

SESSION GA1302

## Améliorez vos applications d'IA générative avec RAG sur Amazon Bedrock

**Hugues Gendre**  
Chief Information Officer  
UCPA

**Merieme Ezzaouia**  
Solutions Architect  
AWS

**Nizar Kheir**  
Senior Solutions Architect  
AWS

AWS

© 2024, Amazon Web Services, Inc. or its affiliates. All rights reserved.



SESSION GA1302

# Améliorez vos applications d'IA générative avec RAG sur Amazon Bedrock

**Hugues Gendre**

Chief Information Officer  
UCPA

**Merieme Ezzaouia**

Solutions Architect  
AWS

**Nizar Kheir**

Senior Solutions Architect  
AWS



© 2024, Amazon Web Services, Inc. or its affiliates. All rights reserved.



**Pourquoi personnaliser un *Foundation Model (FM)* ?**

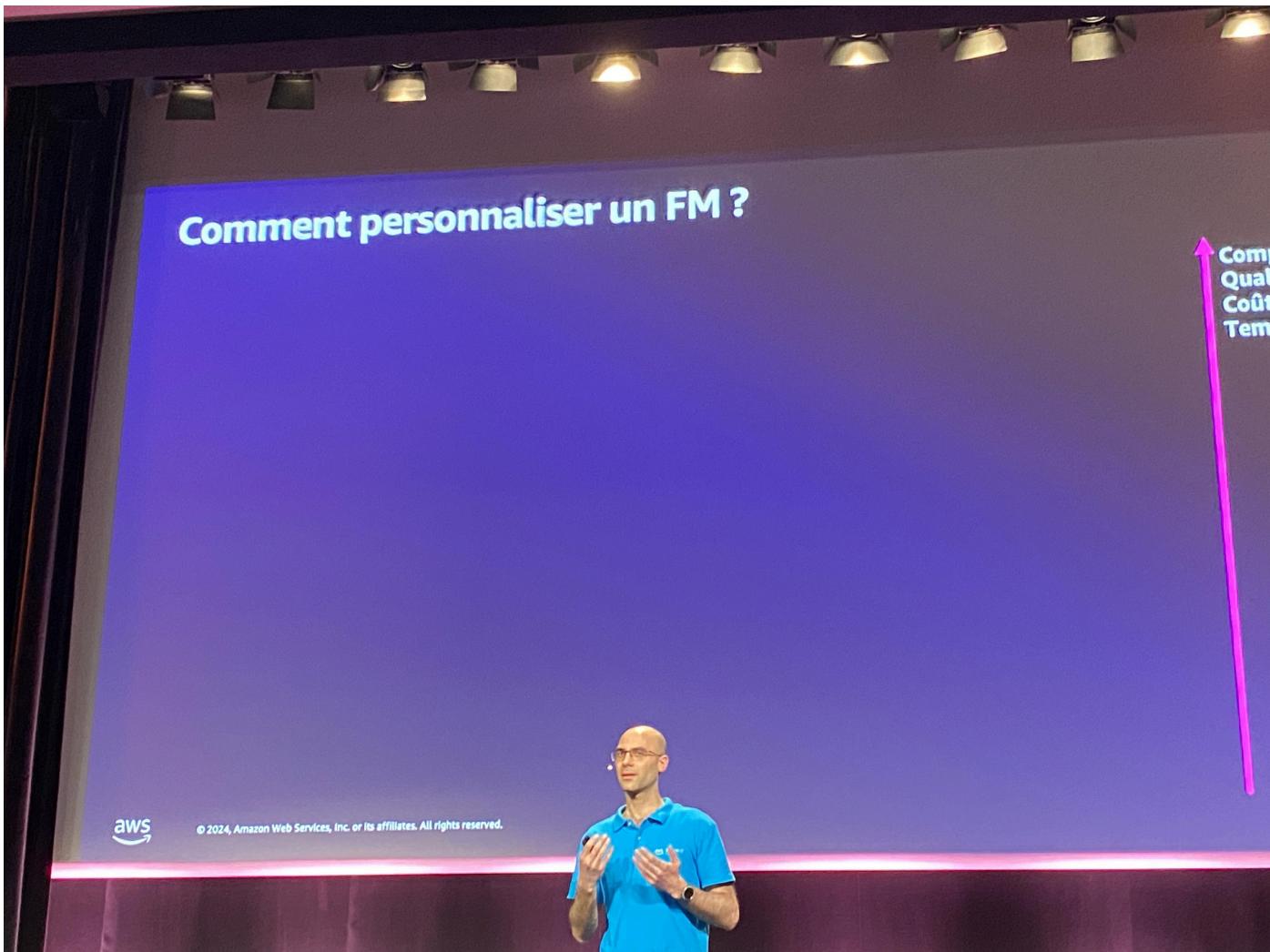
*"Là où il y a de la fumée de données, il y a un feu d'opportunités pour les entreprises"*

T. Redman

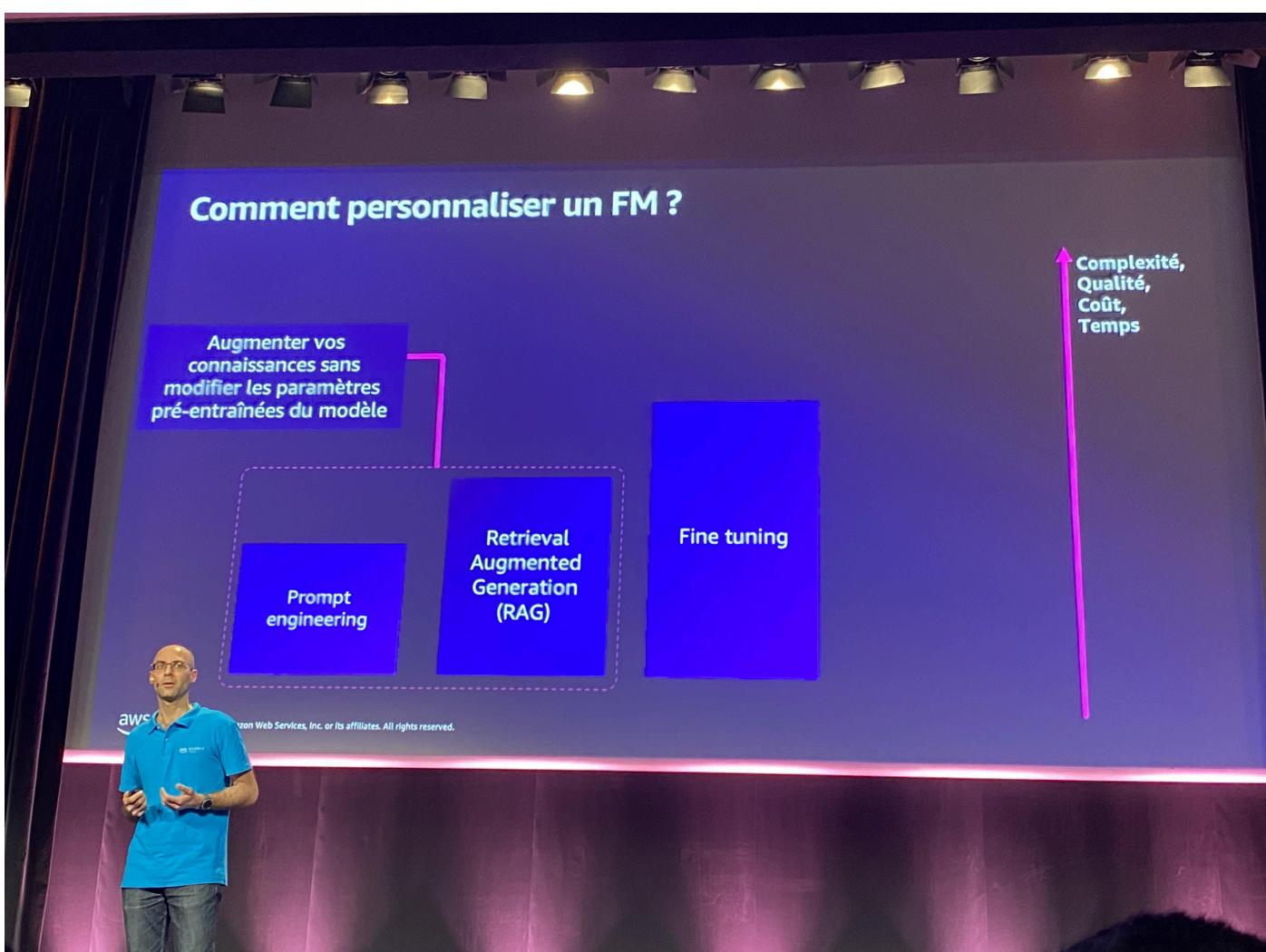
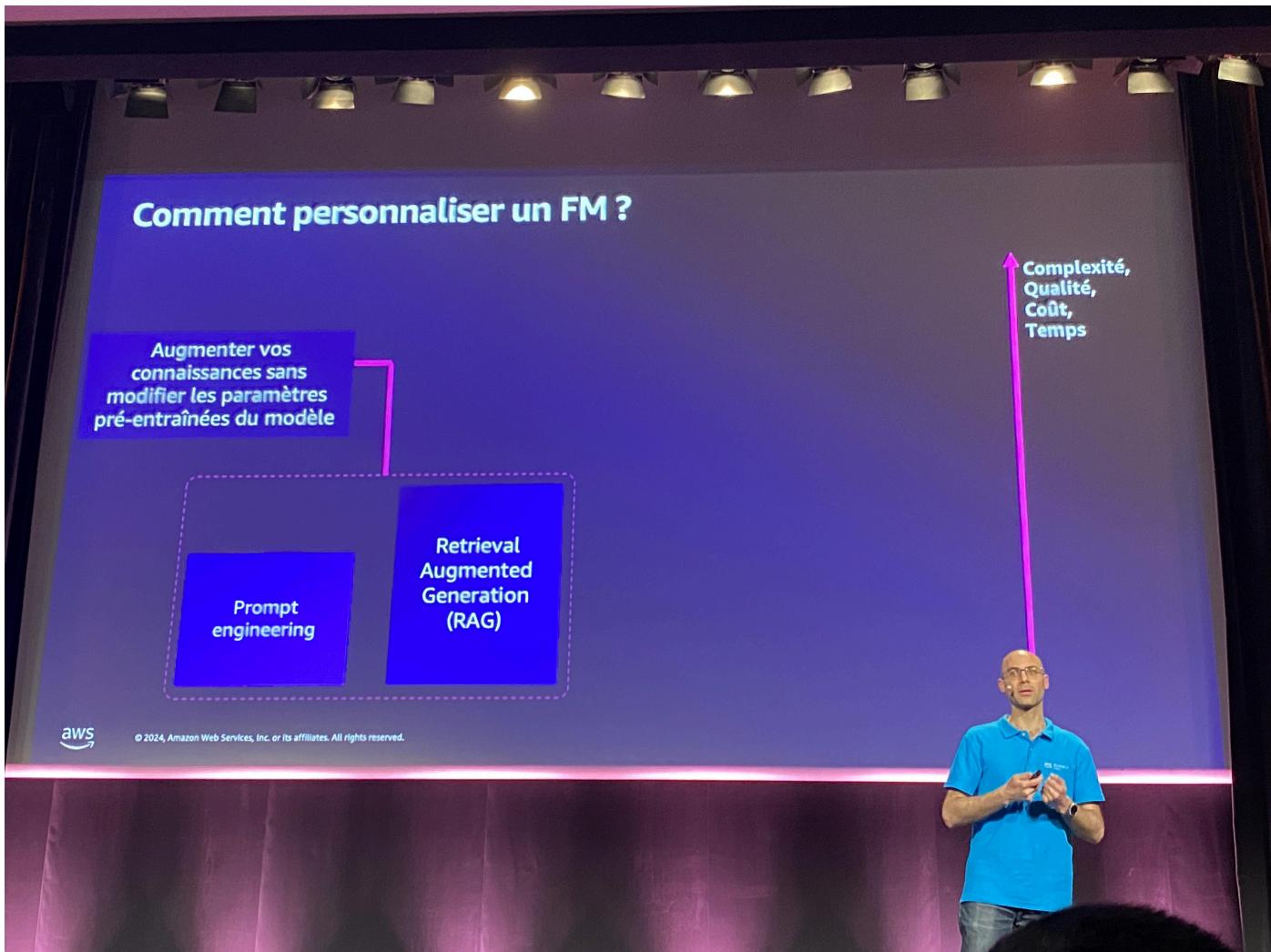


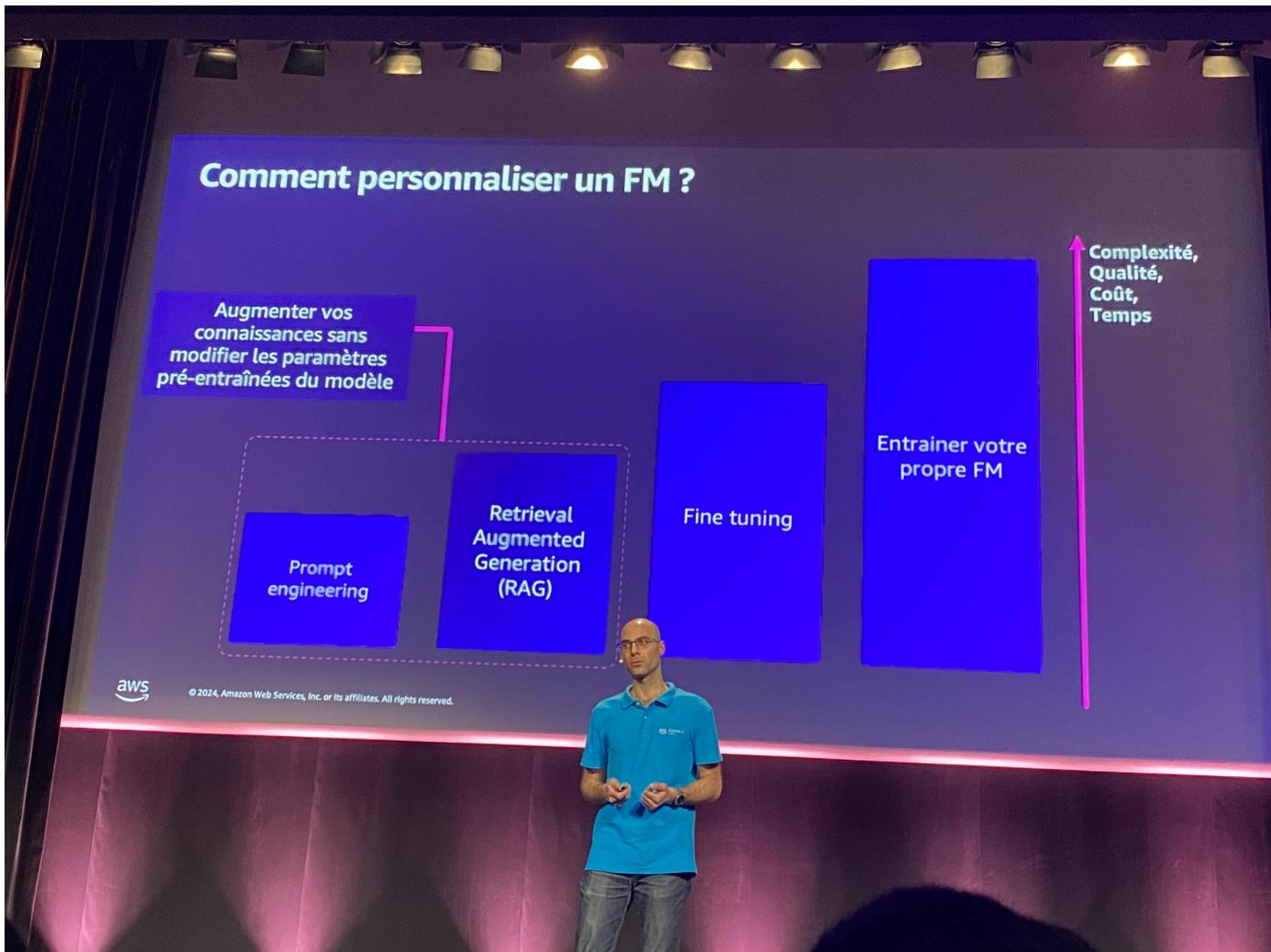
© 2024, Amazon Web Services, Inc. or its affiliates. All rights reserved.



