

Q1) Design a cloud agnostic solutions for below scenario, you can make reasonable assumptions and list the same in your response.

Scenario

As a customer I would like retain sensitive data(big data) in my custody all times and only grant read-access to data to other party to run analytics over data and to note data at rest is encrypted. I'm looking for a system and architecture design that allows my business to securely share data with my customers with proper security, on-demand data access controls, usage per license, and customers shouldn't be able to copy my dataset. Expect a high level design solution with description of each components and data flows, focus on below.

- 1) Security of Solution
- 2) Data Lineage
- 3) Data updates and versions
- 4) Modern serverless or managed technology and frameworks used
- 5) Resilience of systems
- 6) Any other considerations you like to highlight

Q2) Write a ETL solution leveraging AWS EMR (not AWS Glue) using your language of choice that extracts data from below mentioned open dataset, transforms, and loads into big data store of your choice.

Output: (AWS)

- Share your ETL setup, code and documentation for the setup
- Open dataset in a JSON

OR

Write a ETL solution leveraging Google Cloud DataFlow for building data pipeline that extracts data from below mentioned open dataset, transforms, and loads into BigQuery.

Data Source

<https://data.nsw.gov.au/data/dataset/formal-gipa-access-application-2016-2017-fa-13>

Resource

<https://data.nsw.gov.au/data/dataset/4e51f1d3-4b72-48b8-96ef-8493b7aa9c37/resource/ca5c4a11-64f8-4583-91c9-d9436d38e2e9/download/fa13-2016-17.xlsx>

Output: (Google Cloud)

- Share your Google Project setup, code and documentation for the setup
- BigQuery Query or Open dataset in a JSON on cloud storage.