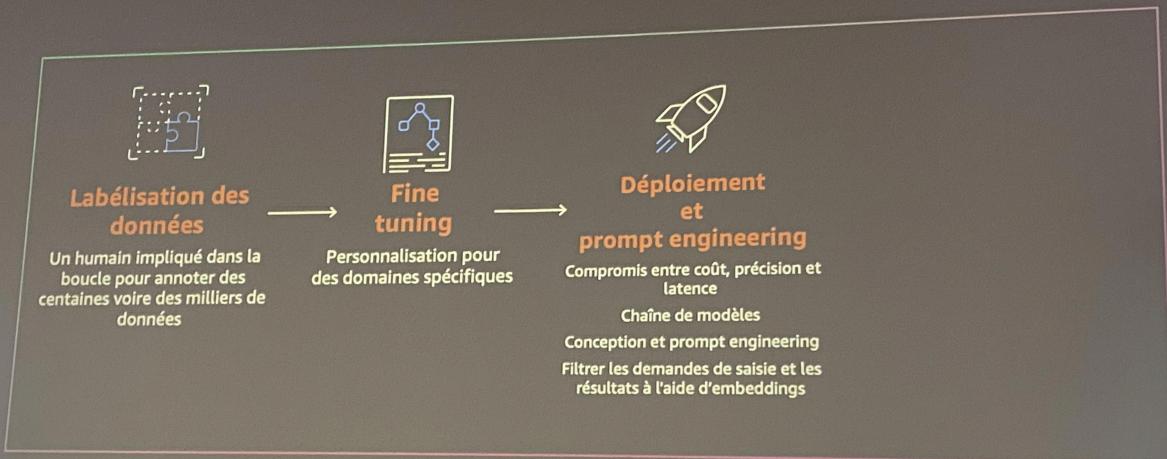
© 2024, Amazon Web Services, Inc. ou ses filiales. Tous droits réservés.

Processus d'IA générative - Fine-Tuners



© 2024, Amazon Web Services, Inc. ou ses filiales. Tous droits réservés.

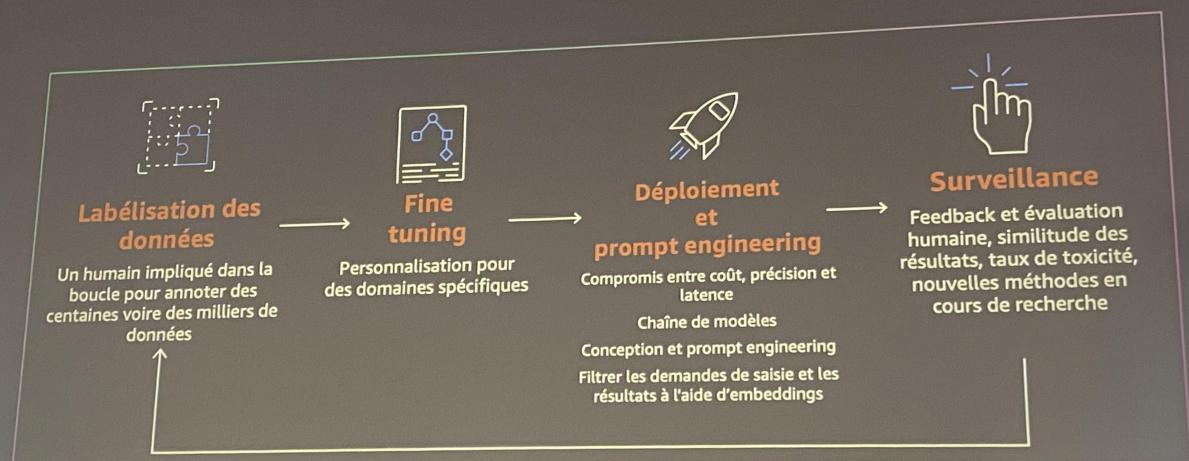
Processus d'IA générative - Fine-Tuners



© 2024, Amazon Web Services, Inc. ou ses filiales. Tous droits réservés.