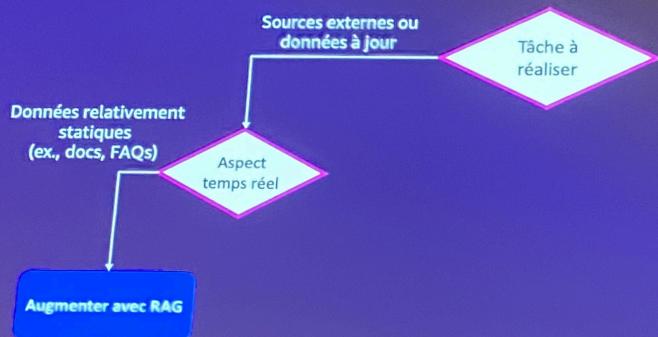








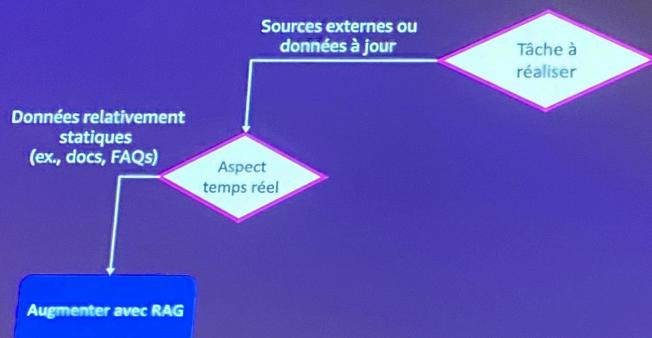
## Comment faire mon choix ?



© 2024, Amazon Web Services, Inc. or its affiliates. All rights reserved.

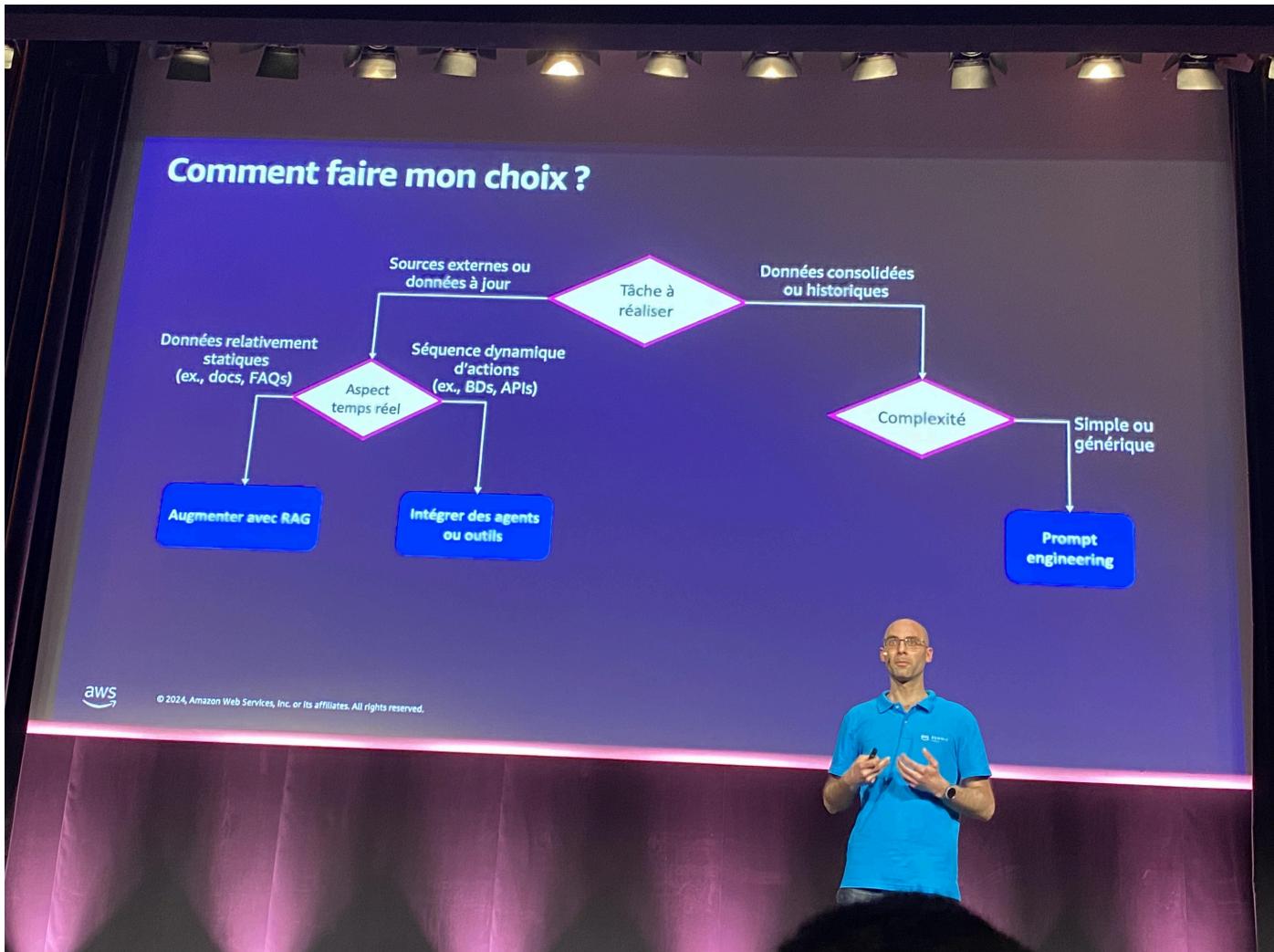
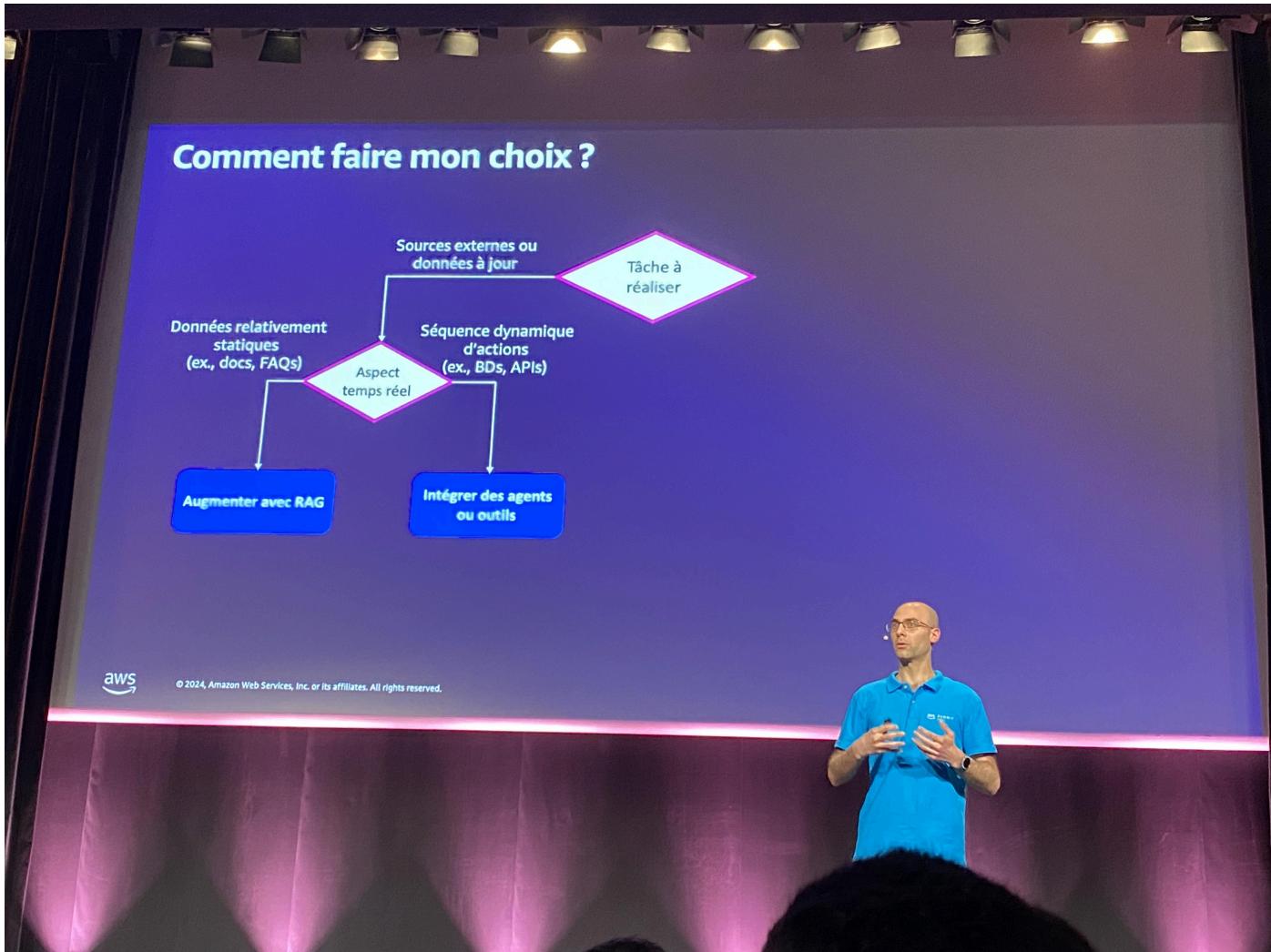


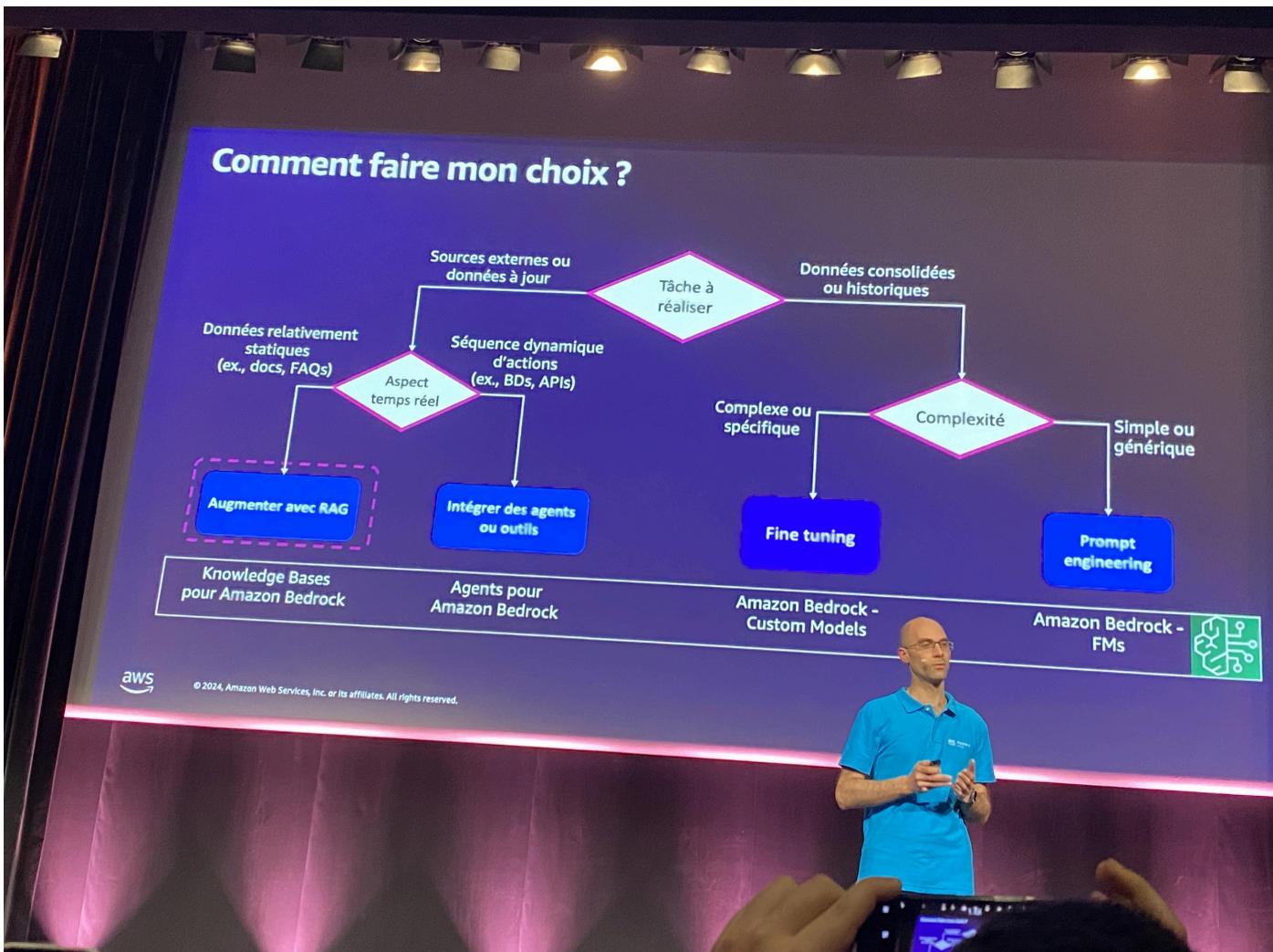
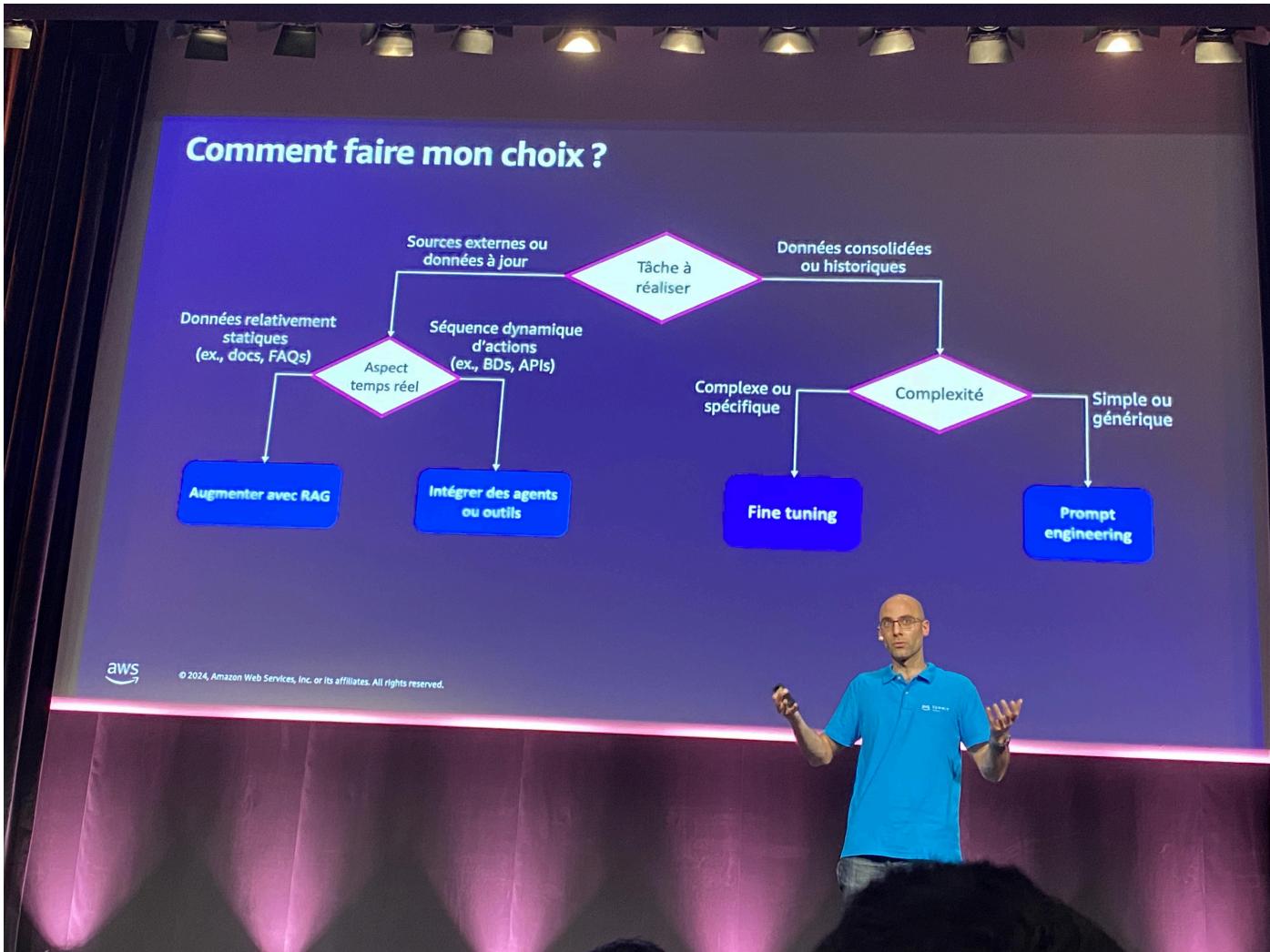
## Comment faire mon choix ?



