A cloud-based pipeline for analysis of FHIR and long-read data



Tim Dunn

▲ University of Michigan

Introduction

As genome sequencing becomes cheaper and more accurate, it is becoming increasingly viable to merge this data with electronic health information to inform clinical decisions. In this work, we demonstrate a full pipeline for working with PacBio sequencing data and clinical FHIR data, from initial data to tertiary analysis. The electronic health records are stored in FHIR - Fast Healthcare Interoperability Resource - format, the current leading standard for health care data exchange. For the genomic data, we perform variant calling on long read PacBio HiFi data using Cromwell on Azure. Both data formats are parsed, processed, and merged in a single scalable pipeline which securely performs tertiary analyses using cloud-based Jupyter notebooks. We include three example applications: exporting patient information to a database, clustering patients, and performing a simple pharmacogenomic study.

| ##filej | Format=V | CFv4 | . 3 | | | | | | | | | 'resource': {'resourceType' | | | | | |
|---|---|--------|-------|----------|-------|----------|-----------------|-----------|-------------|------------------|-------------------------|---|--|----------------------------------|--|--|--|
| ##fileDate=20210115 ##source=VariantCallerScript | | | | | | | | | | | 'id': 'b5f1da11-3826-48 | | | | | | |
| <pre>##source=VariantCallerScript ##reference=file:///references/HG38.fasta</pre> | | | | | | | | | | | 'meta': {'versionId' | | | | | | |
| ##refer | rence=fi | le:/, | //rej | ference | s/HG3 | 8.fasta | | | | | | 'lastUpdated': '2022-0 | | | | | |
| ##conti | g= <id=2< td=""><td>0, Lei</td><td>ngth</td><td>=624359</td><td>64,as</td><td>sembly=E</td><td>336,species="Ho</td><td>no sapien</td><td>s"></td><td></td><td></td><td>'profile': ['http://h]</td></id=2<> | 0, Lei | ngth | =624359 | 64,as | sembly=E | 336,species="Ho | no sapien | s"> | | | 'profile': ['http://h] | | | | | |
| ##phasi | ing=part | ial | | | | | | | | | | 'text': {'status': 'ger | | | | | |
| ##INFO= | <id=dp,< td=""><td>Numb</td><td>er=1,</td><td>Type=I</td><td>ntege</td><td>r,Descri</td><td>ption="Total D</td><td>epth"></td><td></td><td></td><td></td><td>'div': '<div allele="" frequency"="" xmlns="ht</td></tr><tr><td colspan=9><pre>##INFO=<ID=AF,Number=A,Type=Float,Description="></div></td><td></td><td colspan="4" rowspan="3"><pre>'extension': [{'exten</pre></td></id=dp,<> | Numb | er=1, | Type=I | ntege | r,Descri | ption="Total D | epth"> | | | | 'div': ' <div allele="" frequency"="" xmlns="ht</td></tr><tr><td colspan=9><pre>##INFO=<ID=AF,Number=A,Type=Float,Description="></div> | | <pre>'extension': [{'exten</pre> | | | |
| ##INFO= <id=aa,number=1,type=string,description="ancestral allele"=""></id=aa,number=1,type=string,description="ancestral> | | | | | | | | | | | | | | | | | |
| <pre>##INFO=<id=db,number=0,type=flag,description="dbsnp 129"="" build="" membership,=""></id=db,number=0,type=flag,description="dbsnp></pre> | | | | | | | | | | | | | | | | | |
| ##FILTER= <id=q10,description="quality 10"="" below=""></id=q10,description="quality> | | | | | | | | | | 'display': 'Whit | | | | | | | |
| ##FILTE | R= <id=s< td=""><td>50,D</td><td>escr</td><td>iption=</td><td>"Less</td><td>than 50</td><td>% of samples h</td><td>ave data"</td><td>></td><td></td><td></td><td>{'url': 'text', 'val</td></id=s<> | 50,D | escr | iption= | "Less | than 50 | % of samples h | ave data" | > | | | {'url': 'text', 'val | | | | | |
| ##FORM4 | AT= <id=g< td=""><td>T, Nui</td><td>nber</td><td>=1, Type</td><td>=Stri</td><td>ng,Descr</td><td>ription="Genoty</td><td>oe"></td><td></td><td></td><td></td><td>'url': 'http://hl7.or</td></id=g<> | T, Nui | nber | =1, Type | =Stri | ng,Descr | ription="Genoty | oe"> | | | | 'url': 'http://hl7.or | | | | | |
| ##FORM4 | AT= <id=di< td=""><td>P,Nu</td><td>mber</td><td>=1, Type</td><td>=Inte</td><td>ger,Desc</td><td>ription="Read </td><td>Depth"></td><td></td><td></td><td></td><td>{'extension': [{'url':</td></id=di<> | P,Nu | mber | =1, Type | =Inte | ger,Desc | ription="Read | Depth"> | | | | {'extension': [{'url': | | | | | |
| ##FORM4 | AT= <id=h< td=""><td>Q, Nu</td><td>nber</td><td>=2, Type</td><td>=Inte</td><td>ger,Desc</td><td>ription="Haplo</td><td>type Qual</td><td>ity"></td><td></td><td></td><td>'valueCoding': {'s</td></id=h<> | Q, Nu | nber | =2, Type | =Inte | ger,Desc | ription="Haplo | type Qual | ity"> | | | 'valueCoding': {'s | | | | | |
| #CHROM | POS | ID | REF | ALT | QUAL | FILTER | INFO | FORMAT | NA00001 | NA00002 | | | | | | | |
| 20 | 14370 | rs6 | G | A | 29 | PASS | DP=14;DB; | GT:DP:HQ | 0 0:1:51,51 | 1 0:8:51,51 | | 'code': '2135-2', | | | | | |
| 20 | 17330 | • | Т | A | 3 | q10 | DP=11;AF=0.017 | GT:DP:HQ | 0 0:3:58,50 | 0 1:5:65,3 | | 'display': 'Hispar | | | | | |
| 20 | 1110696 | rs5 | A | G,T | 67 | PASS | DP=10;AA=T;DB | GT:DP:HQ | 1 2:6:23,27 | 2 1:0:18,2 | | {'url': 'text', 'va] | | | | | |
| 20 | 1230237 | | | т. | 47 | PASS | DP=13;AA=T | GT:DP:HQ | 0 0:7:56,60 | 0 0:4:51,51 | | 'url': 'http://hl7.or | | | | | |
| 20 | 1234567 | ms1 | GTC | G,GTCT | 50 | PASS | DP=9;AA=G | GT:DP | 0/1:4 | 0/2:2 | | {'url': 'http://hl7.or 'valueString': 'Julia | | | | | |

