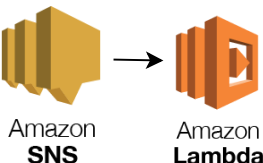
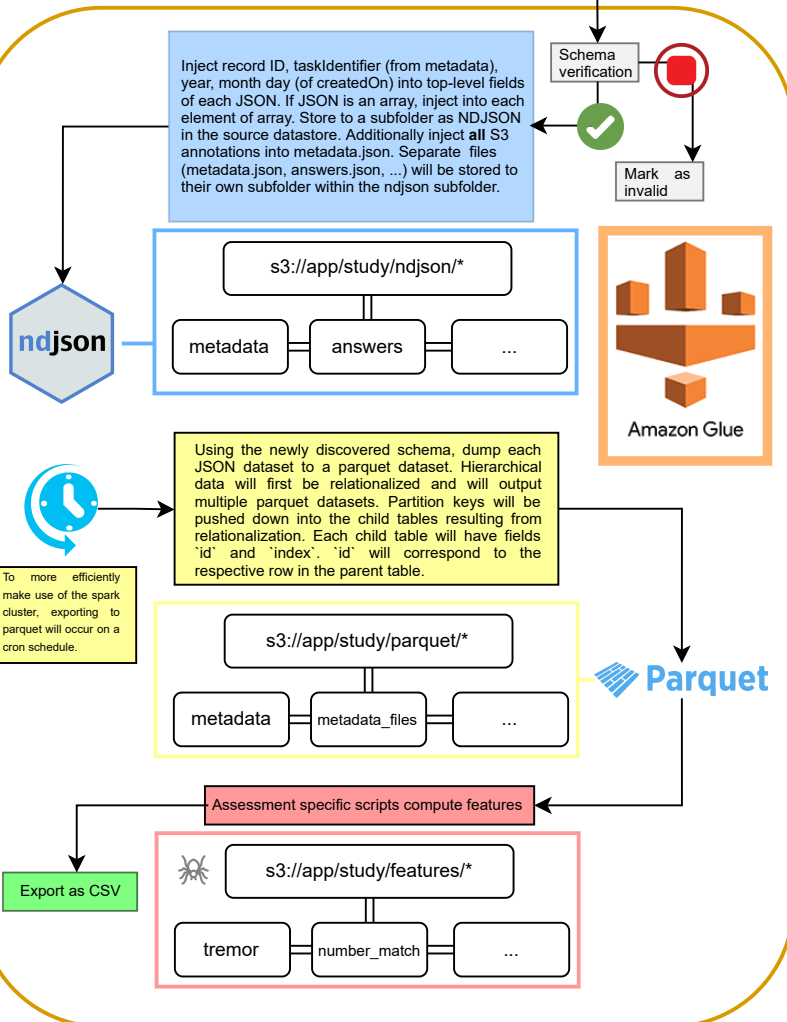


Bridge publishes a message to an app and study-specific SNS topic with information about the received data. A lambda subscribed to this topic triggers a Glue workflow.



"S3 will be used to store raw data, which will generally be unprocessed zip files submitted by the app. We will use [app/study] as prefixes (S3 folders) so that we can grant permissions for entire studies and substudies as needed."

"Record metadata (such as record ID, participant ID (healthcode) / version, assessment ID / revision / guid, user agent, created-on, exported-on, etc) will be stored as S3 metadata. This allows researchers looking at the health data in S3 to know the metadata for a given health record."



Parquet files are accessed via Synapse by setting the parquet S3 location as an external storage location on a Synapse folder. Files will be downloaded directly from S3 using an STS token.

CSV files can be uploaded directly to a Synapse folder (potentially one with an external storage location)



Glue crawlers are responsible for updating the table partitions with the parent table metadata. They can also be used to discover Glue table schemas over JSON or Parquet data. We will not be updating table schemas automatically, so any schema changes in the JSON will need to go through a review process to ensure that they are compatible with the old schema, then the new Glue table schemas will be updated in advance of the schema changes going into production. If a schema change is incompatible with the old schema – for example, if the datatype of a field changes – we will need to act as if the new schema data comes from a new JSON dataset and output to a separate parquet dataset.



There are two crawlers per study: One for JSON objects and one for JSON arrays. They run once a day and do not update table definitions.