© 2024, Amazon Web Services, Inc. or its affiliates. All rights reserved.







# Retrieval Augmented Generation ... C'est quoi ?



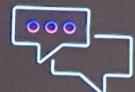
## Retrieval

Récupère le contenu pertinent depuis les sources de données externes, en fonction de la requête de l'utilisateur.



## Augmentation

Ajoute le contexte pertinent récupéré à l'instruction de l'utilisateur, qui sert d'entrée au modèle de fondation.



## Generation

Réponse du modèle de fondation basée sur l'instruction augmentée des données externes.



© 2024, Amazon Web Services, Inc. or its affiliates. All rights reserved.



# Cas d'utilisation du RAG



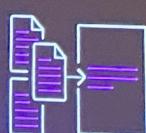
## Qualité des connaissances



## Réponses contextualisées



## Recherche personnalisée



## Synthèse en temps réel



© 2024, Amazon Web Services, Inc. or its affiliates. All rights reserved.



## Types de recherche (retrieval)



### Basé sur des règles

Données non structurées, supports documentaires

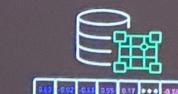
ex., recherche par mots clés



### Données structurées

Bases de données transactionnelles, APIs externes

ex., Select customers from All\_orders where order == 'XYZ'



### Recherche sémantique

Analyse sémantique à travers des modèles de text embeddings

Paris

Metro  
Tour Eiffel  
Champs Elysées



© 2024, Amazon Web Services, Inc. or its affiliates. All rights reserved.



## Types de recherche (retrieval)



### Basé sur des règles

Données non structurées, supports documentaires

ex., recherche par mots clés



### Données structurées

Bases de données transactionnelles, APIs externes

ex., Select customers from All\_orders where order == 'XYZ'



### Recherche sémantique

Analyse sémantique à travers des modèles de text embeddings

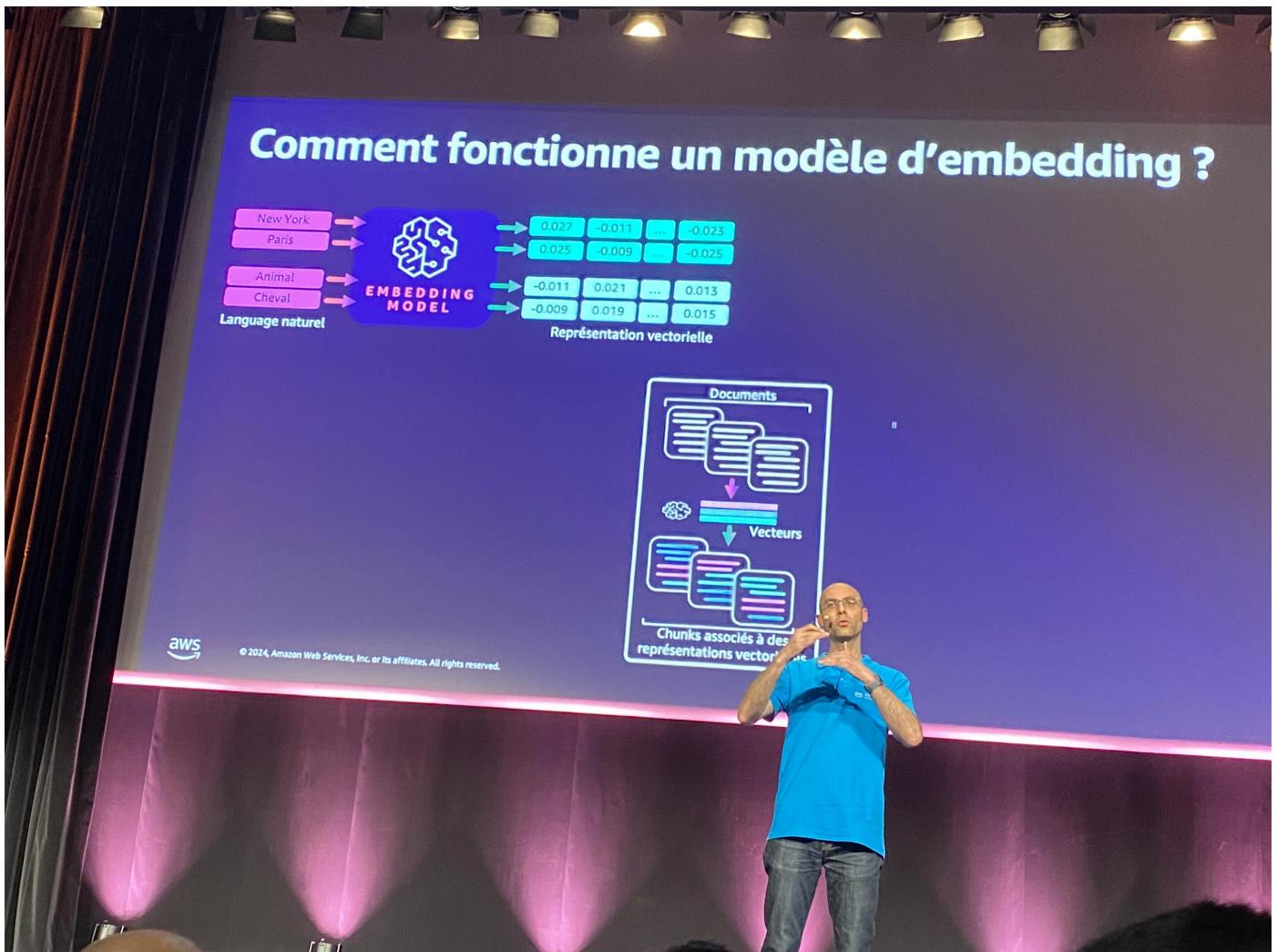
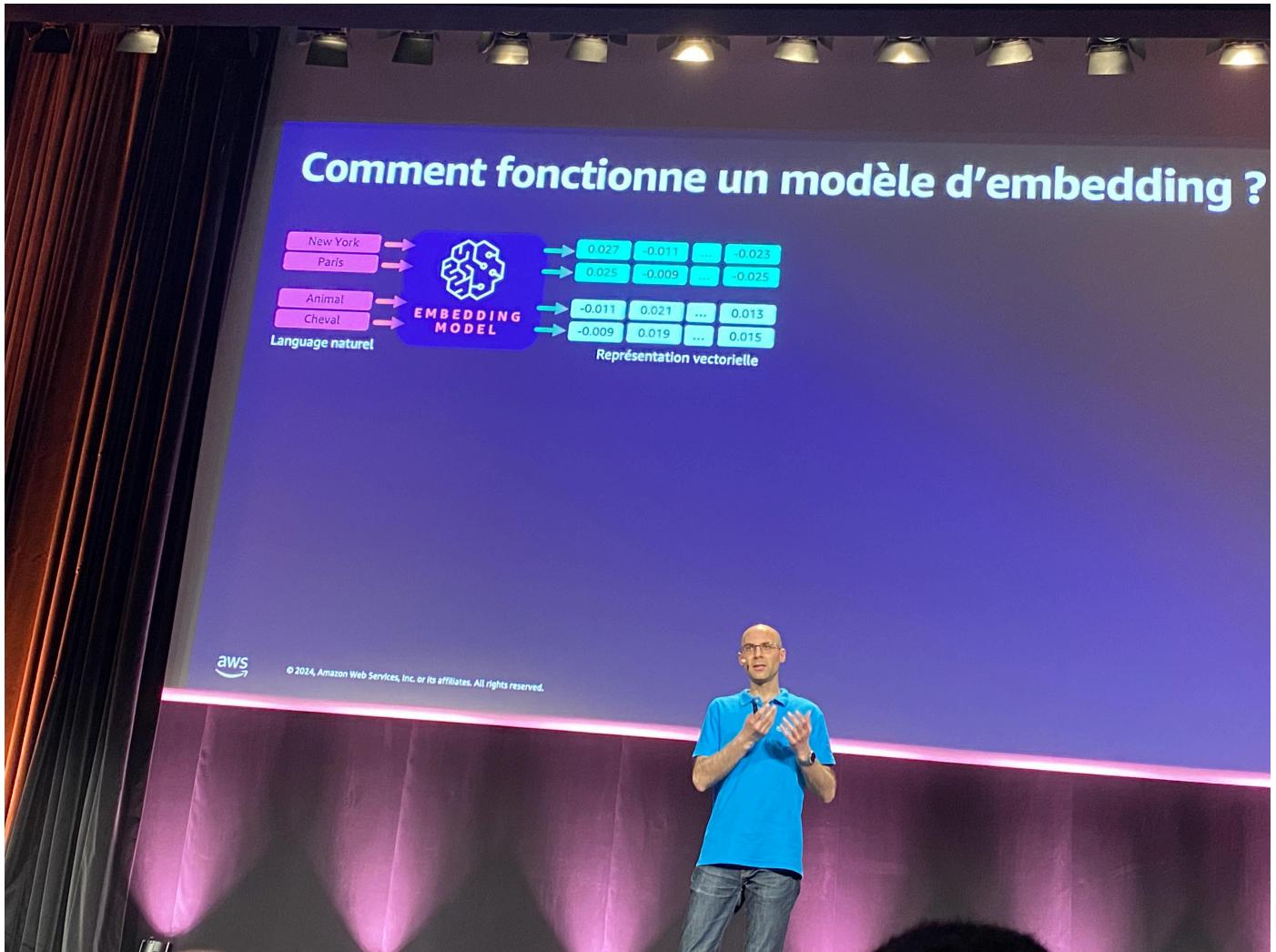
Paris