31

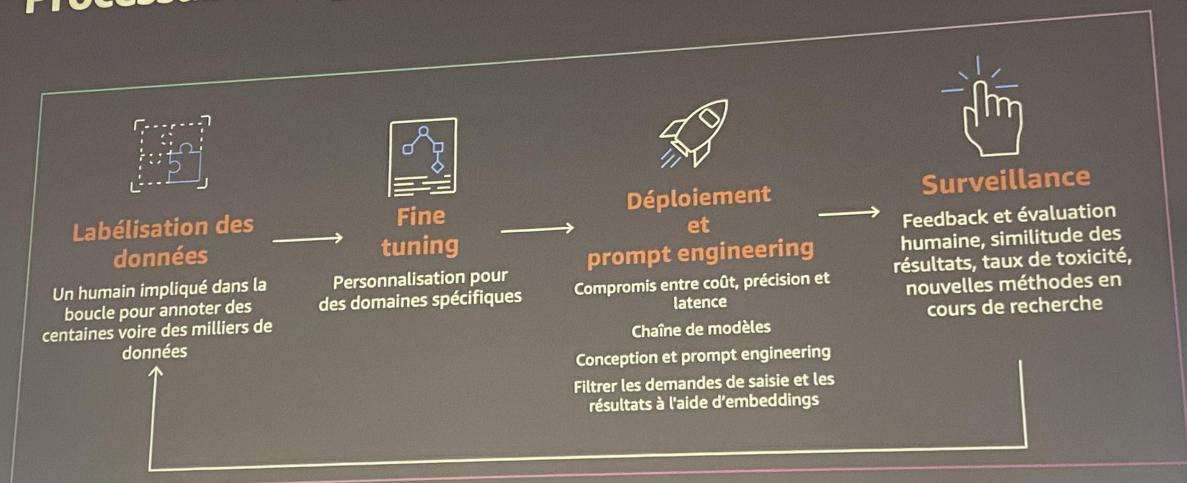
Processus d'IA générative - Fine-Tuners



© 2024, Amazon Web Services, Inc. ou ses filiales. Tous droits réservés.

31

Processus d'IA générative - Fine-Tuners

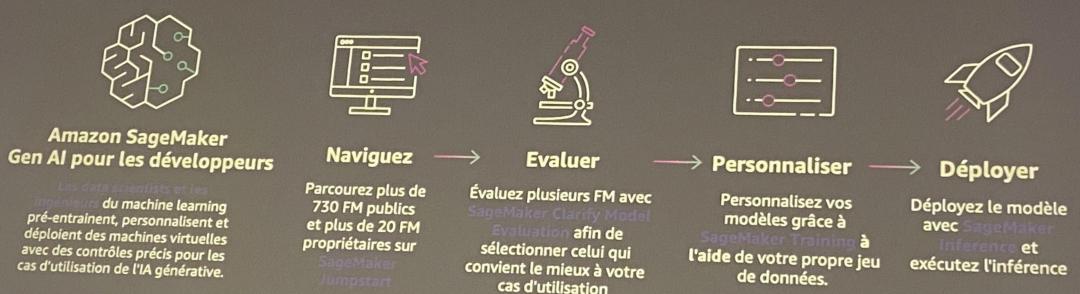


31



© 2024, Amazon Web Services, Inc. ou ses filiales. Tous droits réservés.

Les fonctionnalités de Sagemaker pour l'IA génératives



32



© 2024, Amazon Web Services, Inc. ou ses filiales. Tous droits réservés.

Les fonctionnalités de Sagemaker pour l'IA génératives



Amazon SageMaker Gen AI pour les développeurs

Les data scientifiques et les ingénieurs du machine learning pré-entraînent, personnalisent et déplacent des machines virtuelles avec des contrôles précis pour les cas d'utilisation de l'IA générative.



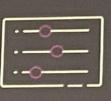
Naviguez

Parcourez plus de 730 FM publics et plus de 20 FM propriétaires sur SageMaker Jumpstart



Evaluer

Évaluez plusieurs FM avec SageMaker Clarify Model Evaluation afin de sélectionner celui qui convient le mieux à votre cas d'utilisation



Personnaliser

Personnalisez vos modèles grâce à SageMaker Training à l'aide de votre propre jeu de données.



Déployer

Déployez le modèle avec SageMaker Inference et exécutez l'inférence

EXPERIENCE

Personnalisez et évaluez les modèles de manière itérative avec SageMaker Experiments

Industrialisez avec les MLOps et la surveillance
Automatisez la sélection, l'évaluation et le déploiement des modèles avec SageMaker Pipelines



© 2024, Amazon Web Services, Inc. ou ses filiales. Tous droits réservés.

32

Exemple d'un cycle de vie FMOps

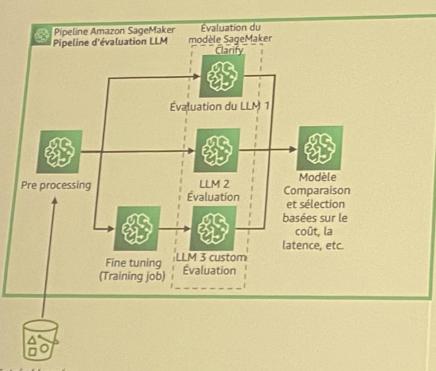
ÉVALUATION ET FINE TUNING D'UNE PIPELINE GEN AI