Fig. 1: Example VCF file, including the file header and variant call data.

': 'Patient',

- nerated',
- ttp://www.w3.org/1999/xhtml">Generated by Synthea</div>'}, ion': [{'url': 'ombCategory' /stem': 'urn:oid:2.16.840.1.113883.6.238',
- lueString': 'White'}],
- rg/fhir/us/core/StructureDefinition/us-core-race'}, : 'ombCategory', ystem': 'urn:oid:2.16.840.1.113883.6.238',
- nic or Latino'}}, lueString': 'Hispanic or Latino'}], rg/fhir/us/core/StructureDefinition/us-core-ethnicity'},
- rg/fhir/StructureDefinition/patient-mothersMaidenName', a241 Luna60'}.
- Fig. 2: Example synthetic FHIR data, generated using Synthea [1].

Methodology

Cromwell is a workflow management system for running genomics analysis scripts at various scales, from a local machine or computing cluster to larger cloud instances. Cromwell was originally developed by the Broad Institute, and is used in the Genome Analysis Toolkit's (GATK) recommended "Best Practices" pipeline for genome analysis. "Cromwell on Azure" is a project developed by Microsoft that configures all the Azure resources necessary to run Cromwell workflows on the Azure cloud. The first stage of our pipeline uses Cromwell on Azure to perform variant calling on long read PacBio HiFi reads. Although this whole genome sequencing pipeline also calls structural variants, at the moment only small germline variants are used in downstream analyses. An overview of how Cromwell on Azure works is shown below, in Figure 4.

