



Revolutionizing Healthcare with Synthetic Clinical Trial Data

Afrah Shafquat

Senior Data Scientist II, Medidata AI

May 10, 2023



Synthetic Data Applications



Data sharing

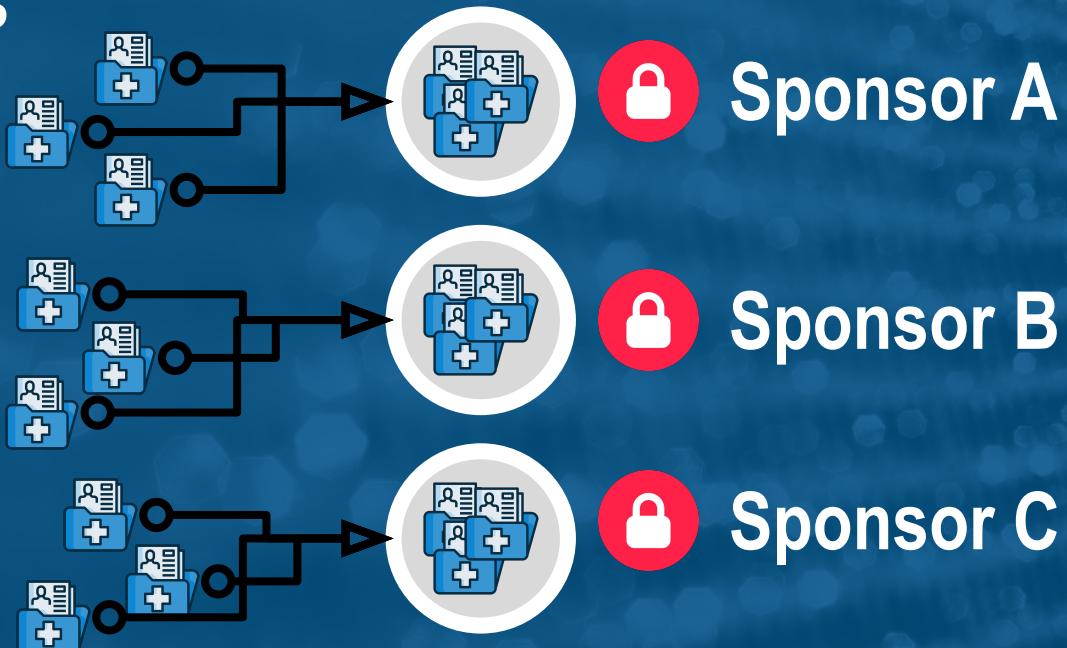


Data augmentation

Data Sharing



Clinical trial data remains siloed due to patient privacy and sponsor privacy concerns



