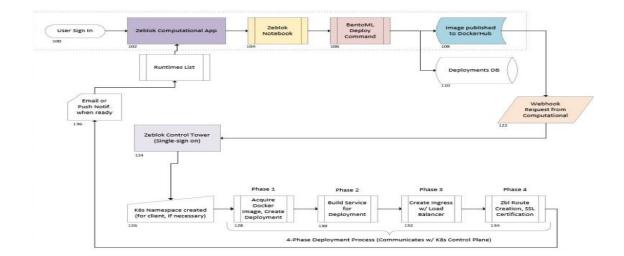# ZEBLOK
## COMPUTATIONAL
# Ai-API™ Engine

## Overview

Getting machine learning models into production requires packaging the model as a docker container and translation into a standard application programming interface (API) – it is complex, with many steps. Data scientists are experts in developing and training AI/ML models, but not at building production services and DevOps best practices and they find it challenging to test and deploy trained models. Generally, they choose to hand the task over to a software engineering to avoid a time consuming and error-prone workflow.

Zeblok's Ai-API™ Engine is an Ai-MicroCloud™ framework, which translates a machine learning model into an API and then enables enterprises to deploy and manage each Ai-API™. It bridges the gap between data science and DevOps, enabling fast, repeatable and scalable delivery of prediction services.

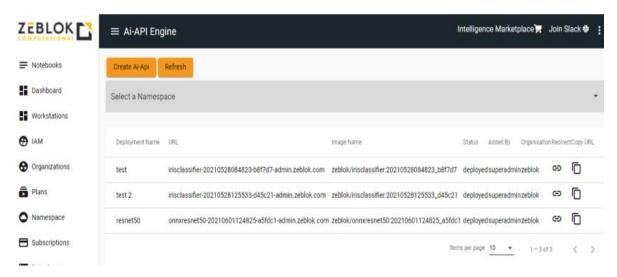Zeblok's Ai-API™ Engine makes moving trained ML models to production easy, performing the following:

- Package models trained with ML framework and then containerize the model server for production deployment as an Ai-API™

- Deploy Ai-API™ anywhere via either online API serving endpoints or offline batch inference jobs

- High-performance API model server, with adaptive micro-batching support

- Ai-API™ server is able to handle high-volume without crashing, supports multi-model inference, API server Dockerization, Built-in Prometheus metric endpoint, Swagger/Open API endpoint for API Client library generation, serverless endpoint deployment etc.

- Central hub for managing models and deployment process via web UI and APIs

- Supports various ML frameworks including: Scikit-Learn, PyTorch, TensorFlow 2.0, Keras, FastAI v1 & v2, XGBoost, H2O, ONNX, Gluon and more

- Supports API input data types including: DataframeInput, JsonInput, TfTensorflowInput, ImageInput, FileInput, MultifileInput, StringInput, AnnotatedImageInput and more

- Supports API output adapters including: BaseOutputAdapter, DefaultOutput, DataframeOutput, TfTensorOutput and JsonOutput

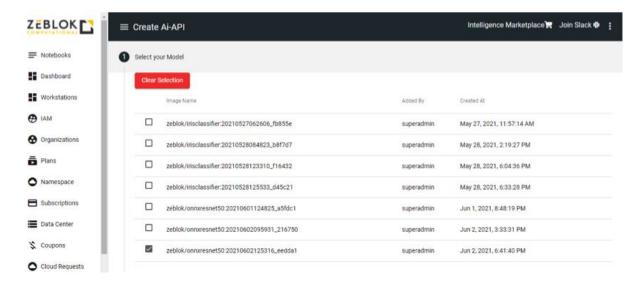### Complete Workflow of Zeblok Ai-API™ Deployment

At the conclusion of the process, each successfully deployed Ai-API™ is listed separately, with its corresponding URL as shown below.
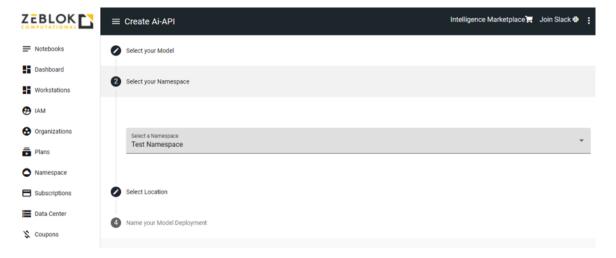


# Process
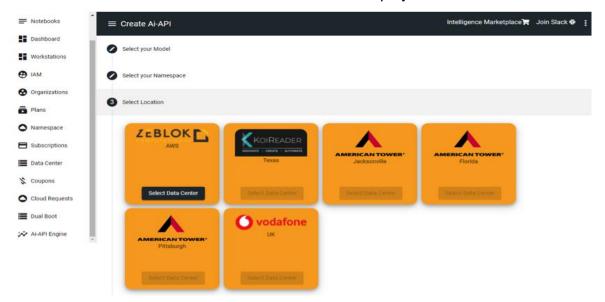
Step 1: Select the model to deploy as an Ai-API™

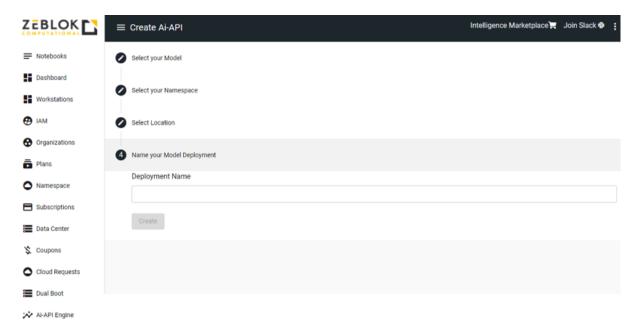Step 2: Option to select the namespace



Step 3: Select the data center locations where model is to be deployed

Step 4: Enter a unique name for the Ai-API™ deployment



**For more information: email Mouli Narayanan**

**Zeblok Computational Inc.**
1500 Stony Brook Road
Stony Brook, NY 11794
www.zeblok.com
mouli.narayanan@zeblok.com
Phone: +1 (631) 223-8233