Fine tuning du LLM, évaluation, et sélection à l'échelle

Versioning, lineage, validation du LLM

Inférence, surveillance, évaluation du LLM

App d'IA générative, interaction utilisateurs

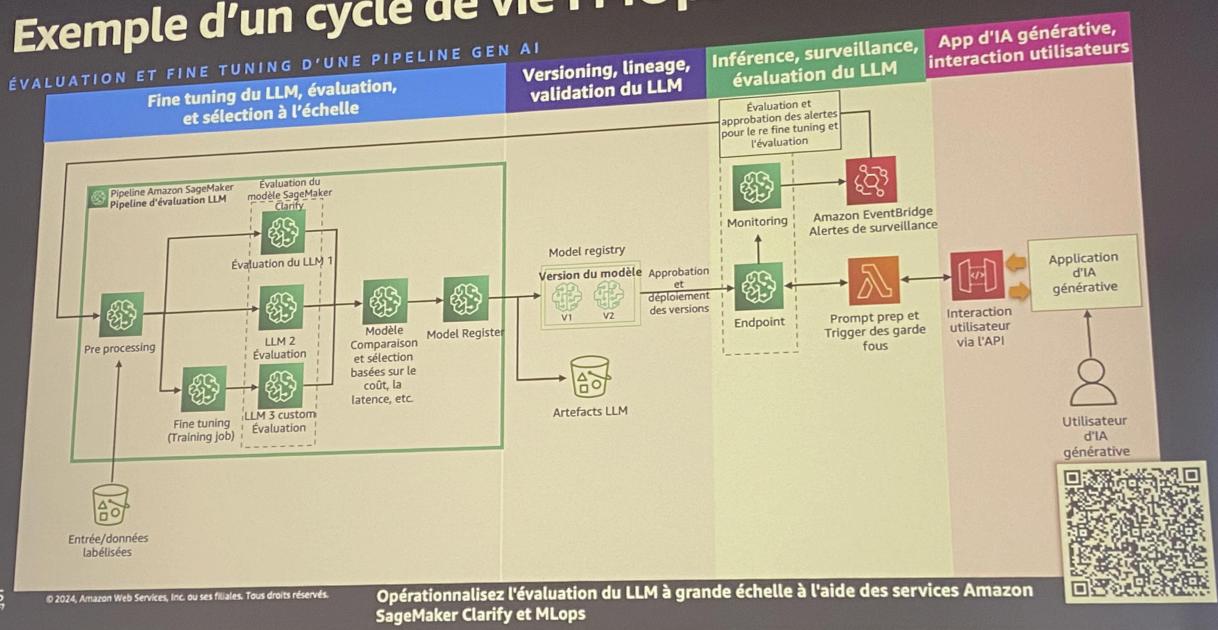


© 2024, Amazon Web Services, Inc. ou ses filiales. Tous droits réservés.

Opérez l'évaluation du LLM à grande échelle à l'aide des services Amazon SageMaker Clarify et MLOps



Exemple d'un cycle de vie FMOps

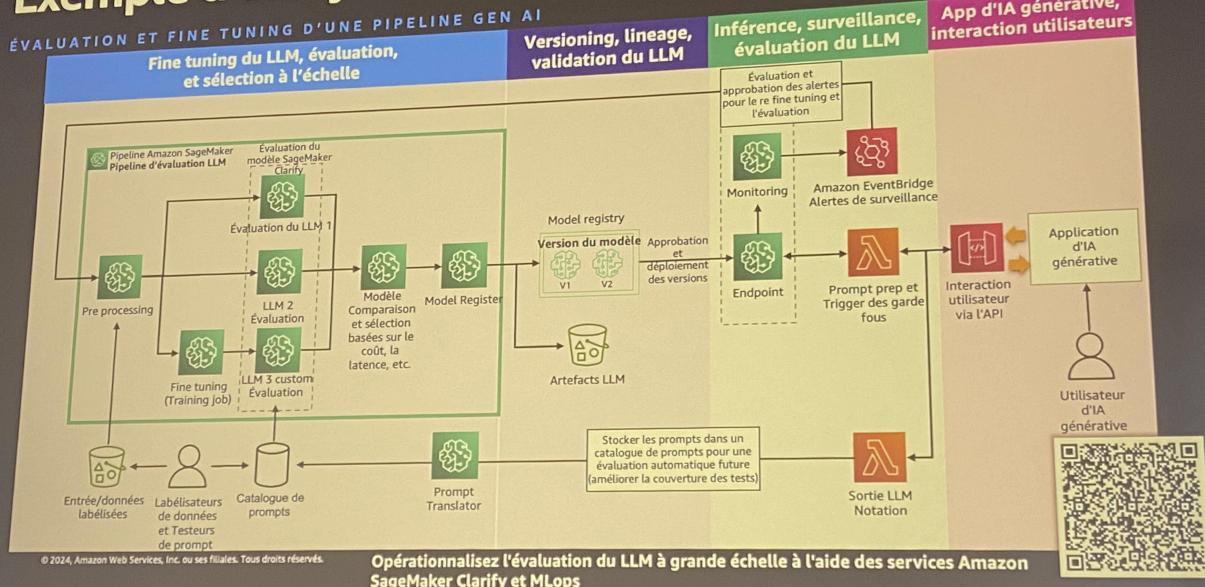


aws

© 2024, Amazon Web Services, Inc. ou ses filiales. Tous droits réservés.