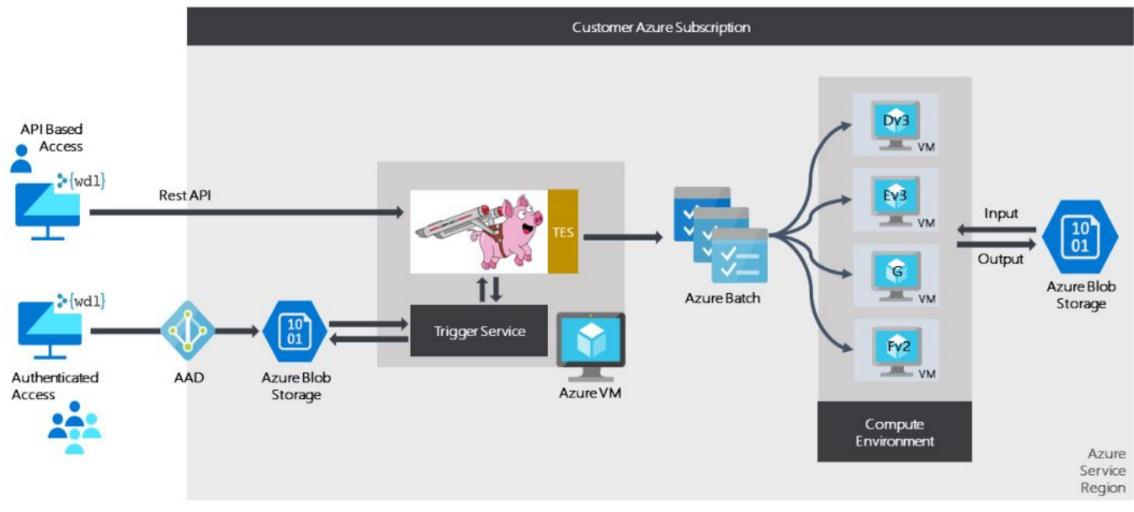


Fig. 4: High-level overview of Microsoft's "Cromwell on Azure" service, used to perform variant calling [2].

In order to ensure that all data processing is secure, our pipeline works within a Jupyter notebook hosted on an Azure "Confidential Compute" virtual machine. Results are then made available on a local machine using SSH tunneling. This setup can be seen below, in Figure 6. Confidential virtual machines have security features such as secure boot, a virtual trusted platform module (vTPM), boot integrity monitoring, and virtualization-based security. Together, these features protect against persistent or advanced threats such as rootkits. Virtualization provides further security by isolating memory address spaces to remove any possibility of memory cross-contamination.

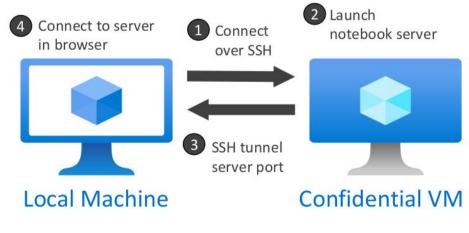
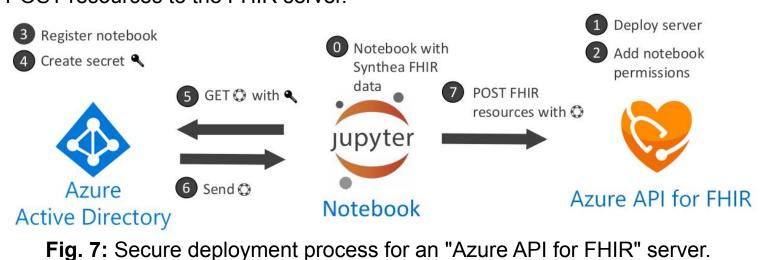


Fig. 6: Confidential computing configuration.

Ensuring the confidentiality of patient health information is paramount. In order to securely upload FHIR data to a server and later query this data, we must have some means of performing authentication. Azure Active Directory is a cloud-based identity and access management service which we selected to use for authentication purposes. After deploying a FHIR server, we added permissions for our notebook script to write data to the server, registering it with Azure Active Directory and providing it with a secret key (generated by Azure AD) as proof-of-identity. In order to upload data to the FHIR server, the notebook requests a secret token from Azure AD using its secret key, and then uses that token to POST resources to the FHIR server.



