

## BII SCIENTIFIC CONFERENCE

**Ming Zhen TAN**

Asst PI, BII

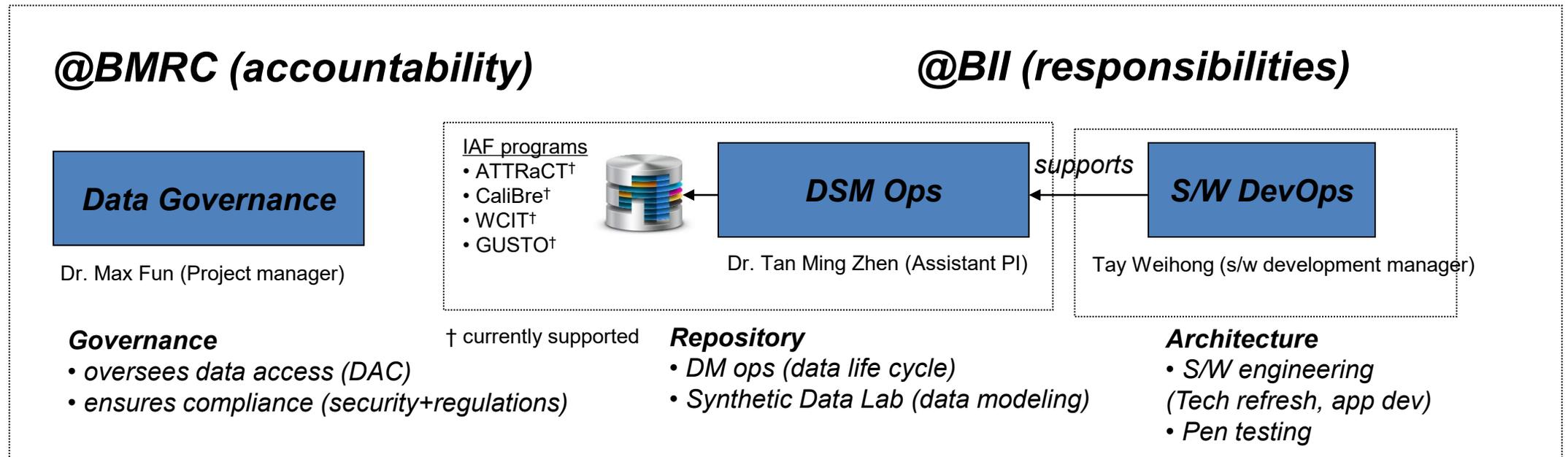
Biomed DAR

29 March 2022

# BioMed DAR: Structure and Operatives

## Mission statement :

To operationalize a SSSO (Standard Systems Support Office) for clinical research data management to support strategic (past, current, future) A\*STAR BMRC programmes and beyond to health clusters



DAC : Data Access Committee  
 DSM : Data Science Management  
 ML : machine learning  
 NLP : Natural Language Processing  
 S/W DevOps: Software Development & Operations



### Accountability

**Introduction** In ethics and governance, accountability is answerability, blameworthiness, liability, and the expectation of account-giving.

**Explanation owed** Yes

### Responsibility

Responsibility may refer to: being in charge, being the owner of a task or event.

Not necessarily

# Synthetic Data Lab



# Extant Clinical Data Sharing and Protection Framework

## Sources of Auxiliary Information

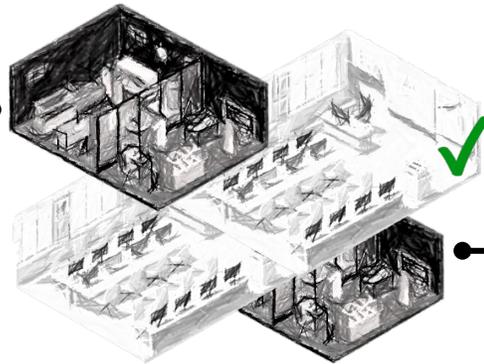
A combination of auxiliary knowledge sourced from social networks and other sources augments the tools at the adversary's disposal. Recovery of pertinent, sensitive information of the data owner from sanitized / deidentified data is already a distinct possibility.