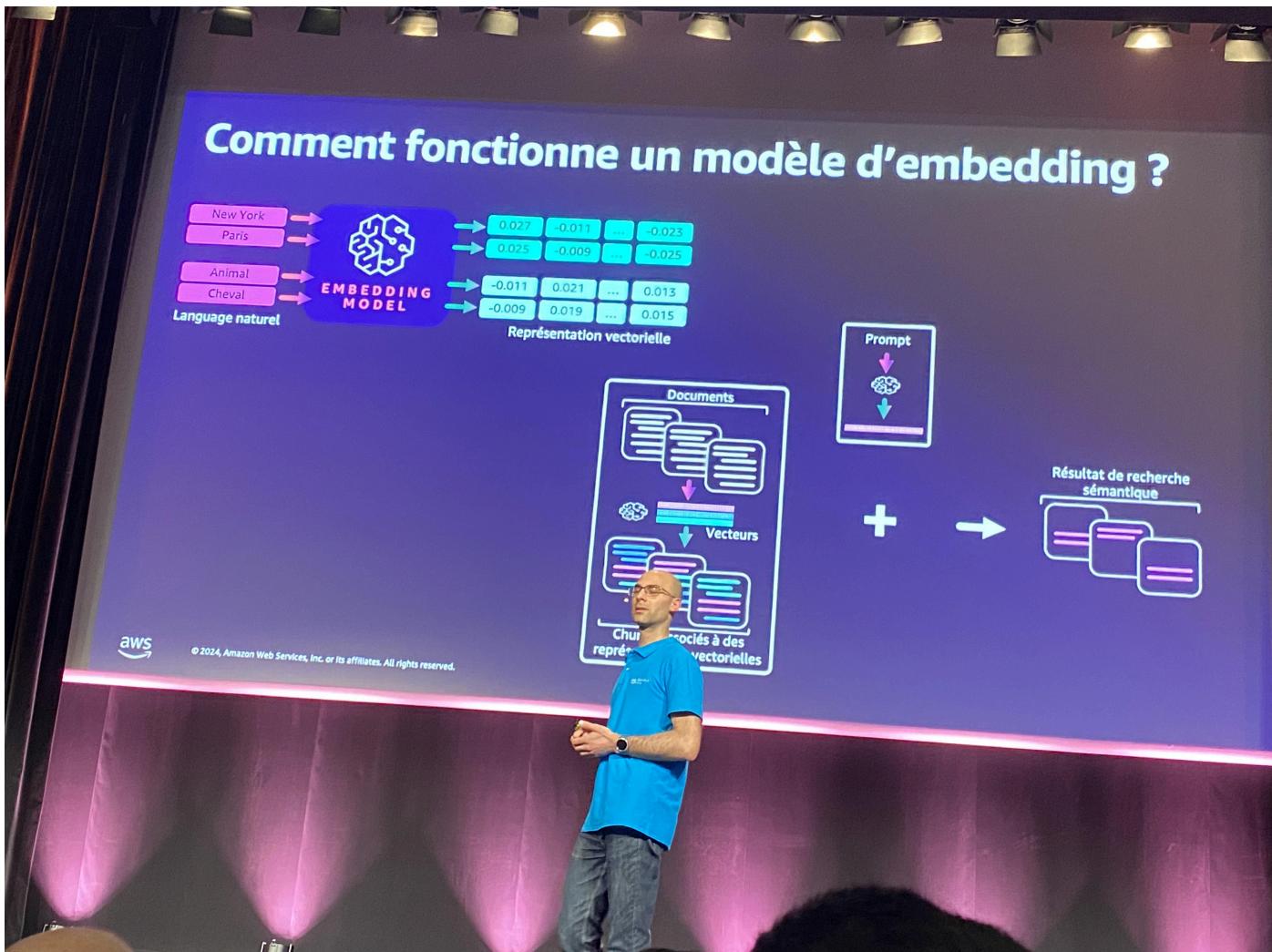
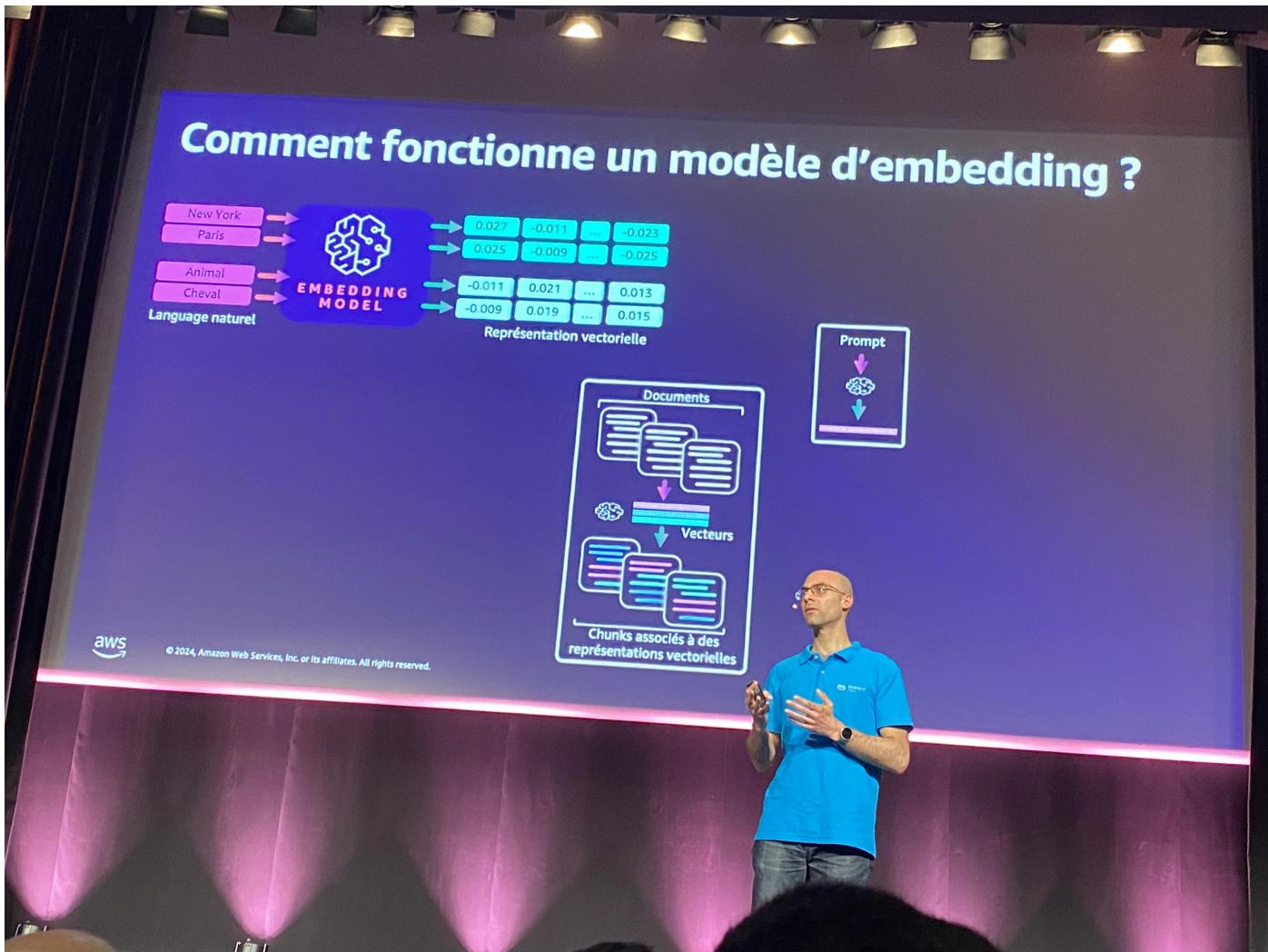
Metro  
Tour Eiffel  
Champs Elysées



© 2024, Amazon Web Services, Inc. or its affiliates. All rights reserved.







# Titan text embeddings model

 **Amazon Titan Text Embeddings**  
V2.0

Traduit les entrées textuelles (mots, phrases) en représentations numériques (embeddings). Comparer les embeddings produit des réponses contextuelles plus pertinentes que la correspondance de mots.

Max Tokens: 8,000  
Output Vectors: 1,536  
Language: Multilingual (25 languages)  
Model ID: amazon.titan-embed-g1-text-02

AWS © 2024, Amazon Web Services, Inc. or its affiliates. All rights reserved.

 **Points forts**

- Titan Text Embeddings propose des embeddings rapides, efficientes, de haute performance, et précis dans 25 langues.
- Optimisé pour les tâches de récupération de texte, de similarité sémantique et de clustering.
- Les applications de ce modèle incluent la recherche sémantique et la personnalisation

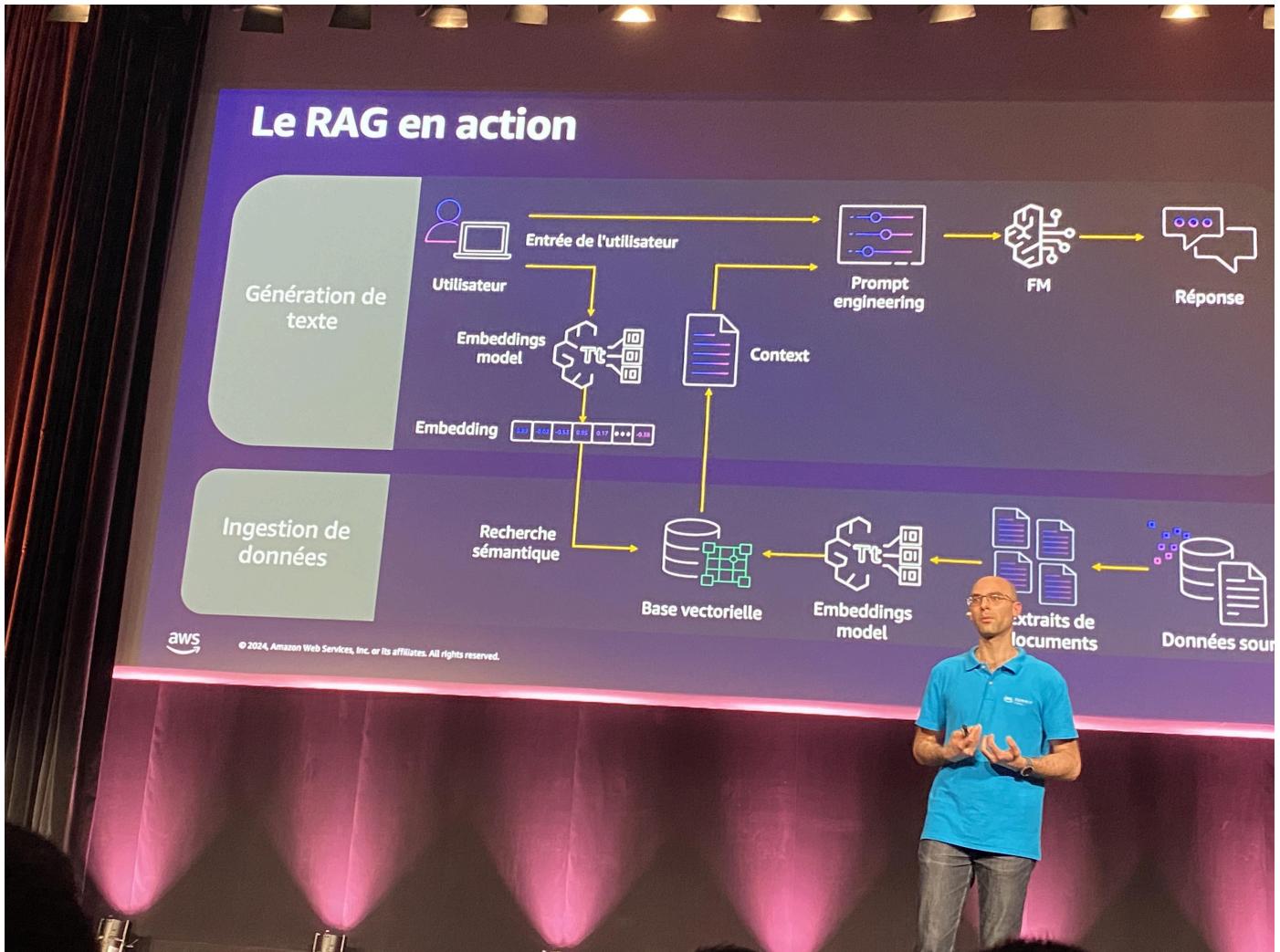
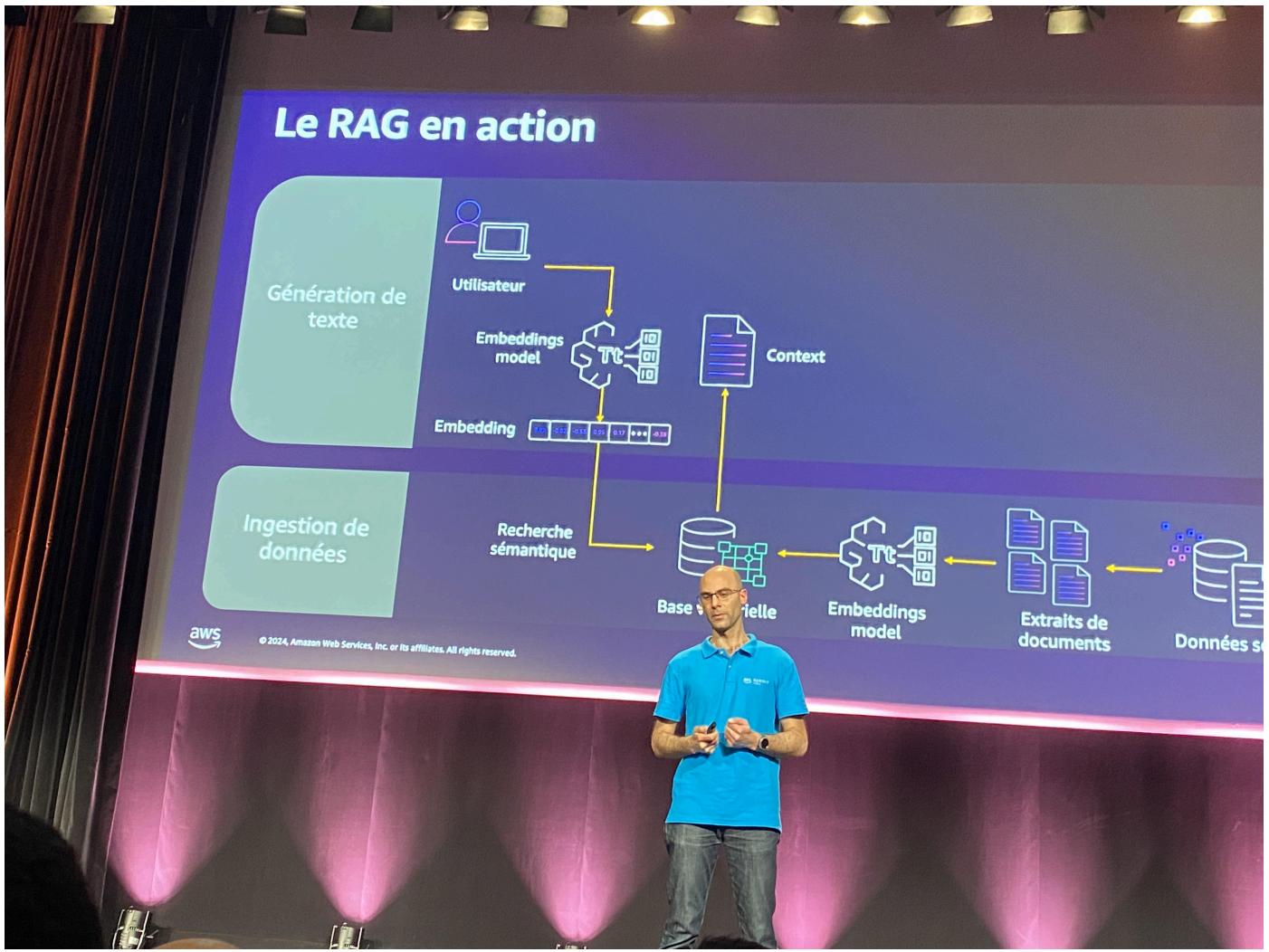


