

**Writing Prompt:**

How would you build a service that offers Stable Diffusion as an API?

The architecture for the system development offering “stable diffusion as an API” would be suggested to be comprised of microservices, and the multiple microservices will be containerized orchestrated. It is assumed that the API will be used by two types of users, normal users, and premium users. Supposedly the premium users will get the inference at the request time and also can request the outputs to be emailed or direct download, but the normal users will be getting the inference after a small time, e.g., 30 minutes and the output will be sent to the email addresses. The system to be developed would be comprised of the following parts,

**1. Model Inference Instance:**

As the main task is to produce Stable Diffusion inferences, therefore a GPU-based compute-optimized instance would be required that can take produce inferences. This will be a separate container and would be responsible for producing inferences. This service will be responsible for model weights availability (download model weights from the web if not available), handling batch inference, and other processes related to model inference such as preprocessing.

**2. APIs Web Framework:**

This would be a separate microservice and main driver for the system that will be providing the interface to the client for communicating with the system, it will be providing API docs, validating premium users, validating the client’s input, registering the new user in the system, requesting model inference container for the inference on the input according to the user type, and storing emails and other required metadata in the database.

There will be three APIs, one for requesting the model inference containing clients metadata in the request, and if the client’s request is of a normal user, the request will be logged in the database, otherwise will be sent for the inference, with the information being the inference request by the premium user, and can also contain metadata information for the model inference.

The other API will be for the registration request or sign-up from the user, this can be provided on the basis of some criteria e.g., buying the API for some time, or providing an email address, etc. Upon signup request, the API will validate the sign-up request and logs the user in the database record with the metadata.

The other API will be for sign-in, which validates the sign-in request by the database, make a session, or provide a token for further requests.

For database access, a database connection would be made if the connection with the database connection is not established first, and neither is stored as a persisted variable.

**3. Persistent storage, Database**

This will be a persistent data storage that will be storing the metadata, username, and passwords for the users, and user\_input\_data for batch inference.

The API web framework will be the interface, it will be responsible for validating the users, registering new users, handling batch and stream inference according to the user type, and communicating with the client and other containers, thus handling all the requests.

The advantage of container orchestration is that it provides auto-scaling, load balancing, restarting containers upon crashing, vpc, etc. The advantage of using microservices instead of a monolithic architecture is that autoscaling would be easier and more efficient, moreover memory utilization due to containerization would also be enhance and well performed.