	A	B	C	D	E	F	G	H	I	J	K
	Avg	Z-	SD	Avg	Z-	SD	Avg	Z-	SD	Avg	Z-
	Value	Value	(%)	Value	Value	(%)	Value	Value	(%)	Value	Value
Sample 1	1.002	0.134	4.260	0.300	0.676	22.371	1.000	-1.130	0.446		
Sample 2	1.001	0.200	3.000	0.322	0.694	21.254	1.700	-1.200	0.201		
Sample 3	1.000	0.000	4.574	0.426	0.933	23.640	0.522	-0.294	0.925		
Sample 4	1.005	0.070	4.574	0.426	0.933	23.640	0.522	-0.294	0.925		
Sample 5	1.000	0.000	3.000	0.322	0.694	21.254	1.000	0.000	0.000		
Sample 6	1.000	0.000	3.000	0.322	0.694	21.254	1.000	0.000	0.000		
Sample 7	1.002	0.102	3.754	0.324	0.712	21.571	1.004	-1.200	0.000		
Sample 8	1.007	0.122	3.941	0.424	0.900	22.100	1.000	-1.200	0.452		
Sample 9	1.001	0.100	4.200	0.400	0.700	21.254	1.000	-1.100	0.000		
Sample 10	1.003	0.200	3.000	0.300	0.600	20.710	1.002	-1.000	0.000		
Sample 11	1.007	0.134	4.260	0.300	0.676	22.371	1.000	-1.130	0.446		
Sample 12	1.000	0.200	4.200	0.300	0.600	21.000	0.000	-0.000	0.000		
Sample 13	1.004	-1.303	3.264	0.304	0.600	20.700	0.000	-0.300	0.000		
Sample 14	1.004	-0.000	4.000	0.401	0.712	22.007	0.007	-1.201	0.450		
Sample 15	1.000	0.017	3.754	0.400	0.800	22.554	0.004	0.000	0.000		
Sample 16	1.022	0.094	4.000	0.402	0.620	21.452	1.001	0.000	0.000		
Sample 17	1.070	0.026	4.200	0.426	0.627	22.227	0.000	0.000	0.000		
Sample 18	1.000	0.000	4.520	0.420	0.500	21.000	0.700	-0.000	0.000		
Sample 19	1.000	0.200	4.527	0.300	0.600	21.201	0.001	0.001	0.000		
Sample 20	1.000	0.200	4.527	0.322	0.600	21.011	0.002	0.002	0.002		

## Accumulation of Data

Methodical accumulation of a huge trove of clinical data heralds the advent of data-centric scientific research in the biosciences.

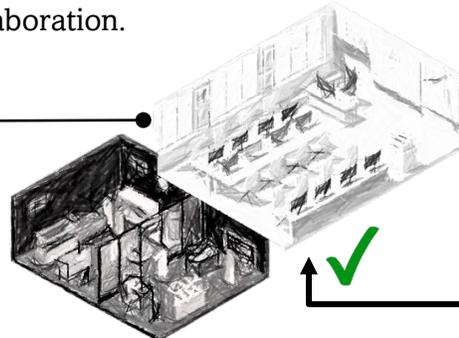


## Wider Collaborative Efforts

Sanitised/Deidentified databases released to the wider community for industrial and academic research. Loose controls over data usage and movement facilitates collaboration.

## Authorised Entities

Authorised Individuals and Entities use data responsibly. Data movements monitored and tracked.



## Cryptographic Protections

Encrypted/Sanitised data cannot be directly reinstated to their original form and can be released for wider collaborative efforts. Improving sanitisation requirements and cryptographic methods keep prying eyes at bay.



	A	B	C	D	E	F	G	H	I	J	K
	Avg	Z-	SD	Avg	Z-	SD	Avg	Z-	SD	Avg	Z-
	Value	Value	(%)	Value	Value	(%)	Value	Value	(%)	Value	Value
Sample 1	1.002	0.134	4.260	0.300	0.676	22.371	1.000	-1.130	0.446		
Sample 2	1.001	0.200	3.000	0.322	0.694	21.254	1.700	-1.200	0.201		
Sample 3	1.000	0.000	4.574	0.426	0.933	23.640	0.522	-0.294	0.925		
Sample 4	1.005	0.070	4.574	0.426	0.933	23.640	0.522	-0.294	0.925		
Sample 5	1.000	0.000	3.000	0.322	0.694	21.254	1.000	0.000	0.000		
Sample 6	1.000	0.000	3.000	0.322	0.694	21.254	1.000	0.000	0.000		
Sample 7	1.002	0.102	3.754	0.324	0.712	21.571	1.004	-1.200	0.000		
Sample 8	1.007	0.122	3.941	0.424	0.900	22.100	1.000	-1.200	0.452		
Sample 9	1.001	0.100	4.200	0.400	0.700	21.254	1.000	-1.100	0.000		
Sample 10	1.003	0.200	3.000	0.300	0.600	20.710	1.002	-1.000	0.000		
Sample 11	1.007	0.134	4.260	0.300	0.676	22.371	1.000	-1.130	0.446		
Sample 12	1.000	0.200	4.200	0.300	0.600	21.000	0.000	-0.000	0.000		
Sample 13	1.004	-1.303	3.264	0.304	0.600	20.700	0.000	-0.300	0.000		
Sample 14	1.004	-0.000	4.000	0.401	0.712	22.007	0.007	-1.201	0.450		
Sample 15	1.000	0.017	3.754	0.400	0.800	22.554	0.004	0.000	0.000		
Sample 16	1.022	0.094	4.000	0.402	0.620	21.452	1.001	0.000	0.000		
Sample 17	1.070	0.026	4.200	0.426	0.627	22.227	0.000	0.000	0.000		
Sample 18	1.000	0.000	4.520	0.420	0.500	21.000	0.700	-0.000	0.000		
Sample 19	1.000	0.200	4.527	0.300	0.600	21.201	0.001	0.001	0.000		
Sample 20	1.000	0.200	4.527	0.322	0.600	21.011	0.002	0.002	0.002		

## Security Frameworks

Security Frameworks and agreements limit data access to authorised individuals and entities. Prohibitive data sharing and usage policies frustrates attempts at data leakage.