Erdal Cosgun

Generation

17.org/fhir/us/core/StructureDefinition/us-core-patient']},

Overview

- Our cloud-based pipeline contains the following steps: 1. Variant Calling: We first run a scalable cloud-based variant
- calling pipeline using "Cromwell on Azure" [2]. **2. FHIR Data Generation:** Synthetic FHIR patient data is generated using Synthea [1].
- **3. Deploy FHIR Server:** A FHIR server is deployed to securely host patient data which can be accessed using an Azure API.
- **4. Upload FHIR Data:** A script is run to upload all synthetic generated patient data to the FHIR server. 5. Deploy Sync Agent to Download FHIR Data: The "Sync Agent
- [3] function app downloads data from the FHIR server and stores it in Parquet format.
- 6. Extract FHIR Data: A script extracts the desired clinical information from the Parquet files into a data frame.
- 7. Merge VCF Files: Variant calls for all patients are merged into a single file for efficient processing.
- 8. Parse VCF Data: A script extracts the desired variant call information from the merged VCF file into a data frame.
- 9. Merge FHIR and VCF Data: The data frames containing patient clinical and genomic data are merged into a single data frame
- **10. Analysis:** Several potential downstream applications are demonstrated.

Synthea is a widely-used open source tool for generating realistic (but synthetic) patient data in FHIR format [1]. This enables researchers to work with realistic clinical datasets without worrying about any of the legal, ethical, or security concerns that would accompany working with real patient data. Synthea works by first using general census demographic data in combination with user-specified configuration information to generate a synthetic world population. Once this population has been generated, disease incidence and prevalence statistics are used in tandem with disease-specific models to simulate patients contracting a given disease and all subsequent interactions with the healthcare system that will be stored in clinical records. This information can then be exported in FHIR format.

In addition to Synthea's default disease modules, for our third example pipeline application (see "Results") we created a custom Synthea module to model a simplified breast cancer patient cohort. A graphical depiction of this module is shown below, in Figure 5. For this cohort, patients had a default survival rate of 50%. Based on the combination of their genetics (the presence of a particular variant) and their prescribed medication (Epirubicin or Doxorubicin), their survival rate either remained the same or increased to 90%. This models the fact that medication effectiveness can depend upon the genetics of the patient involved.

| | Breast Cancer | | | | | | | |
|---|---|--|--|--|--|--|--|--|
| ConditionOnse | SNOMED-CT[254837009]: Malignant neoplasm of breast (disorder) Diagnose at Assign to Attribute: 'simple_breast_cancer' | | | | | | | |
| | 50% | 50% | | | | | | |
| | Doxorubicin Prescribed | Epirubicin Prescribed | | | | | | |
| MedicationOrder | RxNorm[1790099]: 10 ML Doxorubicin Hydrochloride 2 MG/ML Injection Reason: Breast Cancer Assign to Attribute: 'doxorubicin' | MedicationOrder RxNorm[1732186]: 100 ML Epirubicin Hydrochloride 2 MG/ML Injection Reason: Breast Cancer Assign to Attribute: 'epirubicin' | | | | | | |
| | 50% 50% | 50% 50% | | | | | | |
| a un a de la de | ubicin with Variant Doxorubicin without Variant RxNorm[ALT]: variant present MedicationOrder RxNorm[REF]: variant | | | | | | | |
| | 50% | 10% 10% 50% | | | | | | |
| | | Death | | | | | | |

Fig. 5: Custom Synthea module which models breast cancer patients' survival rates as a function of their medications and genetics.

In order to perform data science applications with FHIR data, we first need to convert the hierarchical JSON data into tabular format. An existing open source tool called the "FHIR to Synapse Sync Agent" solves this problem by downloading FHIR data from the server and converting it to Parquet format. Figure 8 shows an overview of how the Sync Agent works. Parquet files store the same information as ordinary CSV ("comma-separated values") or TSV ("tab-separated values") files, but in a more efficient compressed manner. Data is stored in column-major order, which allows compression algorithms such as run length encoding, dictionary encoding, or delta encoding to be applied per column depending upon each column's format and values.

Once downloaded, we parse the Parquet FHIR data to retrieve relevant information and load it into a data frame. Although Parquet is already a tabular format, the converted data contains extraneous information which must be discarded. Furthermore, FHIR data pertaining to a single patient is stored in multiple resource types. These files (such as Patient, Medication, or Condition resources) must be parsed separately and the records associated with one another using the patient's unique Medical Record Number (MRN). In the end, each row stores data for a single patient, and the numerous columns contain all desired information about that patient, extracted from the Parquet files. At this point the FHIR and PacBio data can be merged together, as shown in Figure 9.