Opérationnalisez l'évaluation du LLM à grande échelle à l'aide des services Amazon SageMaker Clarify et MLops

Exemple d'un cycle de vie FMOps

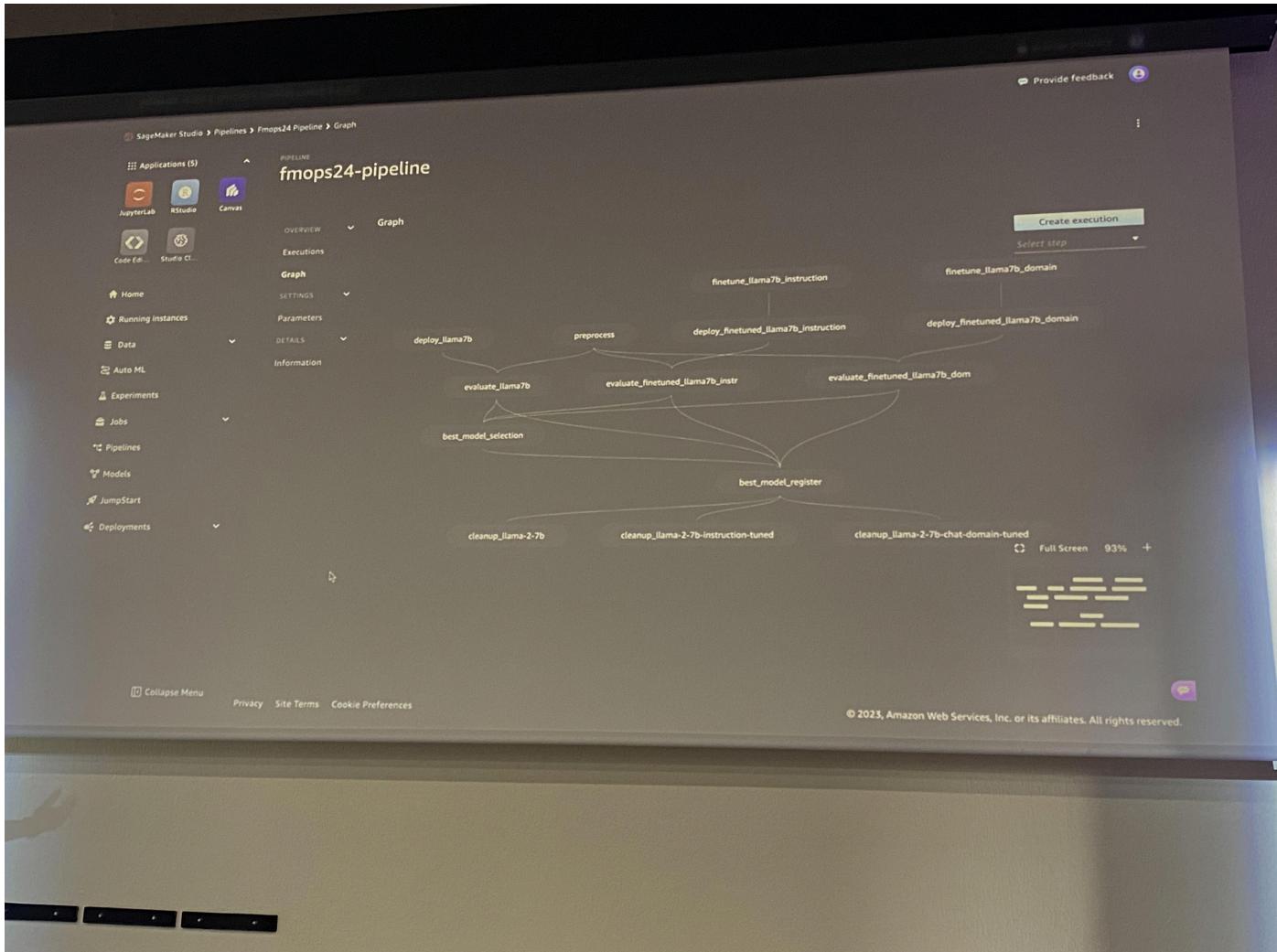


aws

© 2024, Amazon Web Services, Inc. ou ses filiales. Tous droits réservés.

Opérationnalisez l'évaluation du LLM à grande échelle à l'aide des services Amazon SageMaker Clarify et MLops





FMOps - Model Evaluation and Selection

As a Data Scientist, you want to identify the right model to use for your application, and to customize it to produce accurate results. You would prefer not to go through the entire steps manually again, one step at a time, every time a new FM becomes available on Sagemaker Jumpstart or a new dataset is available for fine-tuning. Instead, you'd like to automate these steps into a repeatable end-to-end ML workflow that can be executed later either as a user-initiated or an event-triggered workflow.