# Data Sharing and Protection Framework with SYNTHETIC DATA

## Wider Collaborative Efforts

Collaboration with the wider community can be done at minimal risk of privacy loss.

## Reduced Red Tape

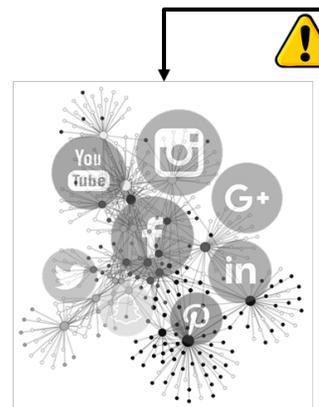
Reduced efforts to adhere to personal data protection requirements because data did not come from any person.

## No Connection to Data Source

Generated synthetic data has no conceivable personal links to the individual from which the raw data was sampled. Auxiliary information contributes little towards privacy infringement.

### Sources of Auxiliary Information

A combination of auxiliary knowledge sourced from social networks and other sources augments the tools at the adversary's disposal. Recovery of pertinent, sensitive information of the data owner from sanitized / deidentified data is already a distinct possibility.

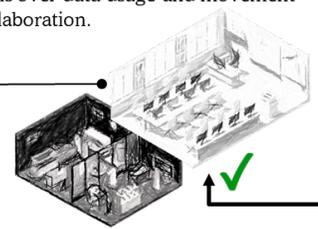


### Wider Collaborative Efforts

Sanitised/Deidentified databases released to the wider community for industrial and academic research. Loose controls over data usage and movement facilitates collaboration.

### Authorised Entities

Authorised Individuals and Entities use data responsibly. Data movements monitored and tracked.



Sample	Age	Sex	Height	Weight	Temp	HeartRate	BloodPressure	Glucose	Cholesterol
Sample 1	1.002	0.134	1.200	66.300	36.875	22.375	1.000	-1.120	0.400
Sample 2	1.001	-0.200	1.200	66.300	36.875	22.375	1.000	-1.120	0.400
Sample 3	1.000	0.000	1.200	66.300	36.875	22.375	1.000	-1.120	0.400
Sample 4	1.005	-0.075	1.200	66.300	36.875	22.375	1.000	-1.120	0.400
Sample 5	1.000	-0.000	1.200	66.300	36.875	22.375	1.000	-1.120	0.400
Sample 6	1.000	0.000	1.200	66.300	36.875	22.375	1.000	-1.120	0.400
Sample 7	1.002	-0.002	1.200	66.300	36.875	22.375	1.000	-1.120	0.400
Sample 8	1.002	0.132	1.201	66.300	36.875	22.375	1.000	-1.120	0.400
Sample 9	1.001	-0.100	1.200	66.300	36.875	22.375	1.000	-1.120	0.400
Sample 10	1.002	0.000	1.200	66.300	36.875	22.375	1.000	-1.120	0.400
Sample 11	1.007	0.134	1.200	66.300	36.875	22.375	1.000	-1.120	0.400
Sample 12	1.009	-0.210	1.200	66.300	36.875	22.375	1.000	-1.120	0.400
Sample 13	1.004	-2.303	1.200	66.300	36.875	22.375	1.000	-1.120	0.400
Sample 14	1.004	-0.000	1.200	66.300	36.875	22.375	1.000	-1.120	0.400
Sample 15	1.000	-0.017	1.200	66.300	36.875	22.375	1.000	-1.120	0.400
Sample 16	1.002	0.000	1.200	66.300	36.875	22.375	1.000	-1.120	0.400
Sample 17	1.000	0.000	1.200	66.300	36.875	22.375	1.000	-1.120	0.400
Sample 18	1.000	0.000	1.200	66.300	36.875	22.375	1.000	-1.120	0.400
Sample 19	1.000	-0.320	1.200	66.300	36.875	22.375	1.000	-1.120	0.400
Sample 20	1.000	0.000	1.200	66.300	36.875	22.375	1.000	-1.120	0.400

