William Chang Liu - 1009048852
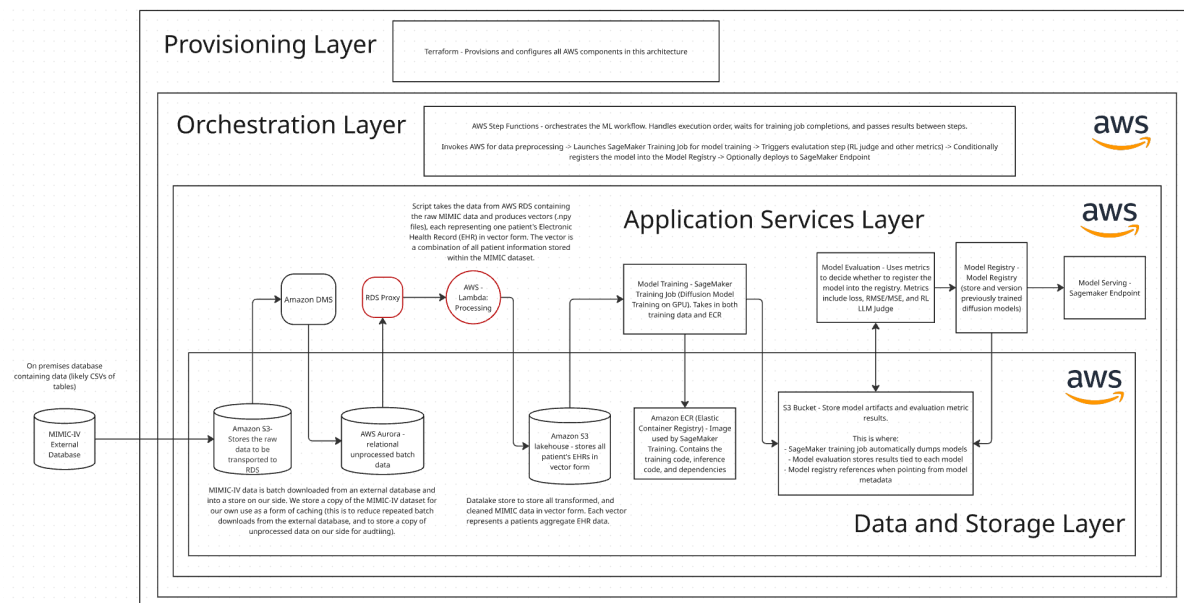
Ryan Zhang - 1009020453

Taiyi Jin - 1009075796

Abdus Shaikh - 1007288533

# Assignment 3: Deleting Prod

## Part 1: Diagramming Infrastructure



Miro Link:
https://miro.com/app/board/uXjVJDh--SI=/?share_link_id=501292272666

The components highlighted in red will be implemented in the future. This was decided because our data processing infrastructure is subject to the most change as we try to tackle ways to add a temporal dimension to our model. Furthermore, the data format of the EHRDiff model, which we used as a placeholder, is not compatible with the current output of the data-processing script. As a result, the data-ingestion pipeline cannot be connected to the model serving and training pipeline at the moment.

Given that the group had access to only AWS free-tier for the majority of the development process some compromises had to be made. Namely, the model training algorithm is currently cpu-only as the group did not have access to EC2 instances with GPUs. Creating a good model necessitates GPU access and more work will be done to create a deployable GPU training script in the future.

Since this is a research-oriented project and we have no client-facing application, we do not have separate environments for development and production. The decision to have one simple

data-processing to model training and evaluation environment was made because it allows changes to be made and applied more quickly.

Part 2: Infrastructure as Code

Terraform was used as the primary framework for Infrastructure as code. No scripts were used and all infrastructure is defined declaratively. All of the terraform code can be found under a3/terraform. a3/terraform/model_train_serve contains all of the infrastructure for model training and serving. a3/terraform/data_ingestion contains the work-in-progress code for our data processing pipeline.

Several pieces of open-source code were modified and used to create our infrastructure. Firstly, an official AWS blog by Oliver Zollikofer provided us with a guide as well as a repository to starter code that served as the foundation for our model training and serving IaC. Secondly, our demonstration model was based on EHRDiff created by Yuan et al.

Part 3: Deleting Production

The video for deleting production can be found on this google drive link: 🎬 a3_demo_video.mp4 .
It will also be included in the quercus submission.

References

Yuan, H., Zhou, S., & Yu, S. (2023). *EHRDiff: Exploring realistic EHR synthesis with diffusion models* (arXiv preprint arXiv:2303.05656). https://doi.org/10.48550/arXiv.2303.05656

Yuan, H., Zhou, S., & Yu, S. (2023). *EHRDiff*  [Source code]. https://github.com/sczzz3/EHRDiff

Zollikofer, O. (2022, May 4). *Deploy and manage machine learning pipelines with Terraform using Amazon SageMaker*. AWS Machine Learning Blog. https://aws.amazon.com/blogs/machine-learning/deploy-and-manage-machine-learning-pipelines-with-terraform-using-amazon-sagemaker/

Zollikofer, O.  (2022, May 4). *amazon-sagemaker-ml-pipeline-deploy-with-terraform*  [Source code]. GitHub. https://github.com/aws-samples/amazon-sagemaker-ml-pipeline-deploy-with-terraform