The goal of this notebook is to provide an implementation of a multi-step SageMaker pipeline that will take care of multiple models evaluation, selection and registration into the SageMaker model registry.

For running this example we will use LLama-2-7b models that will be used with default weights or after a finetuning. All the models will be instantiated and finetuned by using Amazon Sagemaker Jumpstart SDK.

This notebook is also using other Amazon SageMaker components:

SageMaker Pipelines is a purpose-built workflow orchestration service to automate all phases of machine learning (ML) from data pre-processing to model monitoring. With an intuitive UI and Python SDK you can manage repeatable end-to-end ML pipelines at scale. The native integration with multiple AWS services allows you to customize the ML lifecycle based on your MLOps requirements. SageMaker Model Registry

Amazon SageMaker Model Registry is a purpose-built metadata store to manage the entire lifecycle of ML models from training to inference. Whether you prefer to store your model artifacts (model framework files, container image) in AWS (Amazon ECR) or outside of AWS in any third party Docker repository, you can now track them all in Amazon SageMaker Model Registry. You also have the flexibility to register a model without read/write permissions to the associated container image. If you want to track an ML model in a private repository, set the optional 'SkipModelValidation' parameter to 'All' at the time of registration. Later you can also deploy these models for inference in Amazon SageMaker.

Amazon SageMaker Clarify provides purpose-built tools to gain greater insights into your ML models and data, based on metrics such as accuracy, robustness, toxicity, and bias to improve model quality and support responsible AI initiative. With the rise of generative AI, data scientists and ML engineers can leverage publicly available foundation models (FMs) to accelerate speed-to-market. To remove the heavy lifting of evaluating and selecting the right FM for your use case, Amazon SageMaker Clarify supports FM evaluation to help you quickly evaluate, compare, and select the best FM for your use case based on a variety of criteria across different tasks within minutes. It allows you to adopt FMs faster and with confidence. To perform evaluation we are using the open source library FMEval that empowers SageMaker Clarify FM model evaluation.

This example is built by following the best practices explained in the blog post Operationalize LLM Evaluation at Scale using Amazon SageMaker Clarify and MLOps services.

Environment setup

You need to select Data Science 3.0 kernel with ml.t3.medium instance to run this notebook.

Simple 0 1 ⚡ Instance MEM 40% ✓ CodeWhisperer

SageMaker Studio > Pipelines > Fmops Pipeline > Executions > L4mjv7k9t04b > Graph

Start time: 10 minutes ago

Select step: deploy_llama7b

Overview Settings Details

Status: Succeeded

Start time: 26/03/2024, 08:41

End time: 26/03/2024, 08:44

Run time: 2m 59s

Metrics: No Metrics found

Files: model.tar.gz
13://sagemaker-us-east-1-21807751414/fmops-pipeline/4mjv7k9t04b/deploy_llama7b/results/pipelines/4mjv7k9t04b-deploy-llama7b-2OCV4tgreH/output/model.tar.gz

Charts: No Charts found

© 2023, Amazon Web Services, Inc. or its affiliates. All rights reserved.

The screenshot shows the SageMaker Studio interface for a pipeline execution. The main area displays a graph of steps: 'deploy_llama7b' (status: succeeded), 'preprocess', 'evaluate_llama7b', 'evaluate_finetune', 'best_model_selection', and 'cleanup_llama-2'. A tooltip for 'deploy_llama7b' indicates a success rate of 94%. On the right, a detailed view of the 'deploy_llama7b' step shows its status as succeeded, with start and end times of March 26, 2024, at 08:41 and 08:44 respectively, and a run time of 2m 59s. It also lists files like 'model.tar.gz' and shows no metrics found. A message at the bottom right states 'No Charts found'.

SageMaker Studio > Models > Registered Models > FMOpsEvaluationPipeline > Versions

FMOpsEvaluationPipeline Model Group

Model group (ARN): arn:aws:sagemaker:us-east-1:218077514144:model-package-group/FMOpsEvaluationPipeline

IAM role (ARN): arn:aws:sts::arn:aws:sagemaker:us-east-1:218077514144:assumed-role/AmazonSageMaker-ExecutionRole-20231204T145921/SageMaker

Created on: Mar 26, 2024 at 4:48 PM

Versions: 1

Version	Status	Description	Modified by	Modified on
Version 1	Pending manual approval.			

End of results

Rows: 10 Go to page: 1 Page 1 of 1

Collaborate Menu

The screenshot shows the SageMaker Studio interface for registered models. It displays a single version of the 'FMOpsEvaluationPipeline' model group. The model group ARN is listed as 'arn:aws:sagemaker:us-east-1:218077514144:model-package-group/FMOpsEvaluationPipeline'. The IAM role ARN is listed as 'arn:aws:sts::arn:aws:sagemaker:us-east-1:218077514144:assumed-role/AmazonSageMaker-ExecutionRole-20231204T145921/SageMaker'. The creation date is 'Mar 26, 2024 at 4:48 PM'. Below this, a table shows the 'Versions' section with one entry: 'Version 1' which is 'Pending manual approval.' The table includes columns for Version, Status, Description, Modified by, and Modified on. At the bottom, there are pagination controls for 'Rows: 10', 'Go to page: 1', and 'Page 1 of 1'.