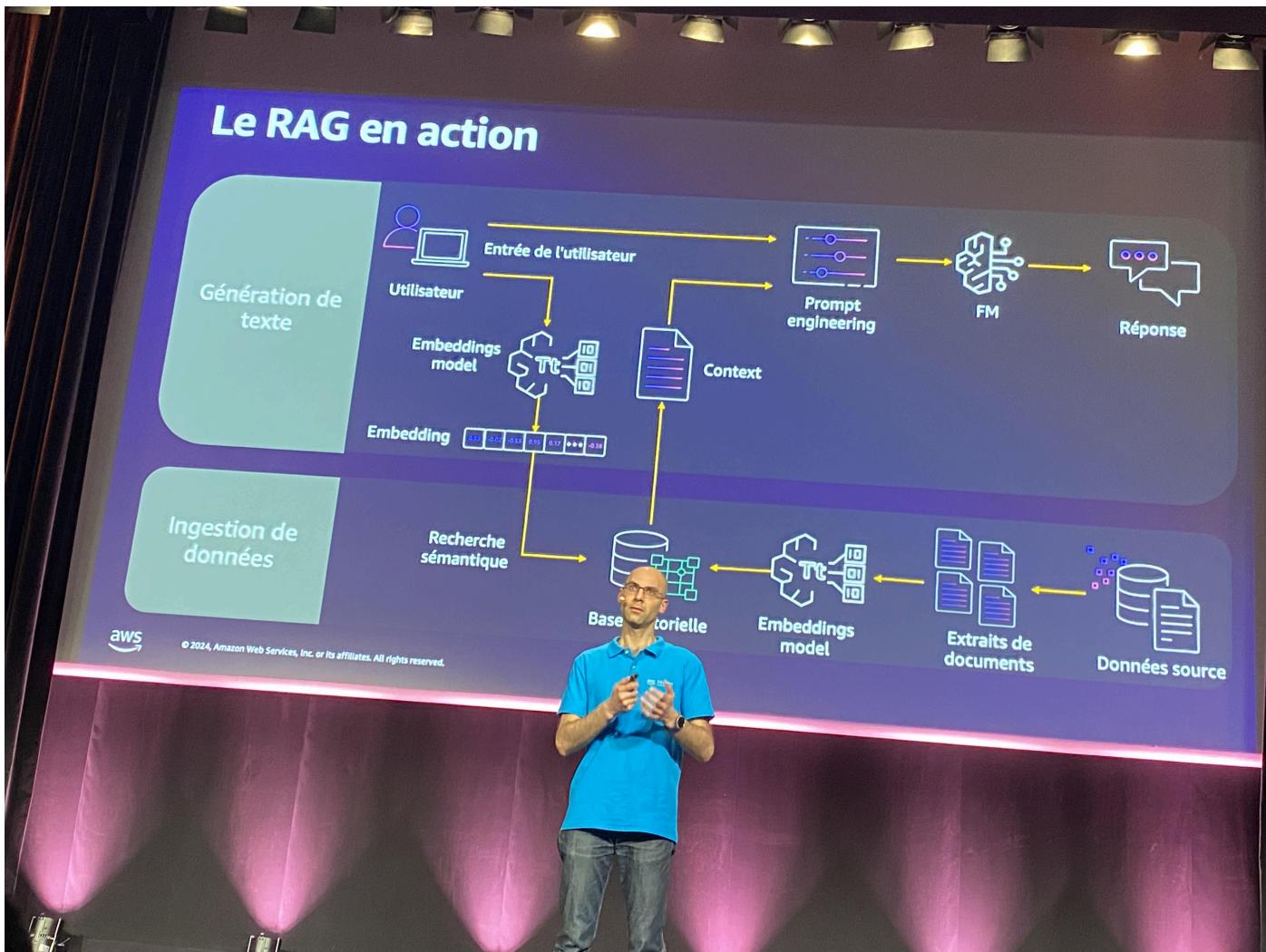
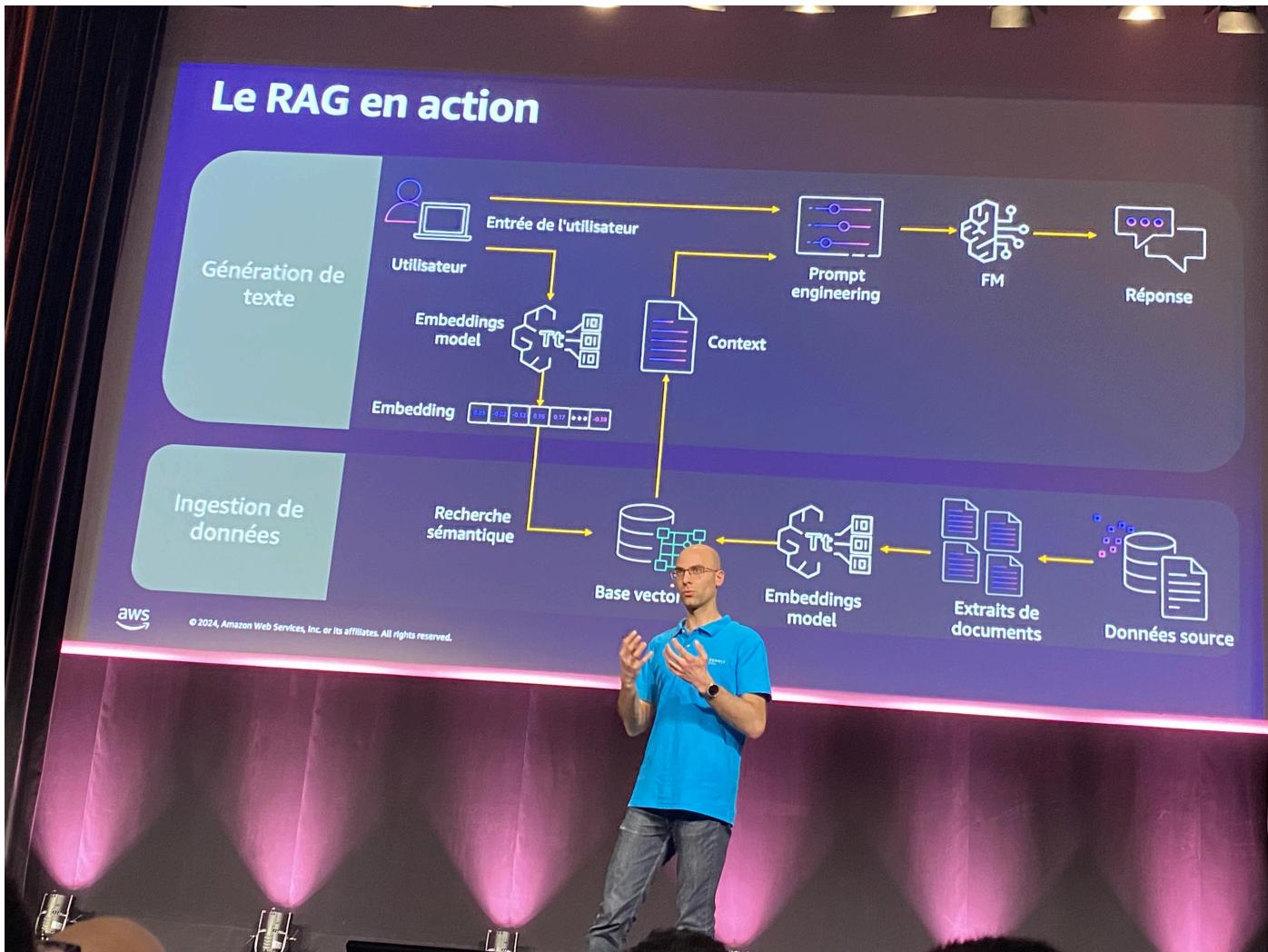
## Le RAG en action

AWS © 2024, Amazon Web Services, Inc. or its affiliates. All rights reserved.