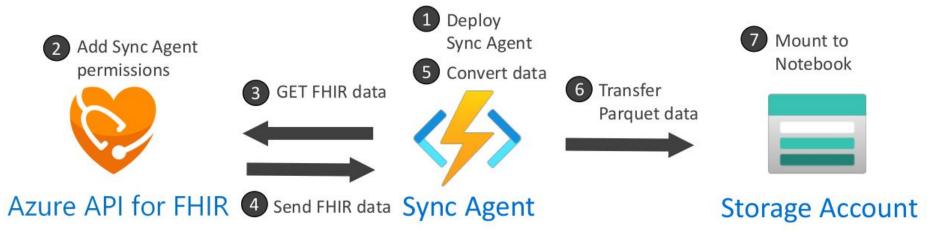


Fig. 8: Deployment process for a "Sync Agent" which downloads data from the FHIR server [3].

Microsoft Biomedical Platforms and Genomics

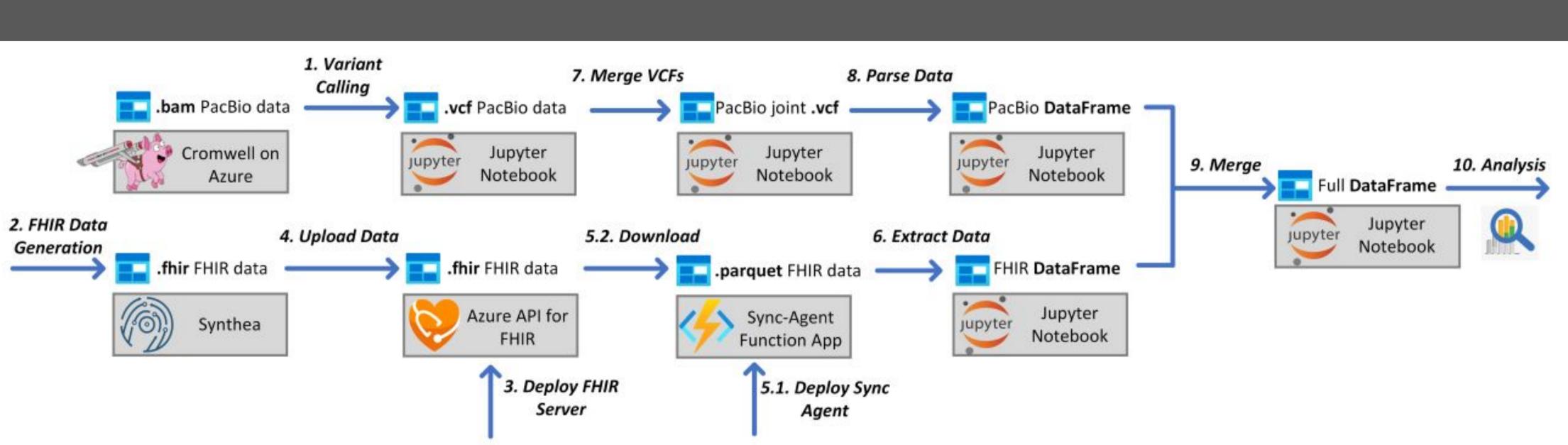


Fig. 3: Overview of data processing pipeline, common to all demonstrated downstream applications.

Results

Application #1 – Data Exportation: The first and simplest application we demonstrated was exporting the merged FHIR and PacBio data to a database for further analyses. We selected Azure Synapse Analytics as our database platform because it ensures scalability and reliability.