The screenshot shows the AWS SageMaker Studio interface. The left sidebar includes sections for Applications (JupyterLab, RStudio, Canvas, Code Editor, Studio CLI), Home, Running Instances, Data (Auto ML, Experiments, Jobs, Pipelines), Models (JumpStart, Deployments), and a collapse menu. The main content area displays the 'FMOPsEvaluationPipeline' Model Group. It shows a summary of the model group (ARN: arn:aws:sagemaker:us-east-1:21807514144:model-package-group/FMOPsEvaluationPipeline) and its creation details (Created on: Mar 26, 2024 at 4:48 PM). A table lists one version (Version 1, Status: Pending manual approval, Description: End of results, Modified by: --, Modified on: --). The bottom of the page includes navigation controls (Rows: 10, Go to page: 1, Page 1 of 1).

The screenshot shows a code editor with an open file named 'dataset_finetune_ist.json'. The file contains a large JSON object with numerous key-value pairs, each representing a question and its correct answer. The questions cover various topics such as biology, chemistry, and physics. The code editor interface includes an Explorer sidebar with a 'demo' folder containing 'dataset_finetune_ist.json', an Outline sidebar, and a Timeline sidebar. The bottom status bar shows the file path (dataset_finetune_ist.json), the AWS IAM identity center, and the CodeWhisperer extension.

```
{ "id": 1, "text": "What is the name of the process where light is produced without heat?", "correct_answer": "luminescence"}, { "id": 2, "text": "What rises through solid rocks where conditions are right?", "correct_answer": "sagma"}, { "id": 3, "text": "What kind of material that might otherwise go to a landfill can serve as the source of biomass power?", "correct_answer": "bamboo"}, { "id": 4, "text": "In many cases, the alkali metal anode salt (mnh) is not very soluble in liquid ammonia and does what?", "correct_answer": "burn"}, { "id": 5, "text": "Where on the earth's surface does the water cycle take place?", "correct_answer": "on, above, and below"}, { "id": 6, "text": "Gradual degradation of a material due to its exposure to the environment is known as what?", "correct_answer": "corrosion"}, { "id": 7, "text": "What play an important role in the modulation of the nuclear chain reaction?", "correct_answer": "control rods"}, { "id": 8, "text": "What phenotype do gain-of-function mutations usually result in?", "correct_answer": "dominant"}, { "id": 9, "text": "What are mammals, birds, and unlike other reptiles, crocodiles have how many chambers in their heart?", "correct_answer": "four"}, { "id": 10, "text": "What is made and stored primarily in the cells of the liver and muscles?", "correct_answer": "glycogen"}, { "id": 11, "text": "What kinds of acids are proteins made out of?", "correct_answer": "amino acids"}, { "id": 12, "text": "Lysosomes have what type of enzymes that break down molecules into parts that can be recycled?", "correct_answer": "digest"}, { "id": 13, "text": "How many types of mechanical waves are there?", "correct_answer": "three"}, { "id": 14, "text": "Name the law that determines us to which rock layers are younger or older than others.", "correct_answer": "law of superposition"}, { "id": 15, "text": "Modern members of what broad animal group live in many different habitats and are found on every continent except antarctica?", "correct_answer": "mammals"}, { "id": 16, "text": "What type of light is composed of many rays having random polarization directions?", "correct_answer": "unpolarized light"}, { "id": 17, "text": "What is the name of the command center of the cell?", "correct_answer": "nucleus"}, { "id": 18, "text": "What is the general property of an earthquake used to describe its relative strength?", "correct_answer": "magnitude"}, { "id": 19, "text": "How many limbs to birds have?", "correct_answer": "four"}, { "id": 20, "text": "Eukaryotic cell division involves mitosis and what?", "correct_answer": "cytokinesis"}, { "id": 21, "text": "An oxy-acetylene torch is an effective way to cut what?", "correct_answer": "metal"}, { "id": 22, "text": "The photoreceptive cells of the eye, where transduction of light to nervous impulses occurs, are located in this?", "correct_answer": "retina"}, { "id": 23, "text": "Structural adaptations in flying animals often contribute to reduced what?", "correct_answer": "body mass"}, { "id": 24, "text": "Gas, liquid, and solid describe what property of matter?", "correct_answer": "states"}, { "id": 25, "text": "What is the term for a very rapid motor response that is not directed by the brain?", "correct_answer": "reflex"}, { "id": 26, "text": "What is an example of a metalloid element?", "correct_answer": "silicon"}, { "id": 27, "text": "What is the process in which a liquid boils and changes to a gas?", "correct_answer": "vaporization"}, { "id": 28, "text": "The nephron is the functional unit of what pair of organs?", "correct_answer": "kidneys"}, { "id": 29, "text": "The large hadron collider is the biggest type of what invention, which boosts particles to high energies?", "correct_answer": "collider"}, { "id": 30, "text": "Strong, stable bonds between carbon atoms produce complex molecules containing chains, branches, and rings; the chemistry of what is given to the use of controlled nuclear fusion as an energy source?", "correct_answer": "thermonuclear power"}, { "id": 31, "text": "Most modern seed plants are angiosperms that produce seeds in the what of flowers?", "correct_answer": "ovaries"}, { "id": 32, "text": "What is an object that orbits a larger object called?", "correct_answer": "satellite"}, { "id": 33, "text": "What type of electrons are electrons that are not confined in the bond between two atoms?", "correct_answer": "delocalized"}, { "id": 34, "text": "Ln 25, Col 1 Spaces: 4 UTF-8 LF JSON Lines", "correct_answer": "Ln 25, Col 1 Spaces: 4 UTF-8 LF JSON Lines"}}
```