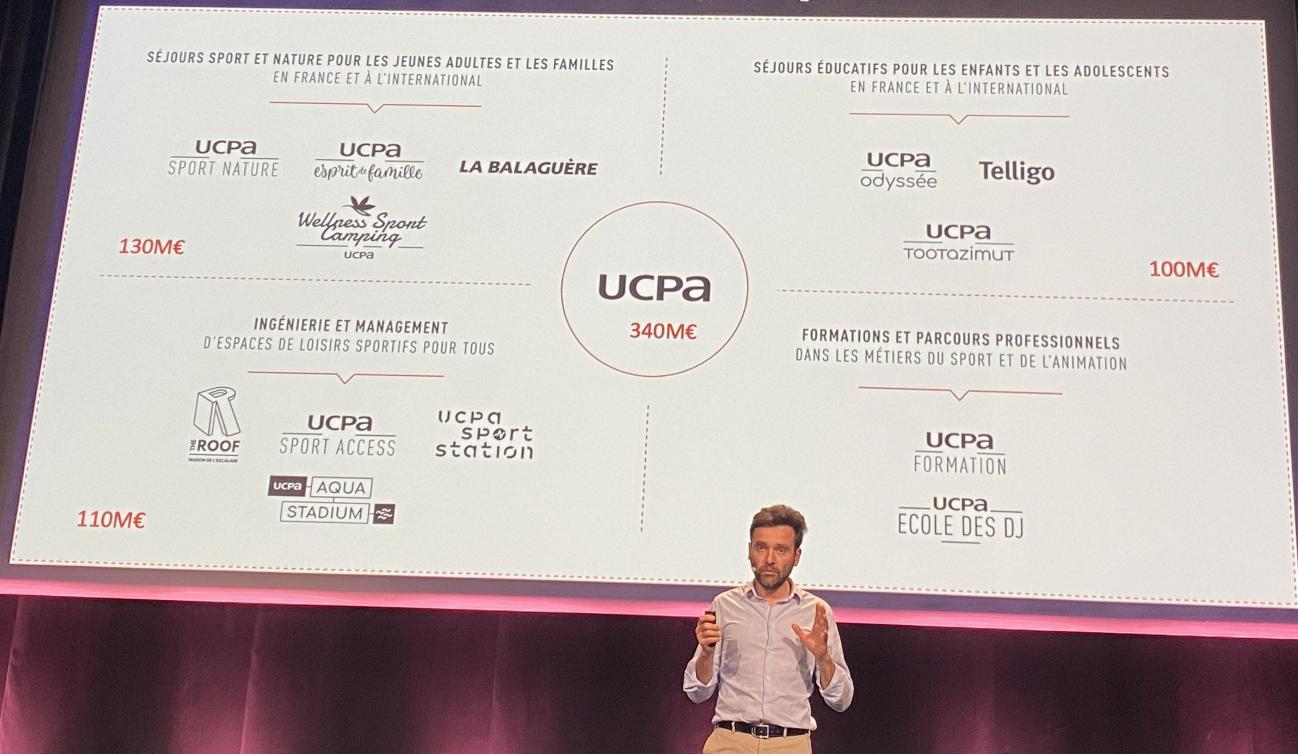








# L'UCPA : l'hyperscaler du sport pour tous



## **Ecosystème technique UCPA depuis 2018**

All in AWS

# Serverless

# Automation

**100% interne**

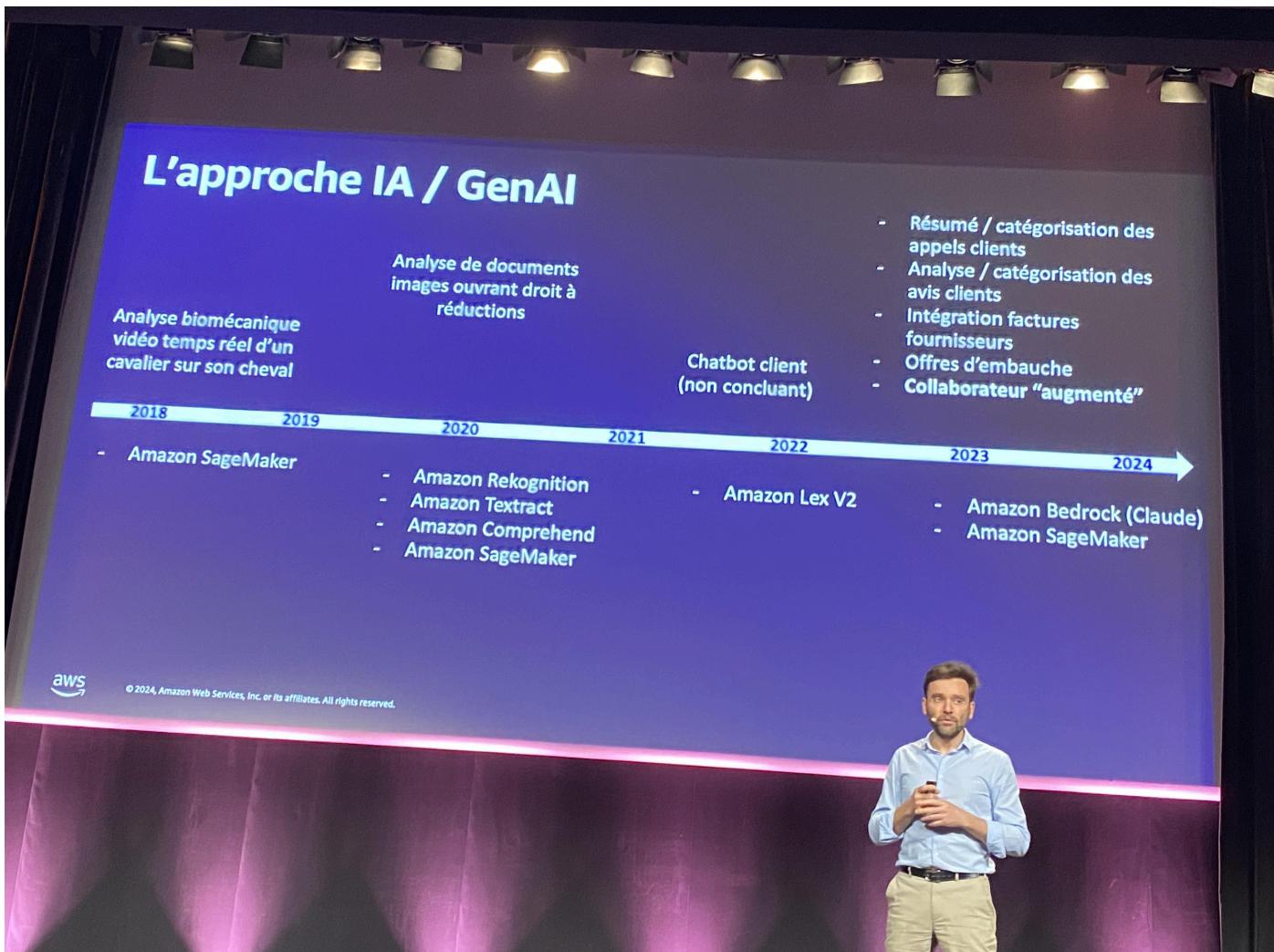
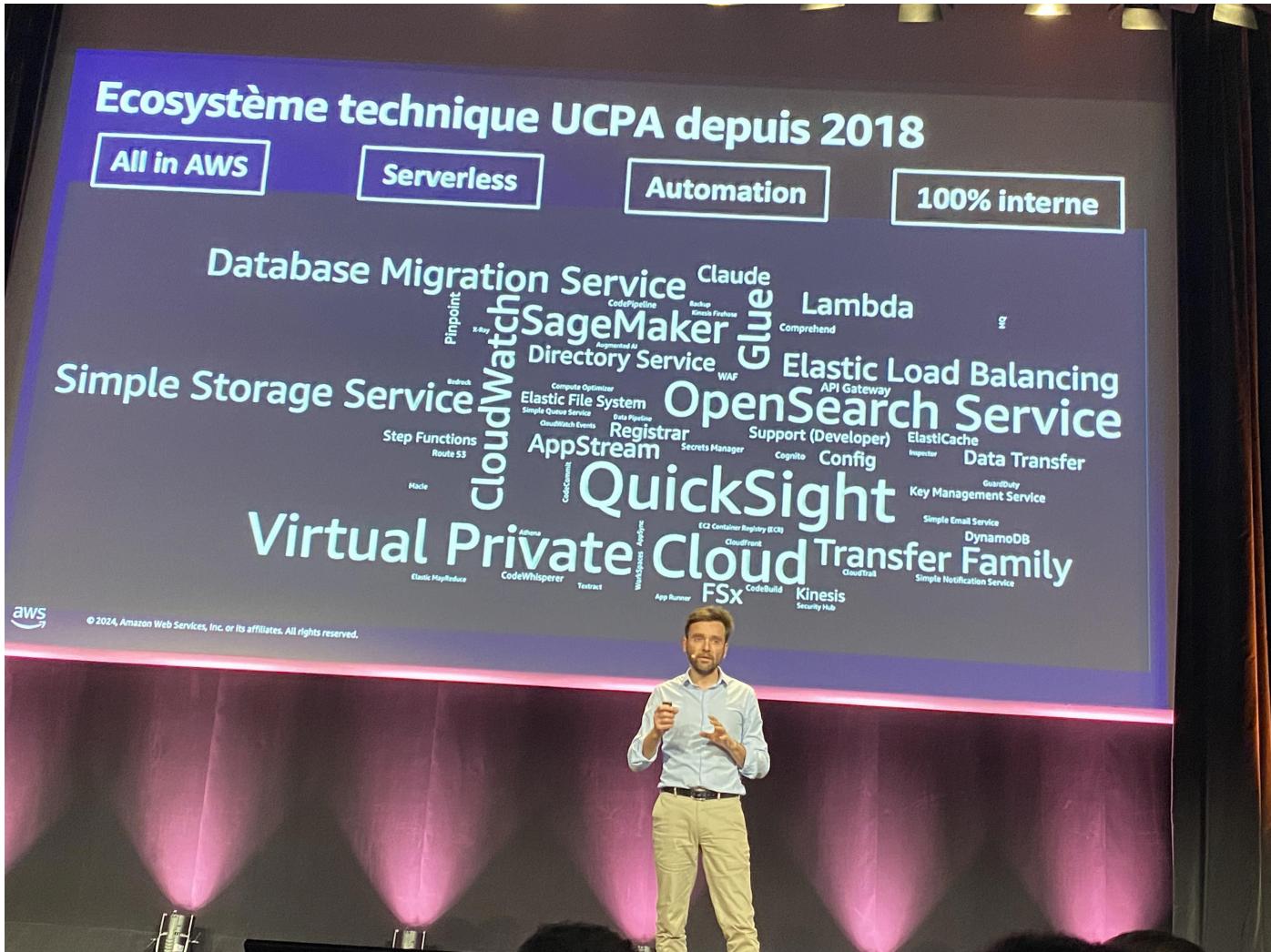
# Simple Storage Service

The word cloud illustrates several AWS services and their features:

- CloudWatch**: Log Insights, Metrics, CloudWatch Events
- SageMaker**: CodePipeline, Keras Training
- Amazon Kinesis**
- Amazon Simple Queue Service**
- Amazon AppStream**: CodeCommit
- Amazon Lambda**: Comprehend
- Amazon Directory Service**: Augmented AI
- Amazon OpenSearch Service**: WAF, API Gateway
- Amazon Elastic File System**: Compute Optimizer
- Amazon QuickSight**: Data Pipeline, Registrar, Secrets Manager, Cognito, Config, Support (Developer)
- Amazon Lambda**: Claude
- Amazon Glue**: Data Pipeline
- Amazon Lambda**: Lambda
- Amazon Lambda**: Lambda

# Virtual Private Cloud Transfer Family

© 2024, Amazon Web Services, Inc. or its affiliates. All rights reserved.



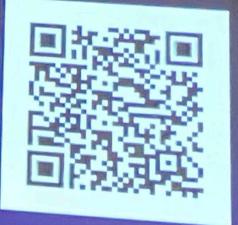
## Le collaborateur "augmenté"

**Le besoin : augmenter la créativité et la productivité des collaborateurs**

- ▣ Mettre à disposition de 200+ collaborateurs un Chatbot privé, en toute sécurité
- ▣ Permettre au collaborateur de dialoguer avec ses propres données
- ▣ Permettre au collaborateur de dialoguer avec les données de l'entreprise

**La solution : intégrer une chaîne RAG sur Bedrock**

- ▣ D'abord en construisant tout nous même :
  - 6 mois de POC itératif
- ▣ Puis en intégrant un projet OpenSource : GENAI Chatbot
  - 1 mois de développement / intégration



aws  
© 2024, Amazon Web Services, Inc. or its affiliates. All rights reserved.



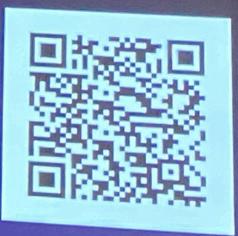
## Le collaborateur "augmenté"

**Le besoin : augmenter la créativité et la productivité des collaborateurs**

- ▣ Mettre à disposition de 200+ collaborateurs un Chatbot privé, en toute sécurité
- ▣ Permettre au collaborateur de dialoguer avec ses propres données
- ▣ Permettre au collaborateur de dialoguer avec les données de l'entreprise

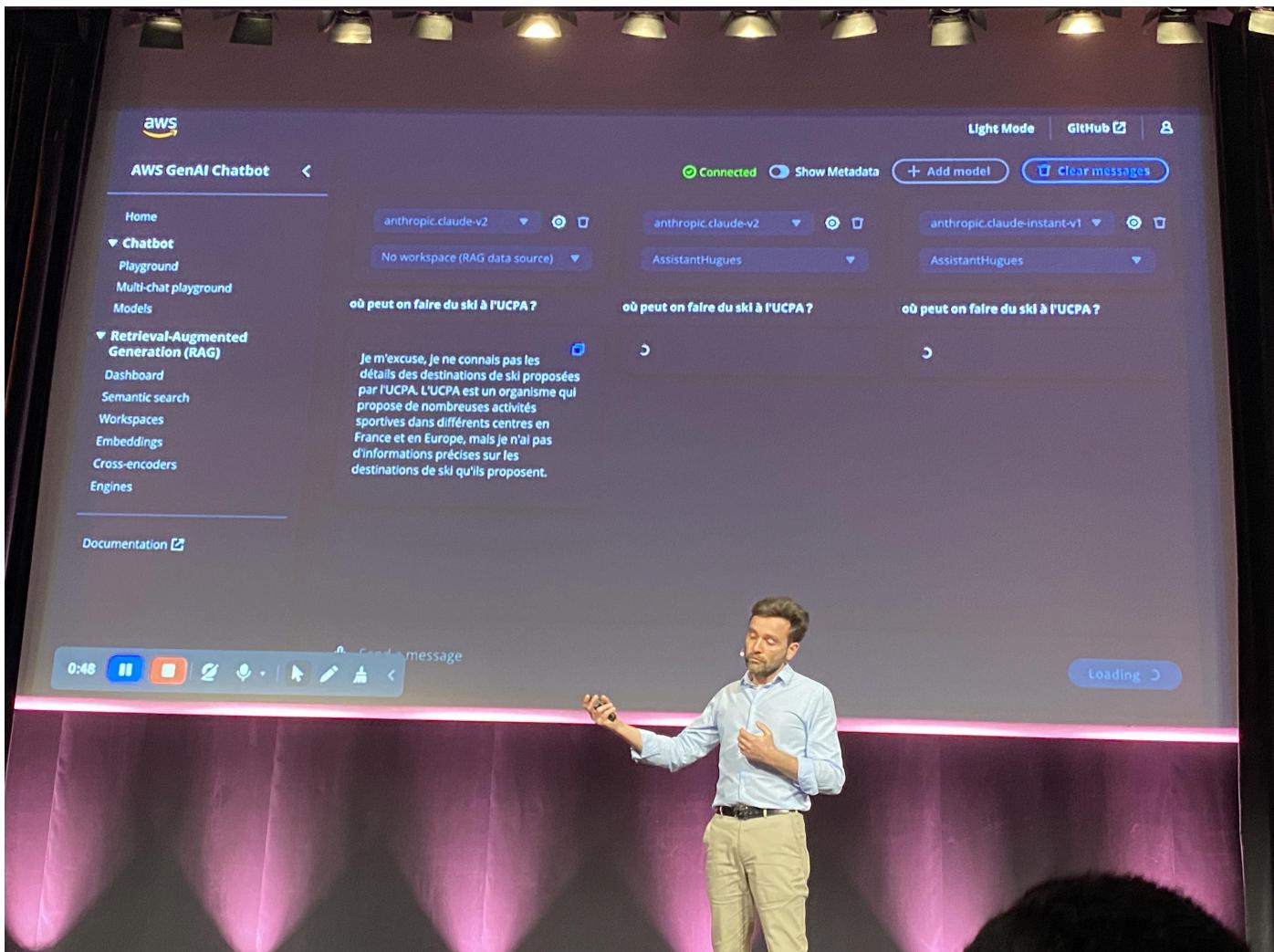
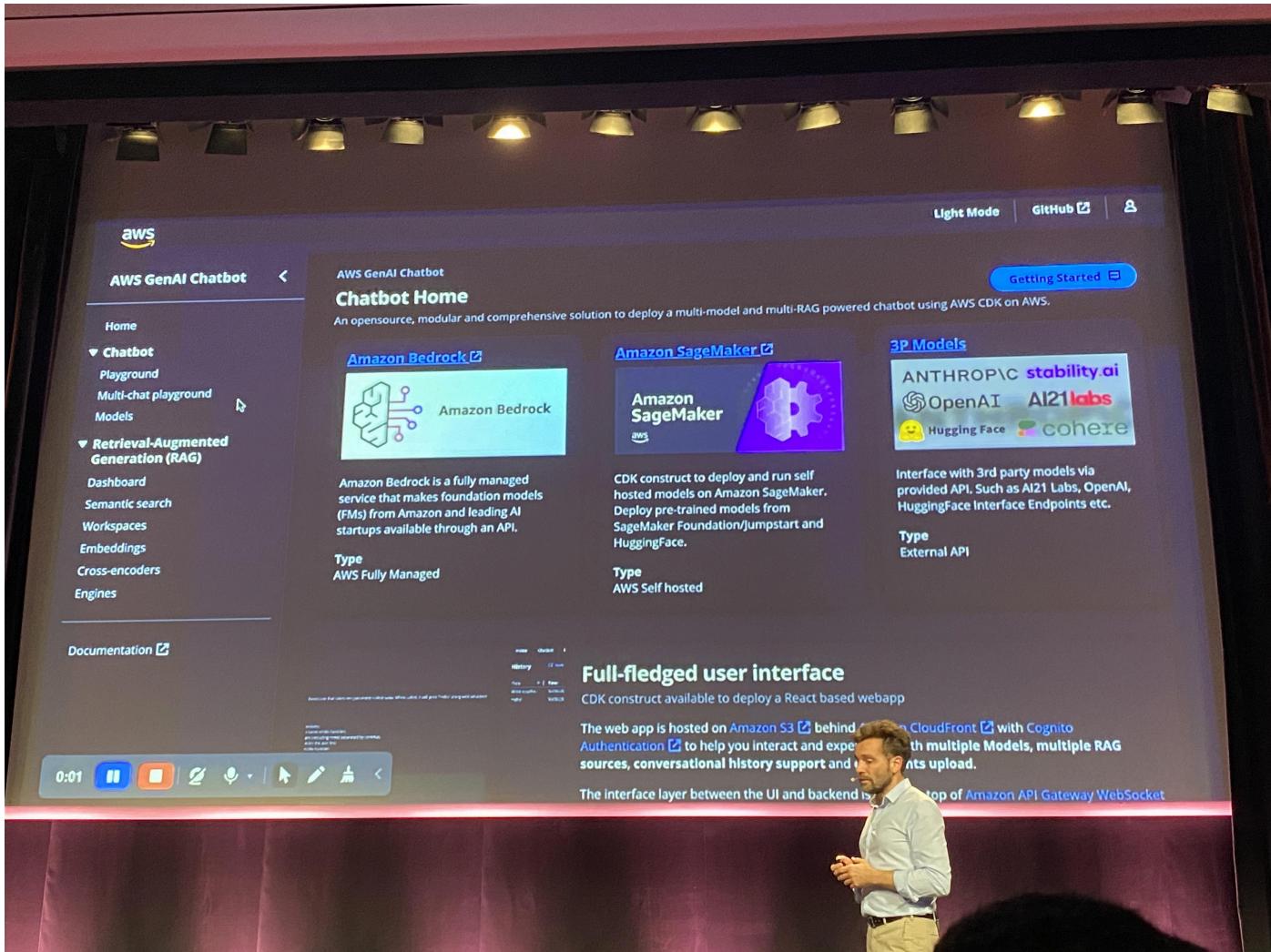
**La solution : intégrer une chaîne RAG sur Bedrock**

- ▣ D'abord en construisant tout nous même :
  - 6 mois de POC itératif
- ▣ Puis en intégrant un projet OpenSource : GENAI Chatbot
  - 1 mois de développement / intégration



aws  
© 2024, Amazon Web Services, Inc. or its affiliates. All rights reserved.





AWS GenAI Chatbot

Home Chatbot Playground Multi-chat playground Models

Retrieval-Augmented Generation (RAG)

Dashboard Semantic search Workspaces Embeddings Cross-encoders Engines

Documentation

Je m'excuse, je ne connais pas les détails des destinations de ski proposées par l'UCPA. L'UCPA est un organisme qui propose de nombreuses activités sportives dans différents centres en France et en Europe, mais je n'ai pas d'informations précises sur les destinations de ski qu'ils proposent.