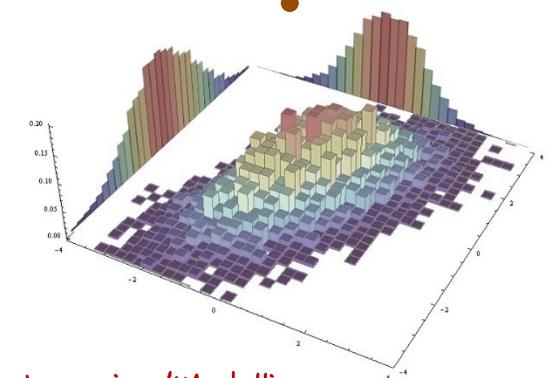
Sample	Age	Sex	Height	Weight	Temp	HeartRate	BloodPressure	Glucose	Cholesterol
Sample 1	1.002	0.134	1.200	66.300	36.875	22.375	1.000	-1.120	0.400
Sample 2	1.001	-0.200	1.200	66.300	36.875	22.375	1.000	-1.120	0.400
Sample 3	1.000	0.000	1.200	66.300	36.875	22.375	1.000	-1.120	0.400
Sample 4	1.005	-0.075	1.200	66.300	36.875	22.375	1.000	-1.120	0.400
Sample 5	1.000	-0.000	1.200	66.300	36.875	22.375	1.000	-1.120	0.400
Sample 6	1.000	0.000	1.200	66.300	36.875	22.375	1.000	-1.120	0.400
Sample 7	1.002	-0.002	1.200	66.300	36.875	22.375	1.000	-1.120	0.400
Sample 8	1.002	0.132	1.201	66.300	36.875	22.375	1.000	-1.120	0.400
Sample 9	1.001	-0.100	1.200	66.300	36.875	22.375	1.000	-1.120	0.400
Sample 10	1.002	0.000	1.200	66.300	36.875	22.375	1.000	-1.120	0.400
Sample 11	1.007	0.134	1.200	66.300	36.875	22.375	1.000	-1.120	0.400
Sample 12	1.009	-0.210	1.200	66.300	36.875	22.375	1.000	-1.120	0.400
Sample 13	1.004	-2.303	1.200	66.300	36.875	22.375	1.000	-1.120	0.400
Sample 14	1.004	-0.000	1.200	66.300	36.875	22.375	1.000	-1.120	0.400
Sample 15	1.000	-0.017	1.200	66.300	36.875	22.375	1.000	-1.120	0.400
Sample 16	1.002	0.000	1.200	66.300	36.875	22.375	1.000	-1.120	0.400
Sample 17	1.000	0.000	1.200	66.300	36.875	22.375	1.000	-1.120	0.400
Sample 18	1.000	0.000	1.200	66.300	36.875	22.375	1.000	-1.120	0.400
Sample 19	1.000	-0.320	1.200	66.300	36.875	22.375	1.000	-1.120	0.400
Sample 20	1.000	0.000	1.200	66.300	36.875	22.375	1.000	-1.120	0.400

## Synthetic Data

Synthetic data sampled from the learned probability distributions. Process is inexpensive as compared to traditional clinical data collection and a new set can be generated for different users or uses.

## Major Utility

Synthetic data shares the same probability distributions with the raw data and can be used as a substitute for algorithm development, etc.



### Accumulation of Data

Methodical accumulation of a huge trove of clinical data heralds the advent of data-centric scientific research in the biosciences.

### Cryptographic Protections

Encrypted/Sanitised data cannot be directly reinstated to their original form and can be released for wider collaborative efforts. Improving sanitisation requirements and cryptographic methods keep prying eyes at bay.

## Learning/Modelling

The underlying data distributions (univariate & joint) are learned/modelled from the raw data. This forms the basis from which new data could be synthesised.

## Major Cost Savings

Synthetic data does not require encryption-at-rest protocols, secure database warehousing, auditing, cryptographic cleansing, approvals, collaboration agreements, etc.

### Security Frameworks

Security Frameworks and agreements limit data access to authorised individuals and entities. Prohibitive data sharing and usage policies frustrates attempts at data leakage.



# Possible Applications / Collaborations for Synthetic Data Lab

- **Types of Synthetic Data**

- Tabular Clinical Data
- Medical Images (CT scans / MRI)
- Time Series (fMRI signals, ECG signals)
- Medical Records

- **Applications of Synthetic Data**

- Training Datasets to accelerate development of algorithms
- Independent Test Datasets to validate algorithms
- Mock data generation for systems development
- Mock data generation for improved statistical power
- Combination of multiple data sources to create joint probability distribution

## Synthetic Data Lab Real Working Example





# REAL TABULAR DATA EXAMPLE: ATTRACT DATA SET

## Number of Variables

315

## Number of Data Points

3964

## Longitudinal Data

### Collection (#Variables)

Baseline: 85  
Baseline/6 months: 45  
Baseline/6 weeks/6 months: 88  
Baseline/6W/6M/1Y/2Y/3Y: 94  
Reference: 3