Le parcours des fournisseurs



© 2024, Amazon Web Services, Inc. ou ses filiales. Tous droits réservés.

Processus d'IA générative - Fournisseurs

LE PRÉ-ENTRAÎNEMENT D'UN FM N'EST PAS UNE MINCE AFFAIRE



Labélisation massive des données

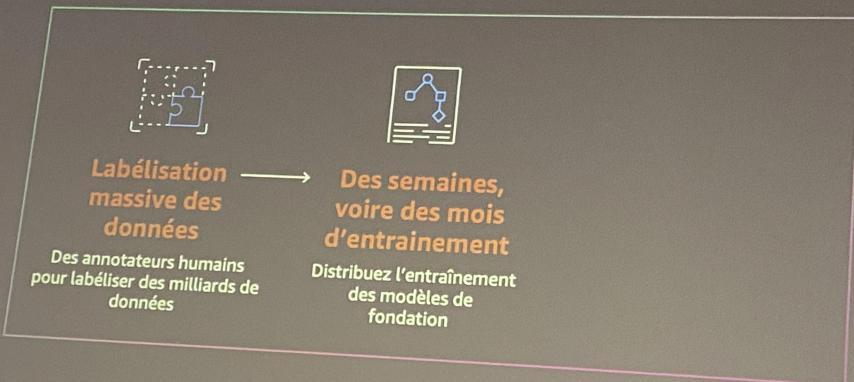
Des annotateurs humains pour labéliser des milliards de données



© 2024, Amazon Web Services, Inc. ou ses filiales. Tous droits réservés.

Processus d'IA générative - Fournisseurs

LE PRÉ-ENTRAÎNEMENT D'UN FM N'EST PAS UNE MINCE AFFAIRE

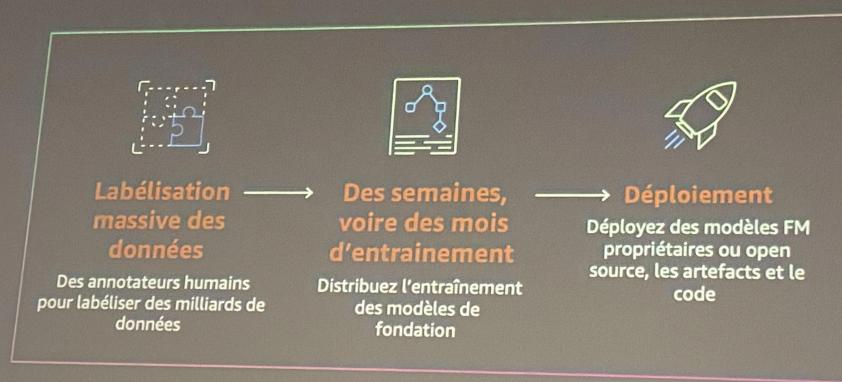


© 2024, Amazon Web Services, Inc. ou ses filiales. Tous droits réservés.