D'après les informations fournies, il semble que l'on peut faire du ski dans les villages sportifs UCPA suivants :

- Pralognan, dans les Alpes, à 5 minutes à pied des remontées mécaniques.
- Les Deux Alpes, dans un chalet au cœur d'une station village préservée, avec 120 km de pistes adaptées à tous les niveaux.
- Paradiski, au carrefour entre Les Arcs, La Plagne et Peisey-Vallandry, avec de nombreux espaces pour le snowboard et le ski.
- Les Arcs, avec un village sportif "skis aux pieds".
- Tignes, avec vue panoramique sur un glacier et neige garantie, ainsi qu'un spa et sauna.

D'après les informations contextuelles fournies, on peut faire du ski à l'UCPA dans plusieurs endroits, notamment :

- À Pralognan-la-Vanoise, où le village sportif UCPA se trouve à 5 minutes à pied des remontées mécaniques et du village.
- Aux Deux Alpes, où l'UCPA accueille dans un chalet au cœur d'une station village préservée.
- À Valmorel, station située entre le Parc Naturel Régional de la Chartreuse et le Parc Naturel Régional du Vercors, avec 120 km de pistes adaptées à tous les niveaux.
- Aux Arcs, village sportif "skis aux pieds" au carrefour d'espaces naturels remarquables.

Donc en résumé, l'UCPA propose des activités de ski dans plusieurs stations des Alpes et des Pyrénées, avec hébergement et encadrement sur place.

Send >

AWS GenAI Chatbot

Home Chatbot Playground Multi-chat playground Models

Retrieval-Augmented Generation (RAG)

Dashboard Semantic search Workspaces Embeddings Cross-encoders Engines

Documentation

Workspace Engine

Vector Engine

- Aurora Serverless v2 PostgreSQL (pgvector)  
Pgvector is open-source vector similarity search for PostgreSQL.
- Amazon OpenSearch Serverless  
The vector engine for Amazon OpenSearch Serverless introduces a simple, scalable, and high-performing vector storage and search capability.
- Amazon Kendra  
Uses Kendra Retrieve API as a retriever for retrieval augmented generation (RAG) systems.

Aurora Workspace Configuration

Workspace Name

My Workspace

Embeddings Model

amazon.titan-embed-text-v1 (1536)

Bedrock

amazon.titan-embed-text-v1 (1536)

OpenAI

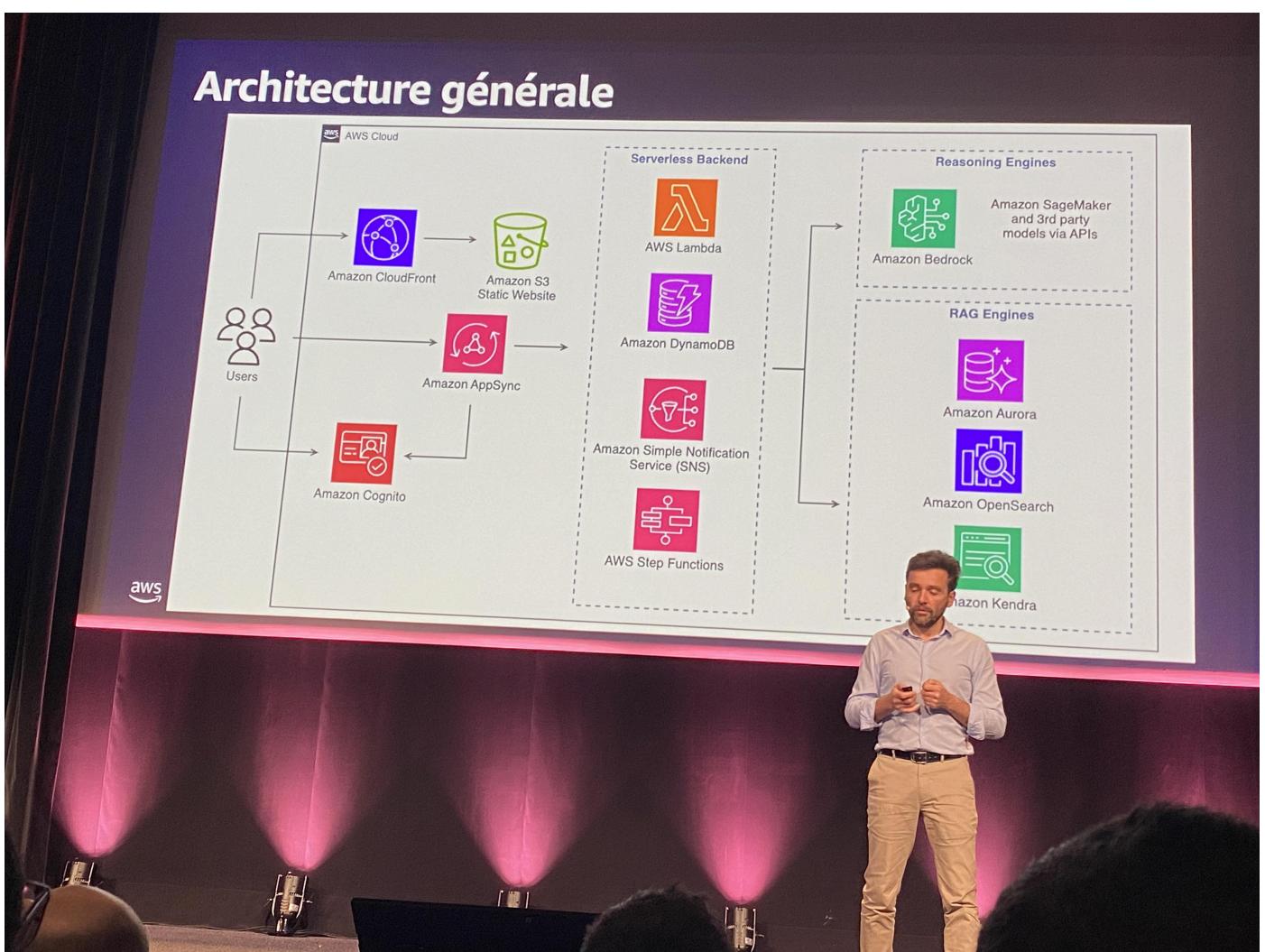
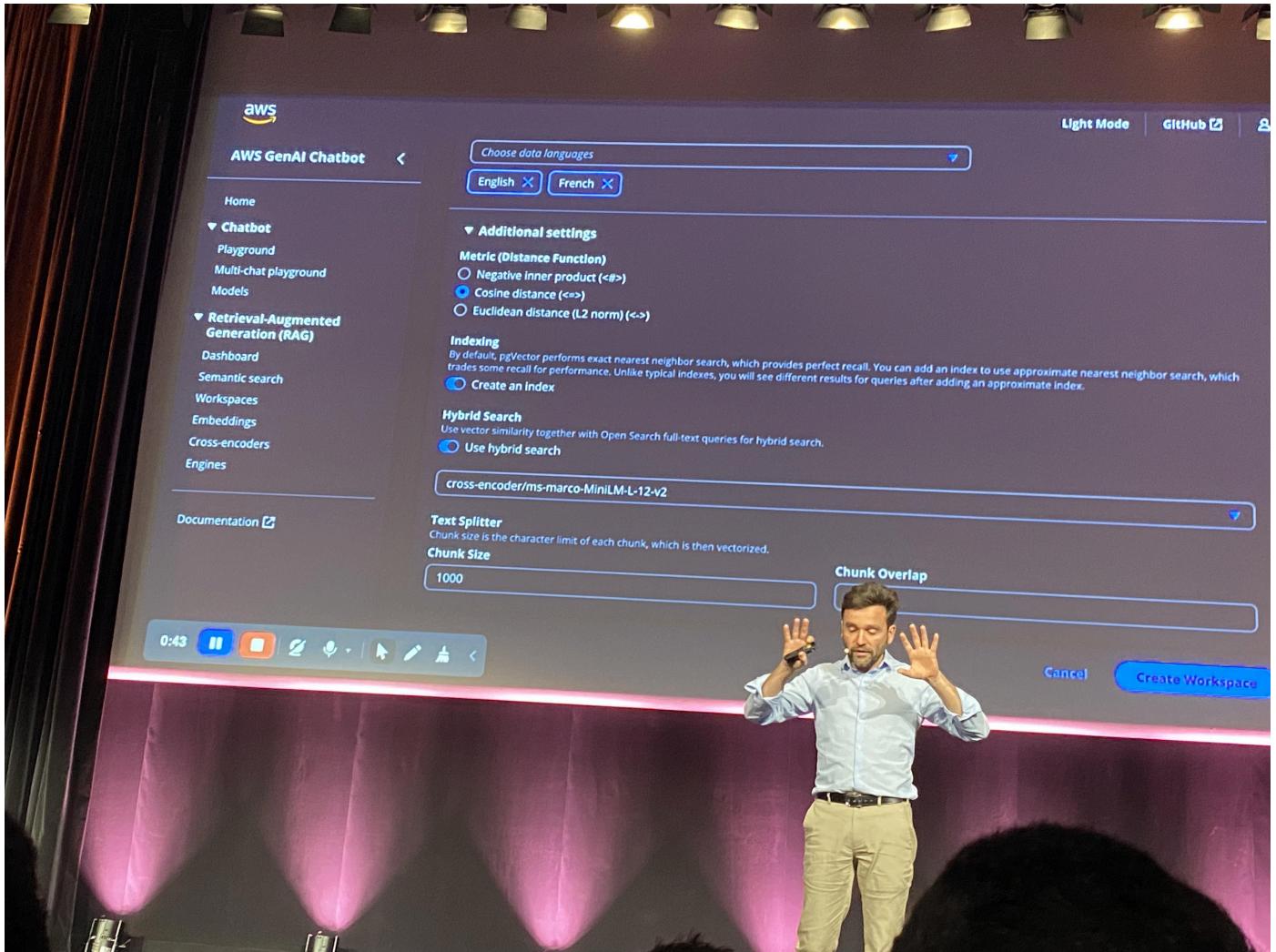
text-embedding-ada-002 (1536)

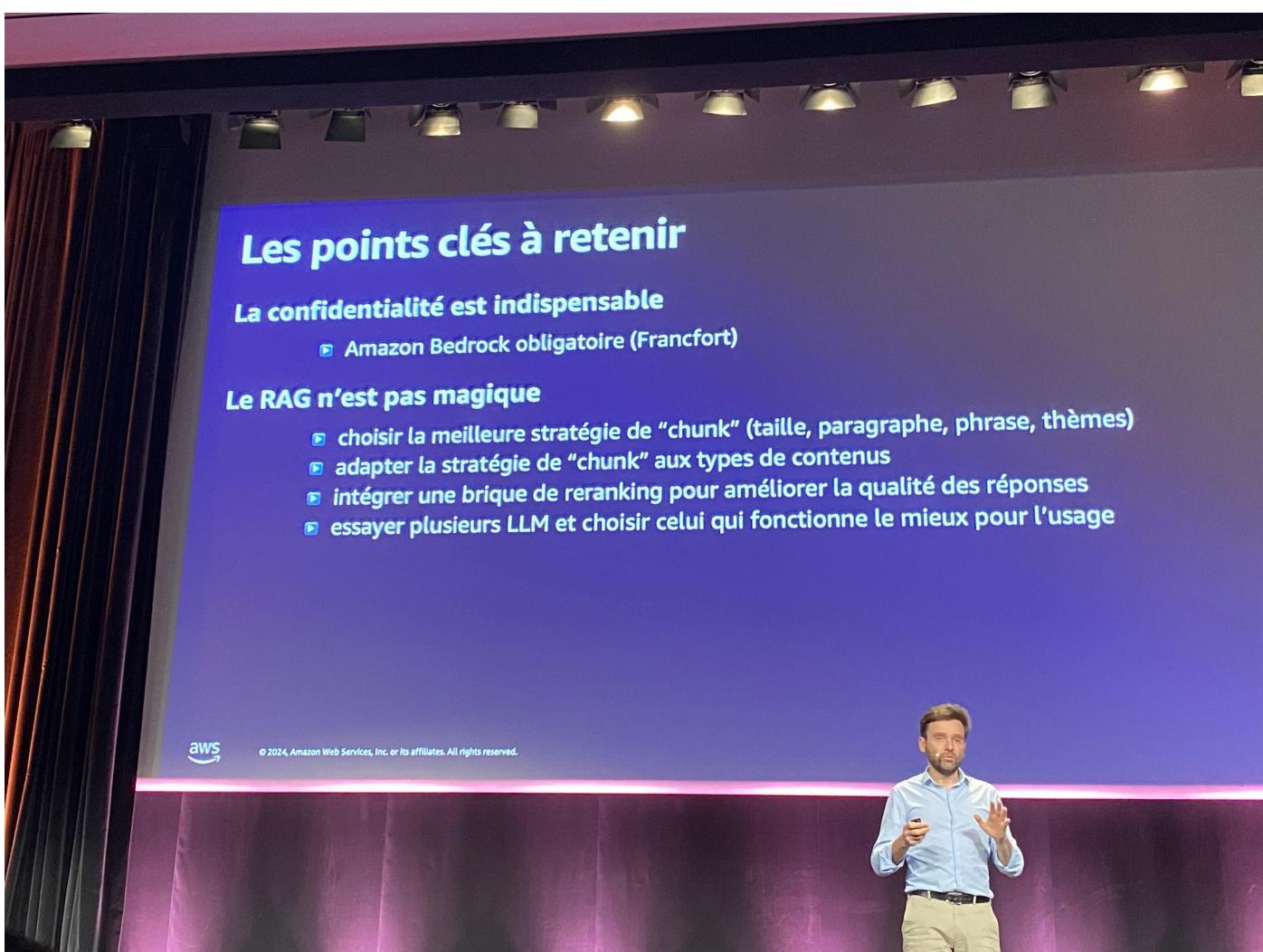
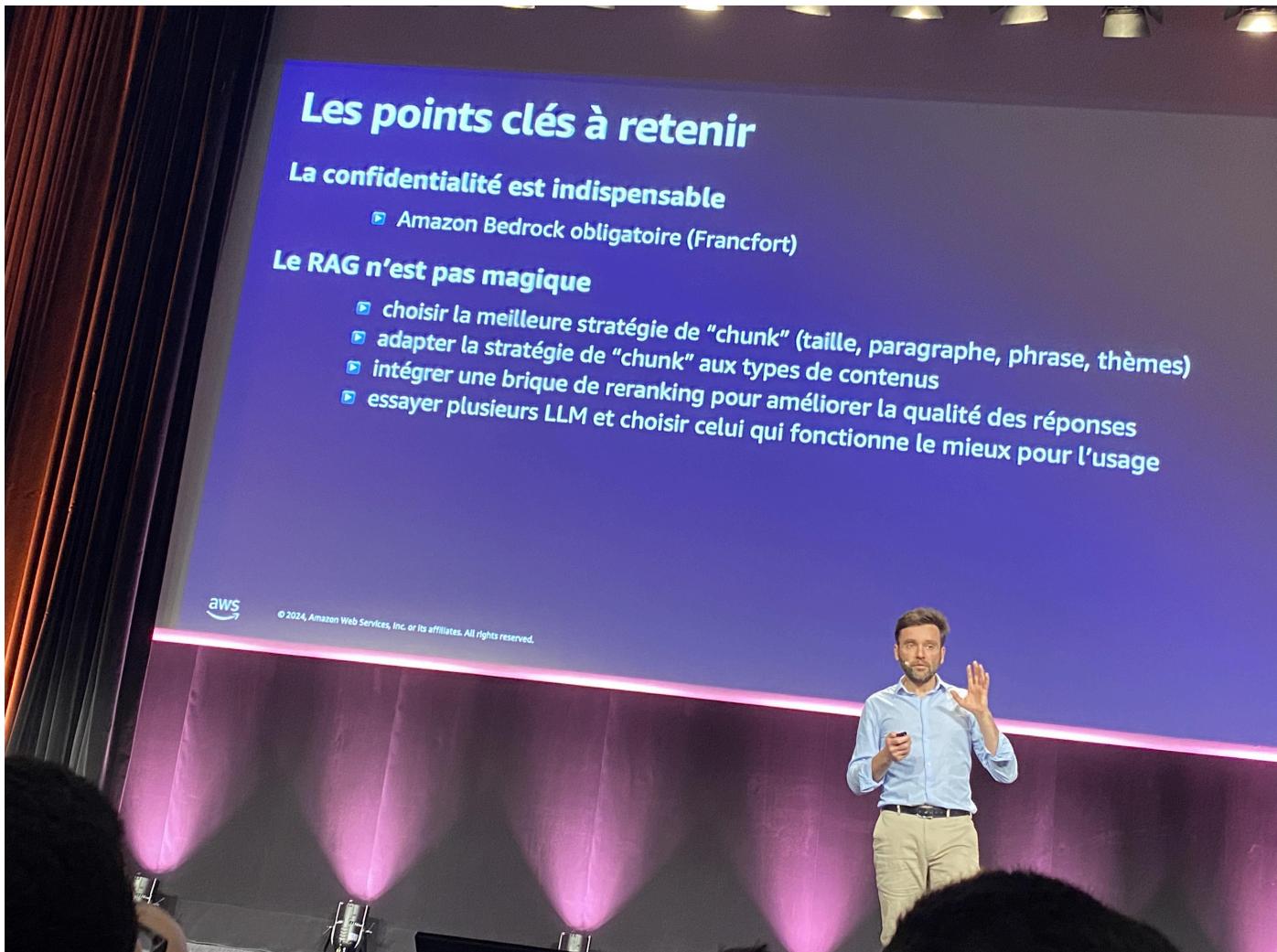
SageMaker