## Longitudinal Data

### Collection (#Data points)

Baseline: 912  
6 weeks: 428  
6 months: 420  
1Y: 816  
2Y: 787  
3Y: 601

## Type of Variables

Date: 41  
Mixture (numeric): 24  
Nominal: 107  
Numeric: 80  
Ordinal: 10  
Percentage: 1  
Range: 1  
String: 5  
Text: 44  
Unknown: 2

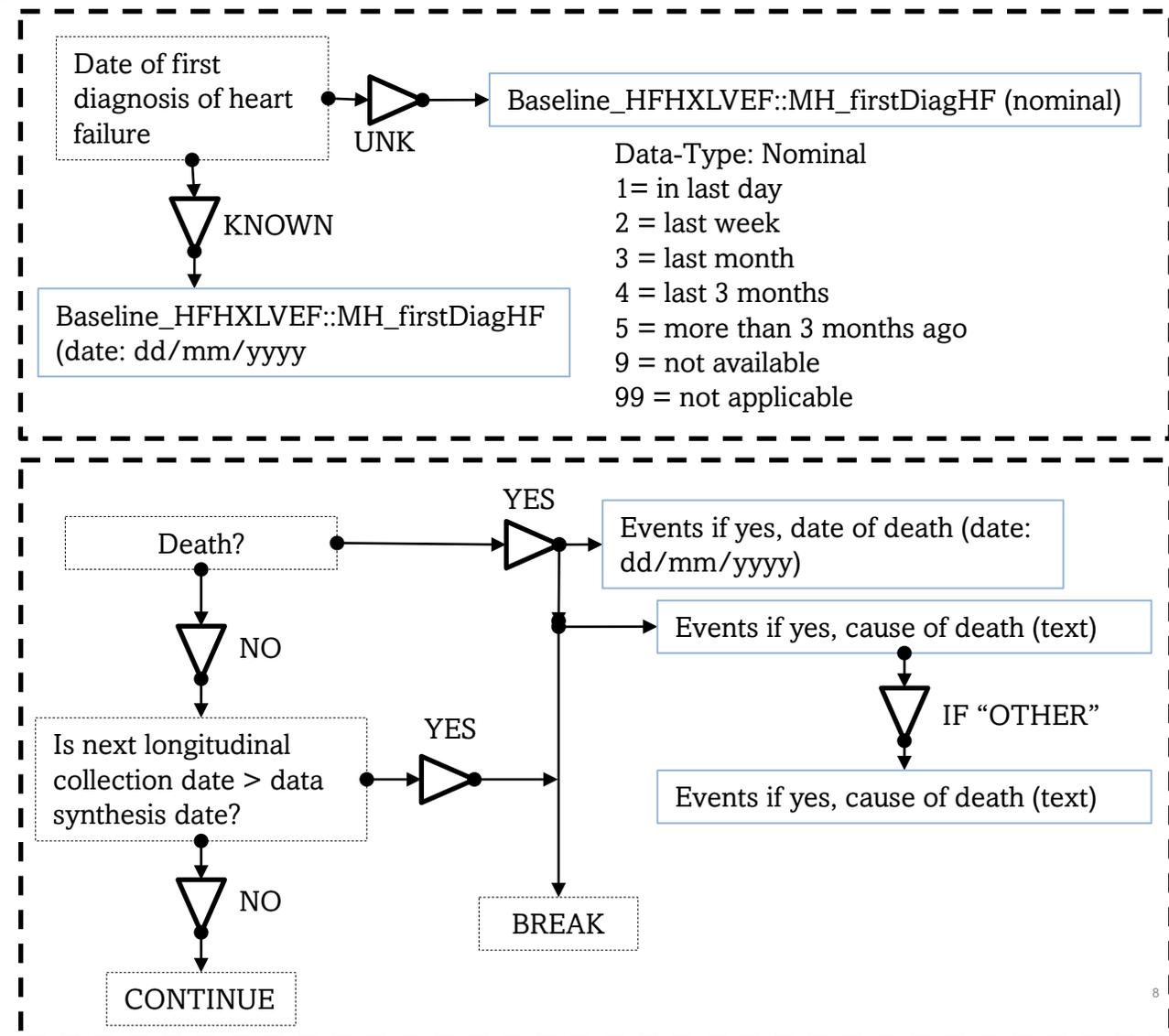
## Missing Datapoints

Yes

## Date Dependencies

Yes

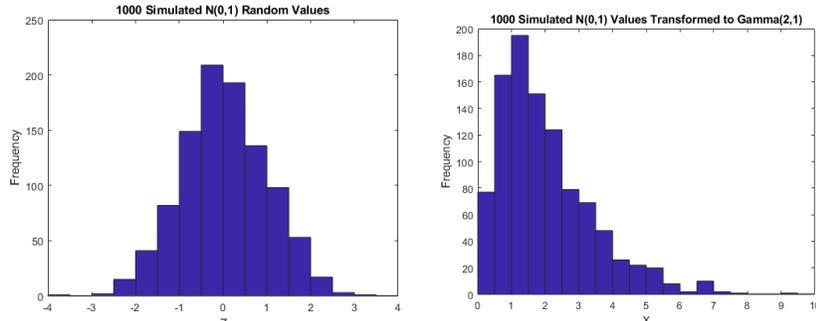
## Tree Dependencies



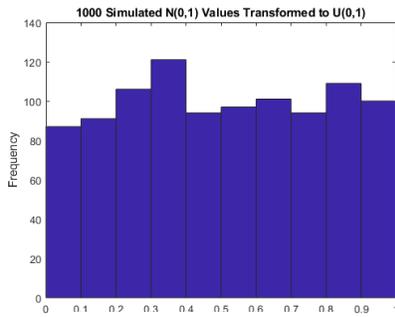


# Marginal Probabilities

The various multivariate probability distributions are a generalisation of univariate probability distributions to one or more variables.



The probability integral transform of a random variable is a random variable that is uniform on [0,1]



$$F_i(x) = \Pr[X_i \leq x]$$

Applying the inverse probability integral transform of a distribution F to a random variable whose distribution is F.

$$(X_1, X_2, \dots, X_d)$$

$$(U_1, U_2, \dots, U_d) = (F_1(X_1), F_2(X_2), \dots, F_d(X_d))$$

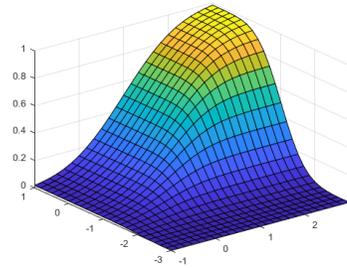
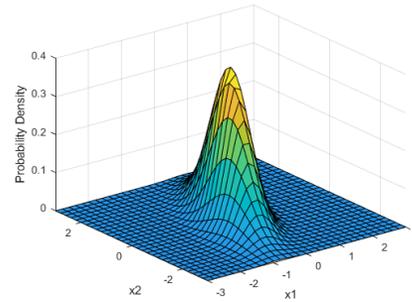
$$(X_1, X_2, \dots, X_d) = (F_1^{-1}(U_1), F_2^{-1}(U_2), \dots, F_d^{-1}(U_d)).$$

