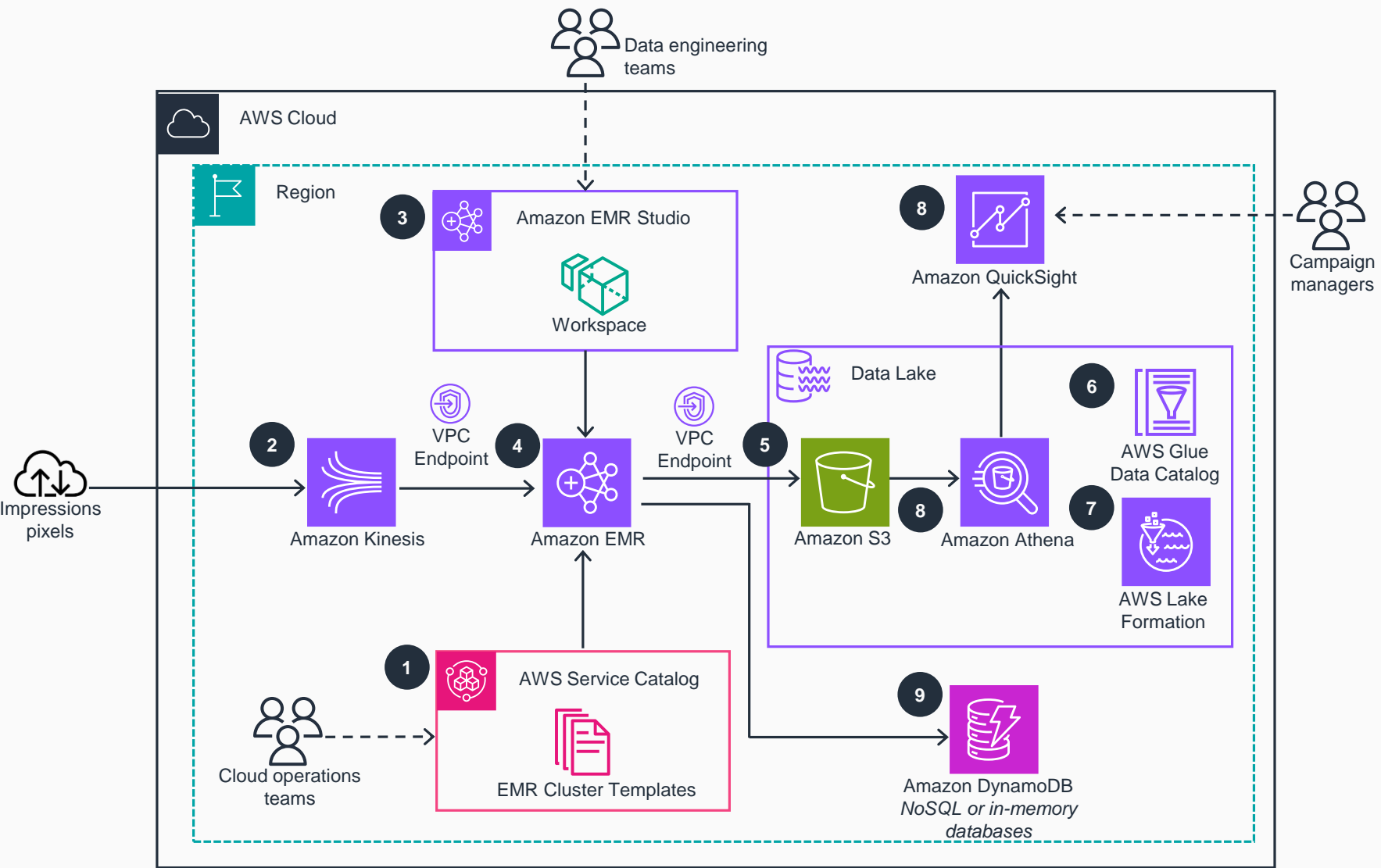


Guidance for Implementing Near Real-Time Analytics with Spark Streaming on AWS

Steps 6-9



- 6 All database and table metadata is registered within an **AWS Glue** Data Catalog, so data can be queried by multiple AWS services like **Amazon Athena** or **Amazon SageMaker**.
- 7 (Optional) Data lake administrators can register the Data Catalog with **AWS Lake Formation** to provide more granular access controls and centralize user management.
- 8 Users can run SQL queries against curated clickstream or impression data in **Amazon S3** in near real-time with **Athena** and visualize dashboards with **Amazon QuickSight**.
- 9 In addition to the **Amazon S3** data lake, **Amazon EMR** workloads can write data to NoSQL databases like **Amazon DynamoDB** or in-memory databases like Aerospike. This supports read workloads requiring fast performance on a large scale, such as bid filtering or operational reporting.