inffloat/multilingual-e5-large (1024)

sentence-transformers/all-MiniLM-L6-v2 (384)

Create Workspace





## Les points clés à retenir

### Gérer le cycle de vie de la "connaissance"

- ☒ éviter les mauvaises données

### Hallucinations

- ☒ à gérer impérativement

### Peu de recul à ce jour sur les coûts au token

- ☒ Attention à la facture !

### Peu de LLM "parlent" français

- ☒ Claude privilégié, en attente de Mistral (sur Francfort)



© 2024, Amazon Web Services, Inc. or its affiliates. All rights reserved.



## Knowledge Bases pour Amazon Bedrock

Fournit aux FM et aux agents des informations contextuelles provenant de vos sources de données privées pour RAG afin de fournir des réponses plus pertinentes, précises et personnalisées.

Support entièrement géré pour le workflow RAG de bout en bout

Connectez en toute sécurité les FM et les agents aux sources de données

Récupérez facilement les données pertinentes et améliorez les prompts

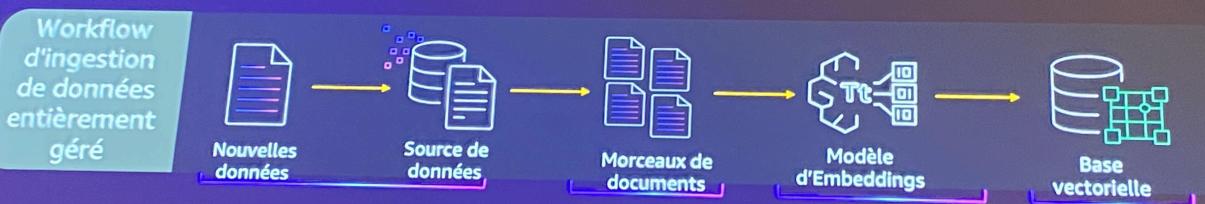
Indiquez l'attribution de la source



© 2024, Amazon Web Services, Inc. or its affiliates. All rights reserved.

# Workflow d'ingestion de données

KNOWLEDGE BASES SUR AMAZON BEDROCK

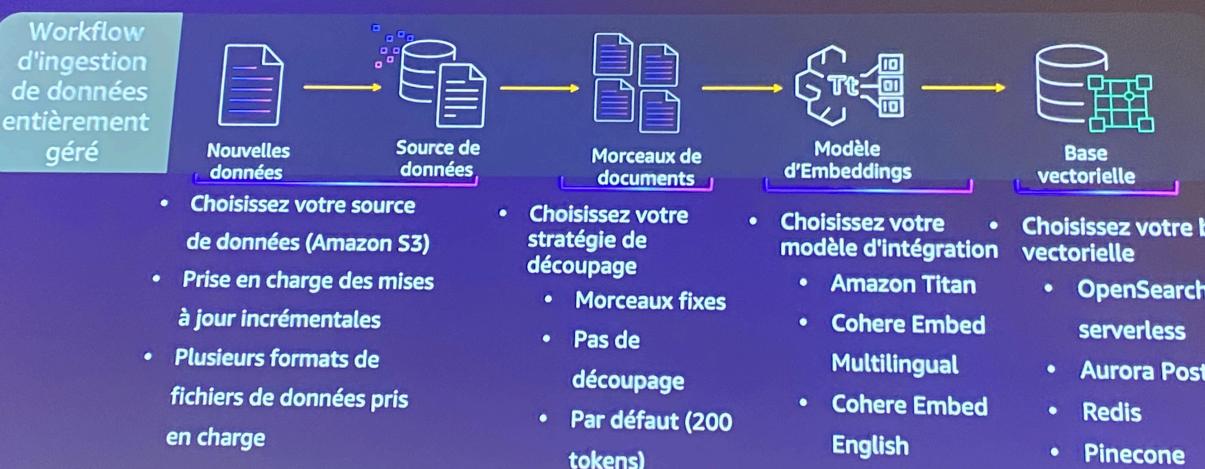


© 2024, Amazon Web Services, Inc. or its affiliates. All rights reserved.



# Workflow d'ingestion de données

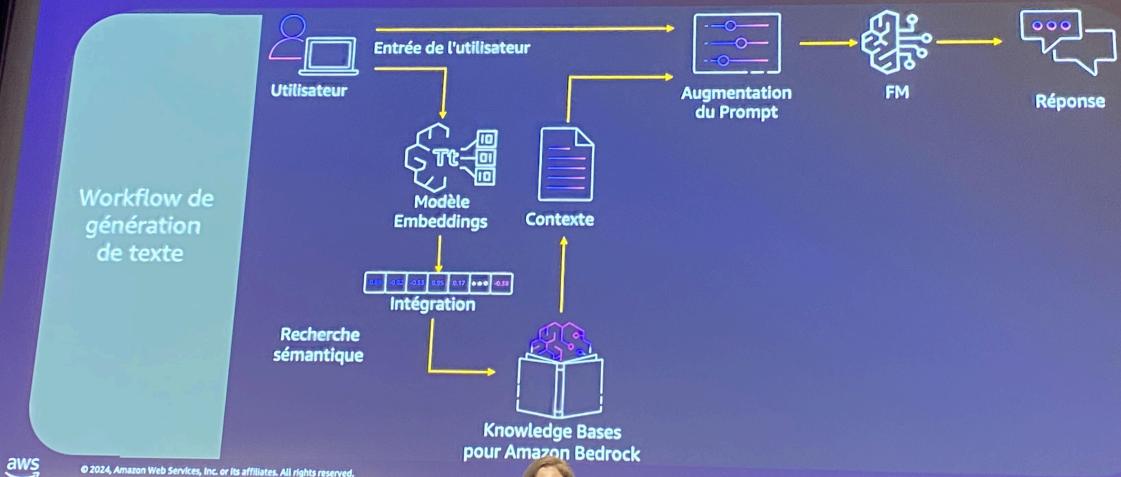
KNOWLEDGE BASES SUR AMAZON BEDROCK



© 2024, Amazon Web Services, Inc. or its affiliates. All rights reserved.



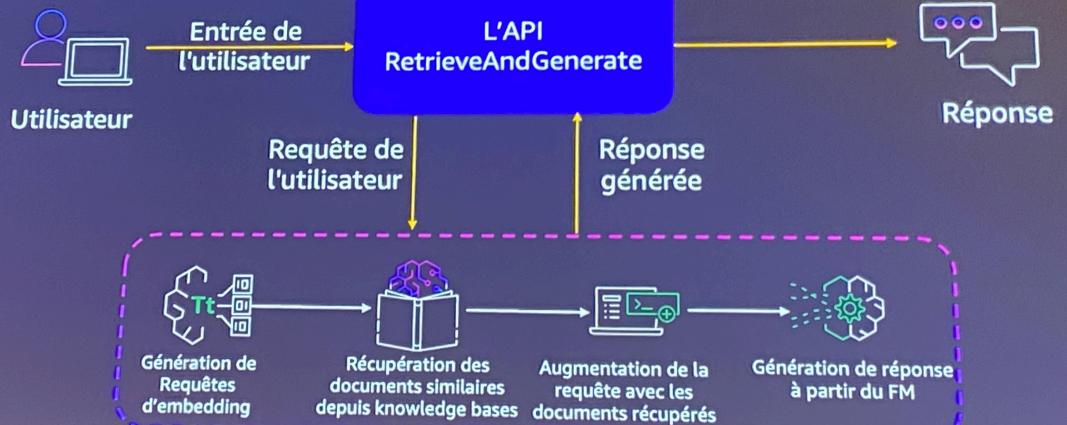
## Récupération et génération

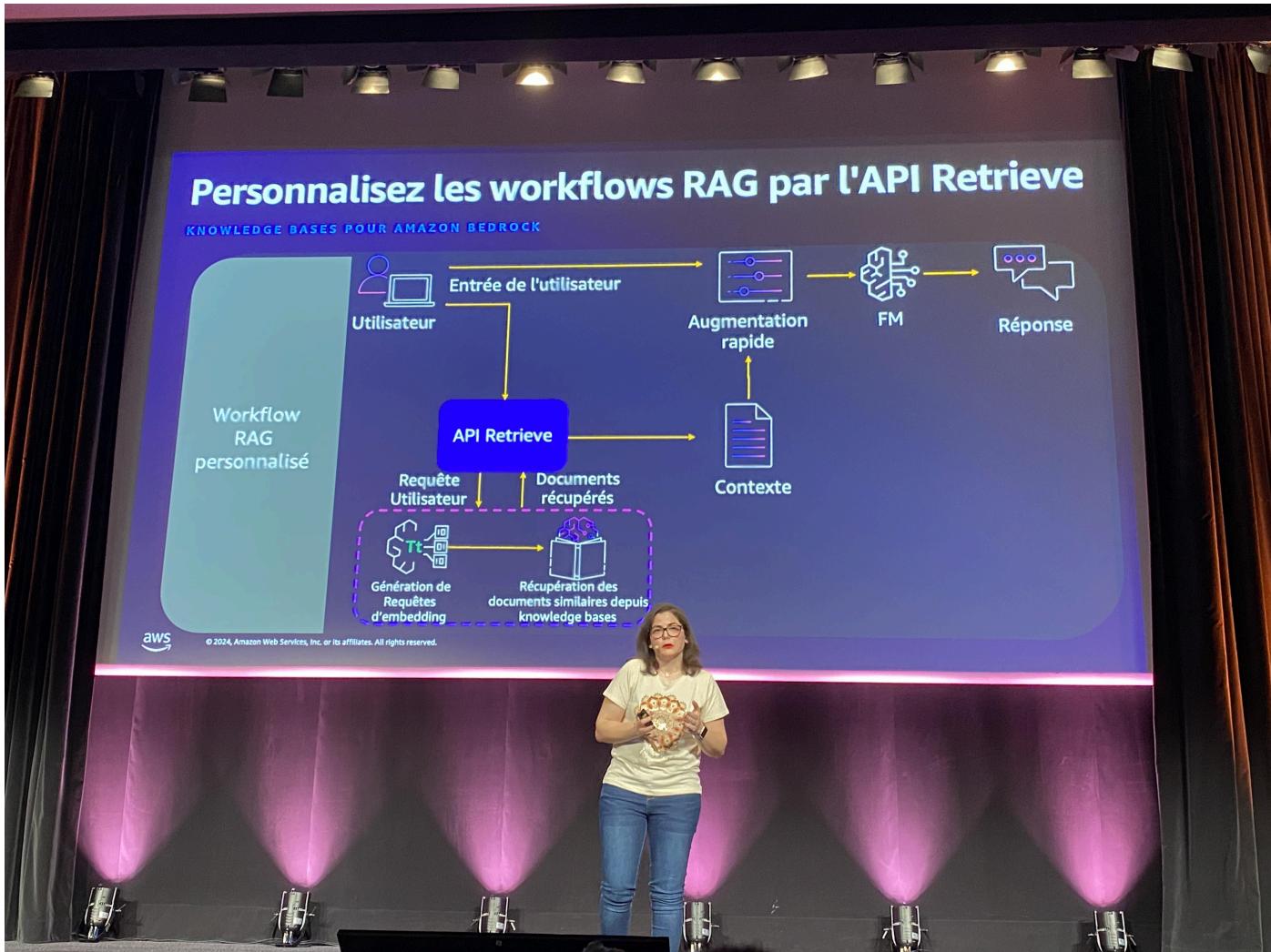


## L'API RetrieveAndGenerate

KNOWLEDGE BASES SUR AMAZON BEDROCK

RAG  
entièrement  
géré





## Comment rechercher ?

$$\text{Distance} = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2 + (z_2 - z_1)^2}$$

$$\text{simCos}(T, U) = (T \cdot U) / (\|T\| * \|U\|)$$

- Petits jeux de données
- Brute force
- Pre-filtre

Exact K-MN

- Jeux de données volumineux
- Post-filtre

Approximate k-NN (AKN)

aws s