# Joint Probabilities

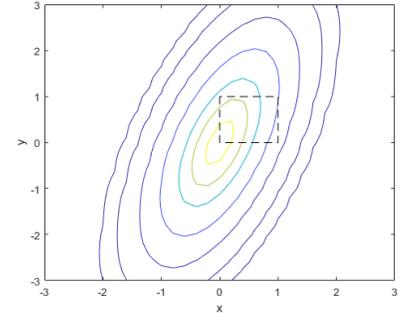
The various multivariate probability distributions are a generalisation of univariate probability distributions to one or more variables.

Closed form of PDF of d-dim multivariate normal dist.

$$y = f(x, \mu, \Sigma) = \frac{1}{\sqrt{|\Sigma|(2\pi)^d}} \exp\left(-\frac{1}{2}(x-\mu) \Sigma^{-1}(x-\mu)'\right)$$



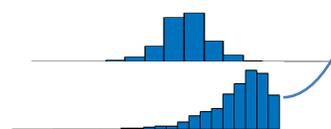
If we have the closed form of a multivariate PDF/CDF, we can generate sets of values from the multivariate PDF.



But we don't have closed forms for varying marginal distributions with varying correlations.

# Copula

$$(X_1, X_2, \dots, X_d)$$

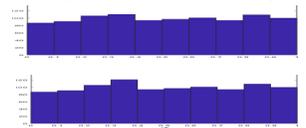


Compute data correlations from training data

Use correlations in some known form of multivariate probability distribution for uniform RV on [0,1].

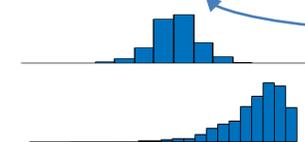
Sample from the multivariate probability distribution of random variables.

$$(U_1, U_2, \dots, U_d)$$



Do inverse transform on sampled variables, based on what I think their marginal distribution is.