| | city | state | country | gender | dead | age | newest_med_code | | 49:GT | 49:AF_0 | 49:AF_1 | 49:PL_0/0 | 49:PL_0/1 | 49:PL_1/1 | 49:DP |
|---|-----------|-------|---------|---------|-------|------------|--------------------|-----|--------|----------|------------|-------------|-----------|-----------|-------|
| 0 | Boston | MA | US | female | False | 22.009037 | 748856 | | 0/0 | 1 | 0 | 0 | 8 | 8 | 7 |
| 1 | Amesbury | MA | US | male | False | 64.588092 | 310798 | | 0/0 | 1 | 0 | 0 | 3 | 3 | 7 |
| 2 | Boston | MA | US | female | False | 17.316361 | 1367439 | | 0/1 | 0.285714 | 0.714286 | 29 | 0 | 23 | 7 |
| 3 | Yarmouth | MA | US | male | True | 69.918688 | 896209 | | 0/0 | 1 | 0 | 0 | 4 | 4 | 5 |
| 4 | Attleboro | MA | US | female | False | 25.182206 | 751905 | | 0/0 | 1 | 0 | 0 | 9 | 9 | 15 |
| | | | Eia | O. Line | data | fromo of r | norged elipical EL | חוו | data a | nd DooD | ie verient | colling dot | | | |

Application #2 – Patient Clustering: The second application we demonstrated was an exploratory clustering of our patient dataset using three different clustering methods: K-Means++, DBSCAN, and Spectral Clustering. There is no ground truth available, so this is only useful for discovering similarities between patients.

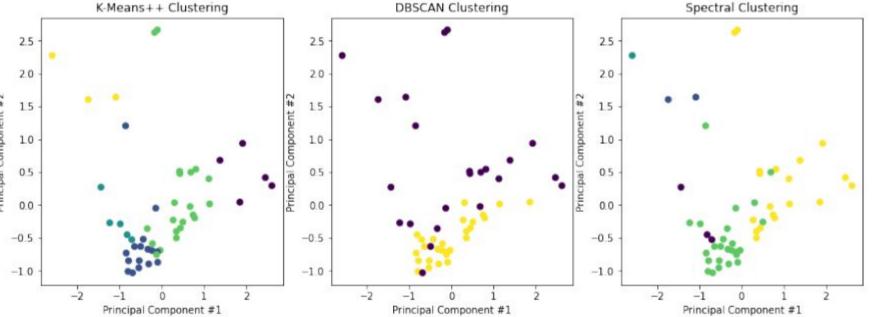


Fig. 10: Patient clusters using merged clinical and genomic data, resulting from several different clustering algorithms.

Application #3 – Pharmacogenomic Study: The final application we demonstrated was a basic pharmacogenomic study to determine the effect of two different medication treatments on breast cancer patients' survival rates. First, a custom Synthea module was designed which modelled a patient population where outcomes were only improved beyond the base 50% survival rate for patients who were prescribed Epirubicin and had a particular genetic variant, and for patients who were prescribed Doxorubicin and did not have the genetic variant. This model is depicted in Figure 5, and uses both clinical and genomic data. Below, Figure 11 shows a breakdown of patient variant presence, treatment type, and survival.

| Patients | Variant |
|----------|------------|
| | No Variant |

Fig. 11: Sankey diagram showing distribution of patient treatments and outcomes.

More Information

References:

[1] Jason Walonoski et al. "Synthea: An approach, method, and software mechanism for generating synthetic patients and the synthetic electronic health care record." Journal of the American Medical Informatics Association, 25(3):230–238, 2018. [2] Microsoft. Cromwell on Azure. https://github.com/microsoft/CromwellOnAzure. [3] Microsoft. FHIR to Synapse Sync Agent. https://github.com/microsoft/FHIR-Analytics-Pipelines.

Acknowledgements: The authors would like to thank Venkat Malladi and Olesya Melnichenko for their support throughout this project.

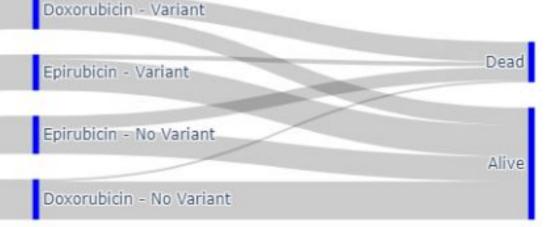
Funding: Microsoft's Biomedical Platforms and Genomics Team supported this work by hosting Tim Dunn through the 2022 Microsoft Research Intern Program. Tim Dunn was also partially supported by the NSF under NSF Graduate Research Fellowship 1841052.

Github Code: https://github.com/microsoft/genomicsnotebook





Fig. 9: Final data frame of merged clinical FHIR data and PacBio variant calling data.







⁴⁸²¹⁻bb84-dd72294c9a4c 07-13T19:23:23.981+00:00'