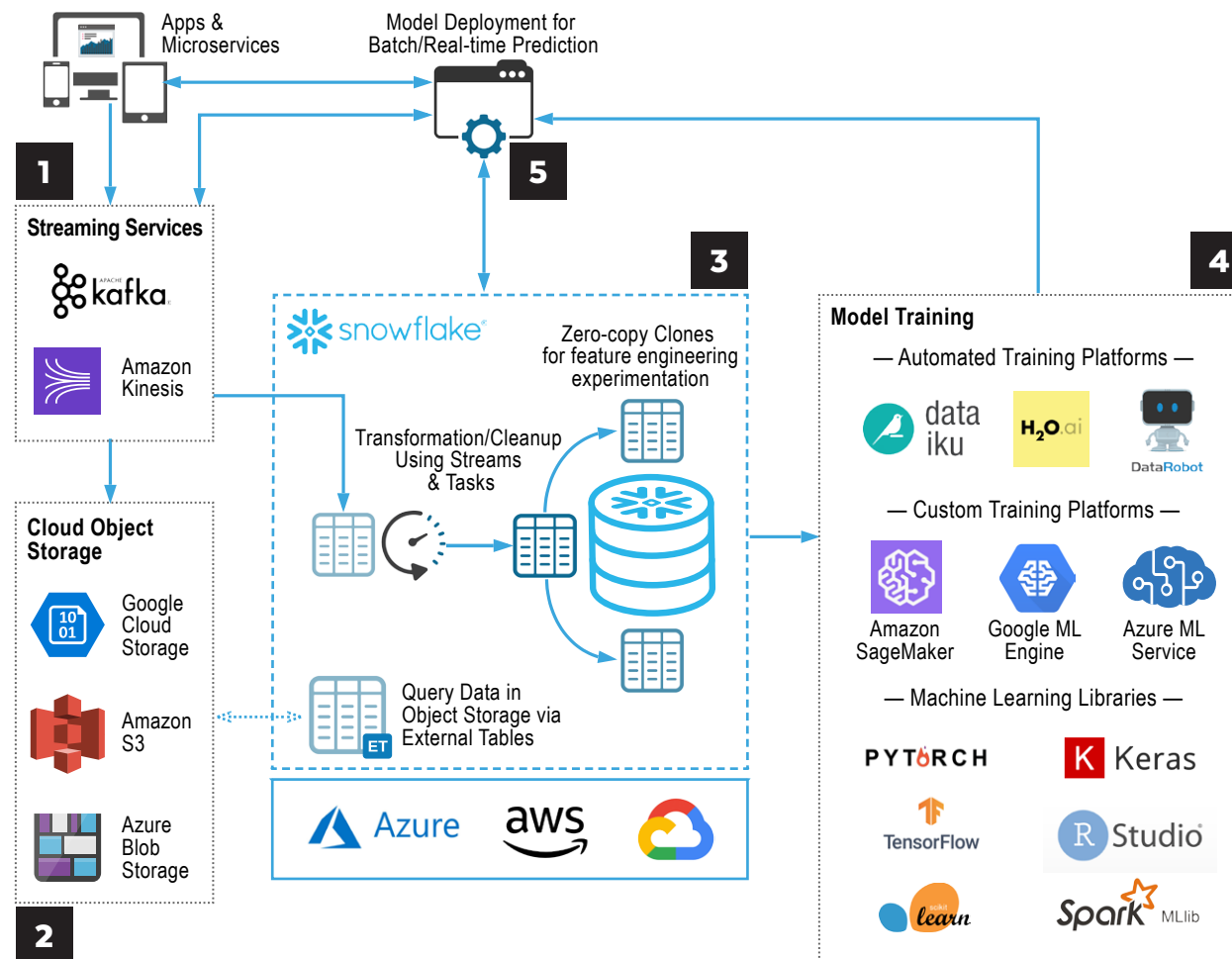


# MACHINE LEARNING AND DATA SCIENCE REFERENCE ARCHITECTURE

## MACHINE LEARNING AND DATA SCIENCE



### OBJECTIVE

Train machine learning (ML) models to build predictive applications, such as recommendation engines.

### DESCRIPTION

- 1** The application produces training data, which Snowflake (3) ingests via the streaming service or via cloud object storage (2). The streaming service buffers the training data to ensure reliable and continuous ingestion.
- 2** When cloud object storage is used, the streaming service batches training data into larger chunks to lower the API storage expenses.
- 3** Snowflake ingests data into a staging table. When new data is detected, the Streams and Tasks feature schedule required transformations. Multiple streams and tasks can be chained to implement a complex data pipeline. External Tables support queries of data in cloud object storage without ingestion. Data scientists can create zero-copy clones of the training data to support feature engineering and experimentation.
- 4** Using the data stored in Snowflake, data scientists train models with ML platforms and available libraries. Once the model artifacts are trained, they are deployed on the training platforms or on a separate process (5) to support predictions.
- 5** The application performs predictions in real time or schedules batch predictions using the deployed models. For batch predictions, data is read from an input table in Snowflake, and the results are stored in an output table where they are available to the application. In cases where subsecond response time is required, predictions can also be performed using input data from the streaming service.