$$\text{Synthetic } (X_1, X_2, \dots, X_d)$$

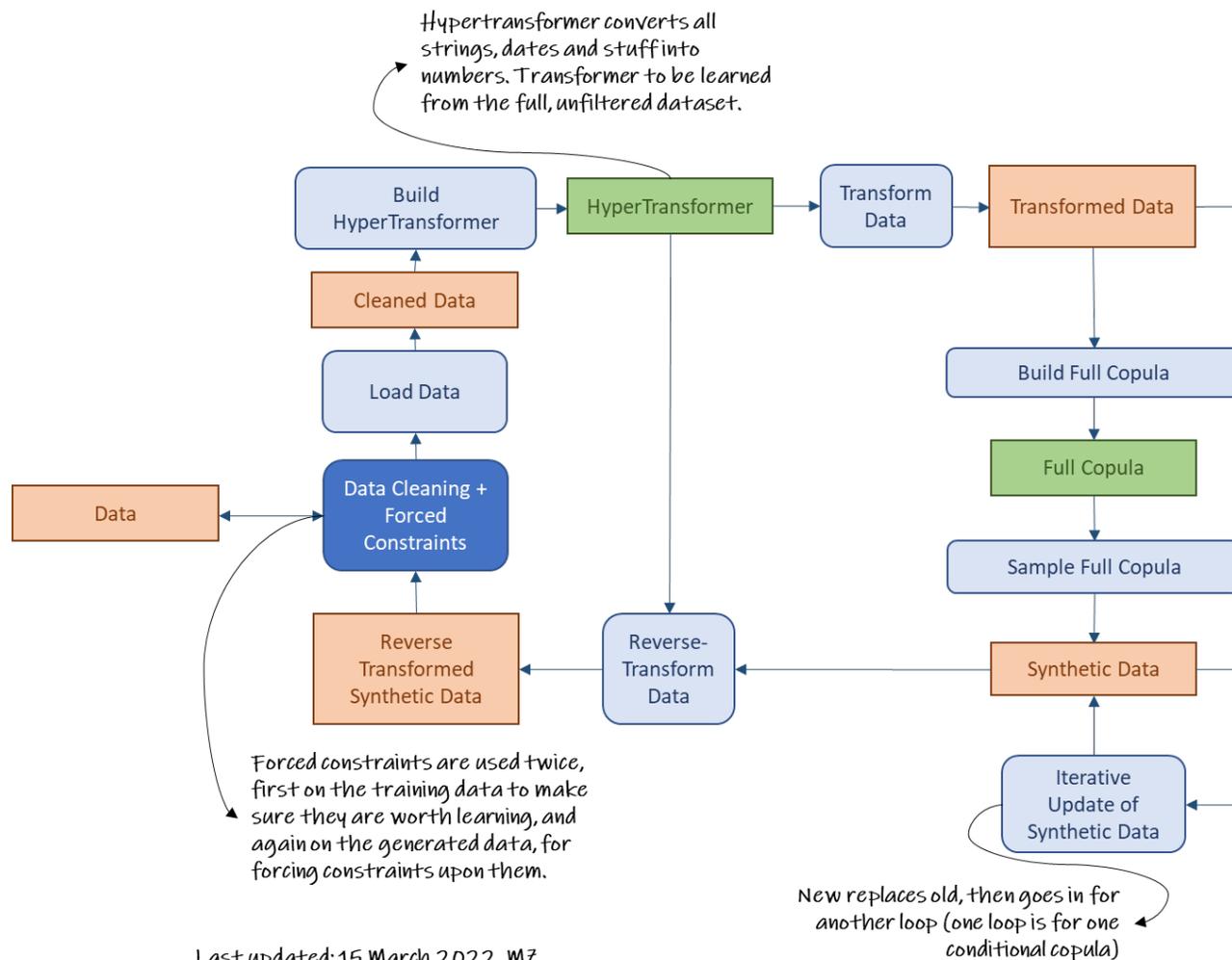


$$(F_1^{-1}(U_1), F_2^{-1}(U_2), \dots, F_d^{-1}(U_d))$$

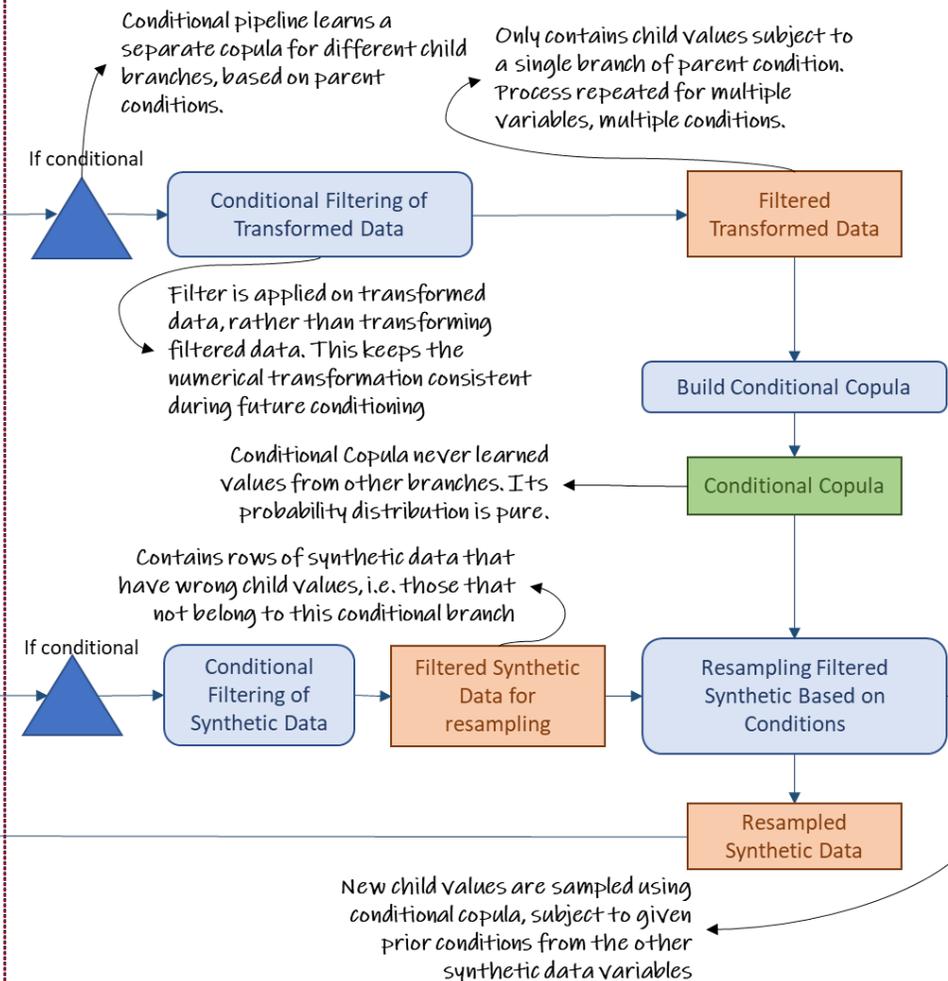
**COPULA:** dependency structure between marginals.  
 $\Pr[U_1 \leq u_1, U_2 \leq u_2, \dots, U_d \leq u_d].$   
 A CDF based on uniform random variables on [0,1]

# MODELLING JOINT DEPENDENCIES

## COPULA LEARNING



## CONDITIONAL-COPULA LEARNING



# QUICK GLANCE

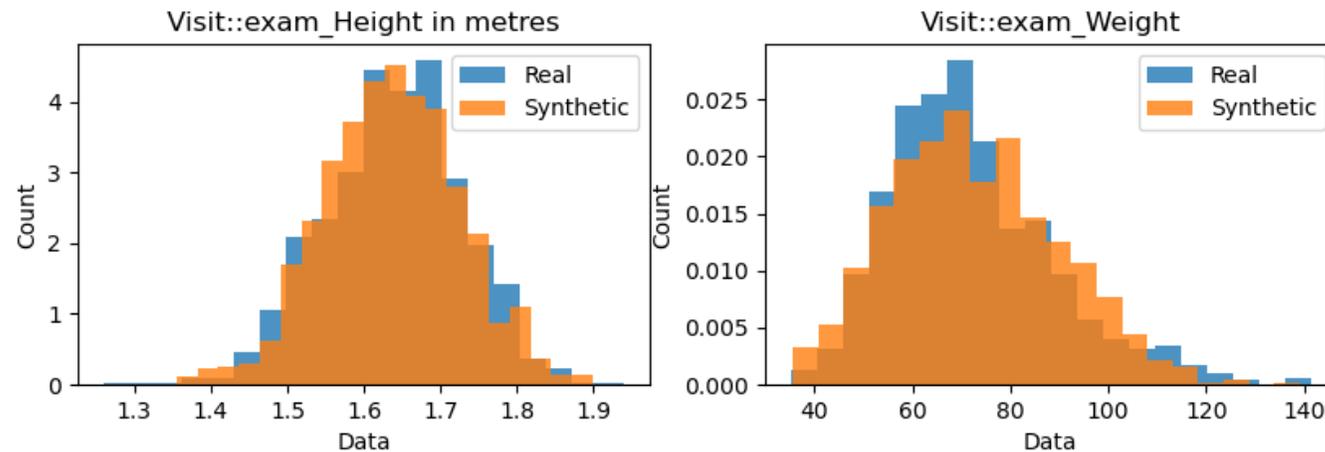
## Experiment

Multivariate copula learning:

- 912 baseline datapoints;
- 315 variables;

## Qualitative Look at Histograms

Qualitative inspection of histogram plots of individual continuous variables. Multivariate copula can model univariate continuous distributions to reasonable accuracy.



## Qualitative Look at Regressions (Linear)

Simple linear regression using height and weight variables from both original and synthetic data. Figure shows superimposed scatterplots and regression trendlines from both the original and synthetic datasets. There is visible high overlap of the two scatterplots, though slope and intercept confidence levels are slightly different.

Original slope (95%): 93.250696 +/- 10.649427  
 Original slope (95%): 96.233674 +/- 10.403114  
 Original intercept (95%): -80.618723 +/- 17.465537  
 Original intercept (95%): -84.849367 +/- 17.052326





CREATING GROWTH, ENHANCING LIVES



**THANK YOU**

---

[www.a-star.edu.sg](http://www.a-star.edu.sg)