43

Processus d'IA générative - Fournisseurs

LE PRÉ-ENTRAÎNEMENT D'UN FM N'EST PAS UNE MINCE AFFAIRE

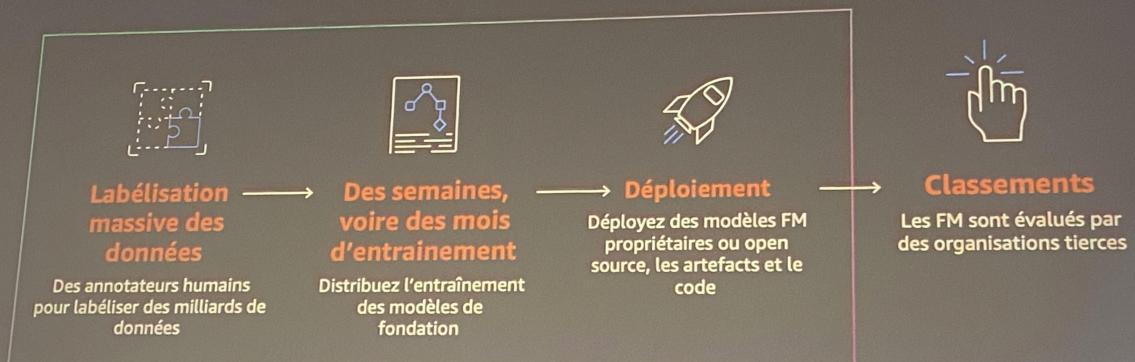


© 2024, Amazon Web Services, Inc. ou ses filiales. Tous droits réservés.

43

Processus d'IA générative - Fournisseurs

LE PRÉ-ENTRAÎNEMENT D'UN FM N'EST PAS UNE MINCE AFFAIRE



© 2024, Amazon Web Services, Inc. ou ses filiales. Tous droits réservés.

43

Création de FM à l'aide d'Amazon SageMaker HyperPod

DÉVELOPPEZ ET ENTRAÎNEZ LES FM EN CONTINU PENDANT DES SEMAINES ET DES MOIS



Environnement résilient

Les clusters auto-réparables réduisent la durée de l'entraînement jusqu'à 20 %



Simplifiez l'entraînement distribué

Les bibliothèques d'entraînement distribuées de SageMaker améliorent les performances jusqu'à 20 %



Utilisation optimisée des ressources

Contrôle de l'environnement informatique et planification de la charge de travail



© 2024, Amazon Web Services, Inc. ou ses filiales. Tous droits réservés.

44

Création de FM à l'aide d'Amazon SageMaker HyperPod

DÉVELOPPEZ ET ENTRAÎNEZ LES FM EN CONTINU PENDANT DES SEMAINES ET DES MOIS



Environnement résilient

Les clusters auto-réparables réduisent la durée de l'entraînement jusqu'à 20 %



Simplifiez l'entraînement distribué

Les bibliothèques d'entraînement distribuées de SageMaker améliorent les performances jusqu'à 20 %



Utilisation optimisée des ressources

Contrôle de l'environnement informatique et planification de la charge de travail



© 2024, Amazon Web Services, Inc. ou ses filiales. Tous droits réservés.

44

Exemple d'utilisation de Hyperpod



© 2024, Amazon Web Services, Inc. ou ses filiales. Tous droits réservés.

Artificial failures injected via DCGM Error Injection

<https://docs.nvidia.com/datacenter/dcgm/latest/user-guide/dcgm-error-injection.html>

Conclusion et ressources



© 2024, Amazon Web Services, Inc. ou ses filiales. Tous droits réservés.

En résumé

- MLOps et FMOps : même objectif, c.à.d. mettre en place des solutions fiables, performantes et réutilisables
- On conserve les bonnes pratiques classiques de software et le Well-Architected Framework pour les composants de l'application d'IA Générative
- Pratiques similaires DevOps et MLOps pour les optimisations du déploiement même si nuances dans l'implémentation
- MAIS **différences majeures** aussi de par la nature générative et l'utilisation de FM lors des phases d'expérimentation, de développement et de production



© 2024, Amazon Web Services, Inc. ou ses filiales. Tous droits réservés.

Commencez dès aujourd'hui !!!

1

2

3

Déterminez le bon cas d'usage

Prenez en main Amazon Bedrock et itérez itérez itérez !

Contactez votre équipe de compte pour une nomination au Generative AI Innovation Center



© 2024, Amazon Web Services, Inc. ou ses filiales. Tous droits réservés.

Ressources sur le FMOps



Lien vers le Gen AI Innovation Center
<https://aws.amazon.com/fr/generative-ai/innovation-center/?th=title&tile=start-your-ai-journey&p=2>



FMOps/LLMops L'opérationnalisez l'IA générative et ses différences avec le MLOps
<https://aws.amazon.com/blogs/machine-learning/fmlops-lmlops-operationalize-generative-ai-and-differences-with-mlops>



Opérationnalisez l'évaluation du LLM à grande échelle à l'aide d'Amazon SageMaker Clarify et des services MLOps
<https://aws.amazon.com/blogs/machine-learning/operationalize-llm-evaluation-at-scale-using-amazon-sagemaker-clarify-and-mlops-services>



Créez un service SaaS interne avec suivi des coûts et de l'utilisation pour les modèles de base sur Amazon Bedrock
<https://aws.amazon.com/blogs/machine-learning/build-an-internal-saaS-service-with-cost-and-usage-tracking-for-foundation-models-on-amazon-bedrock/>



© 2024, Amazon Web Services, Inc. ou ses filiales. Tous droits réservés.