Icons Source: Flaticon; Full credits in the last slide

Collaboration remains limited due to data access restrictions



Sponsor A



Sponsor B



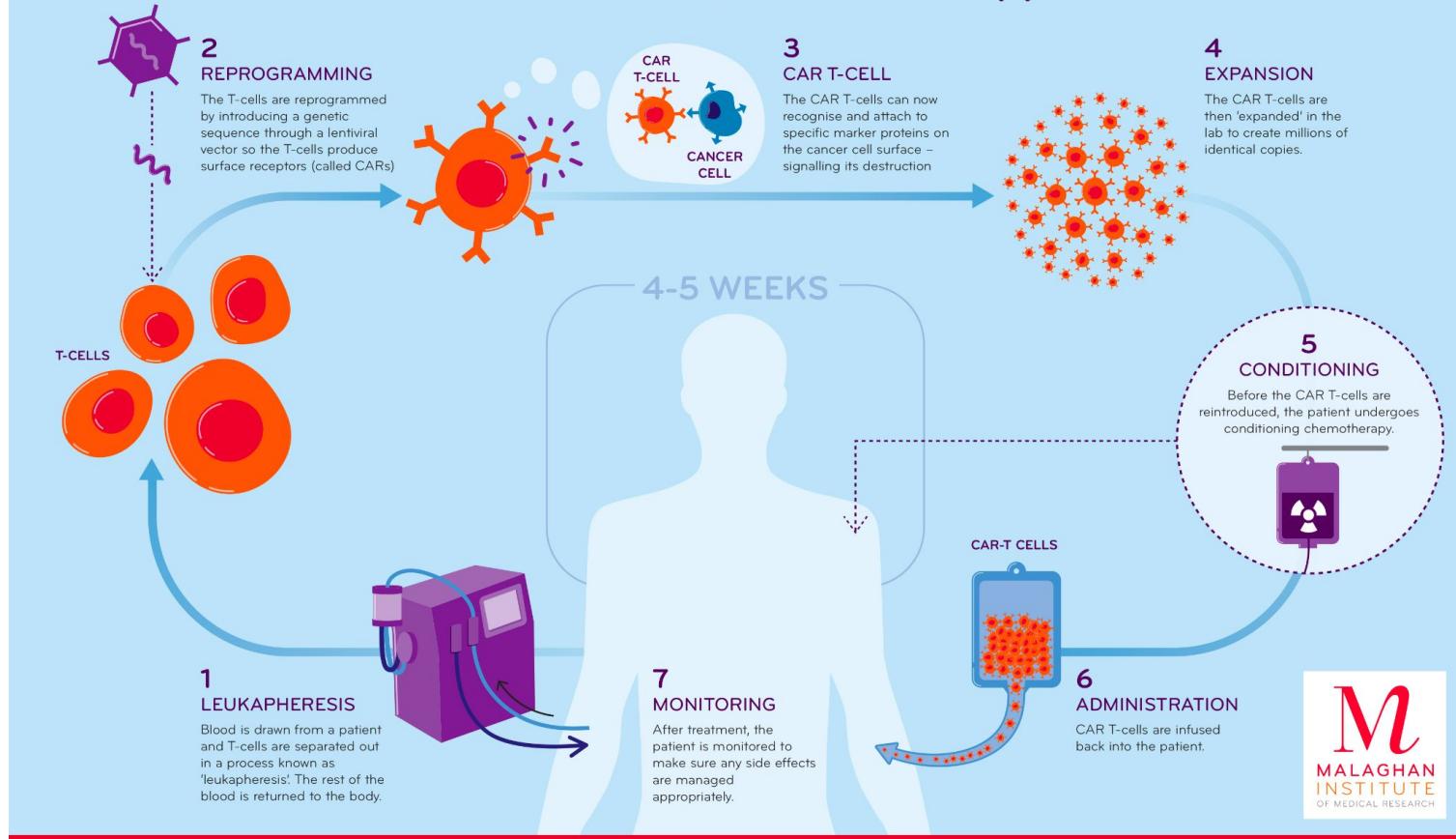
Sponsor C



Why is collaboration necessary?



CAR T-Cell Cancer Therapy



TIME

HEALTH • CANCER

Cancer's Newest Miracle Cure

August 21, 2017 issue of TIME.



“Went looking for a miracle...”

Source

CHOP News (May 2022)



Cell therapy almost always begins its development journey by offering a treatment option to patients with very advanced & difficult to treat cancers.



No one can know
exactly how any one
patient will respond to
these experimental
treatments ...



... but we do know that making a personalized medicine for an individual relies on our experience treating similar people with similar diseases



1 patient



1 doctor

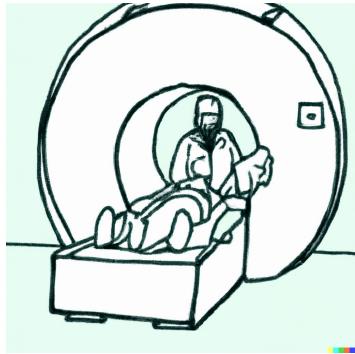


1 trial 1 sponsor **Medidata AI**

If we could learn from the combined experience of the whole world, rather than just relying on the experience of a single doctor or hospital or sponsor ...

How well could we do?





Relapse

The journey of a
CAR-T trial participant
begins with the return
of cancer ...



Relapse



Enrollment

... the failed search
for a viable approved
treatment, and an
introduction to a
clinical investigator ...



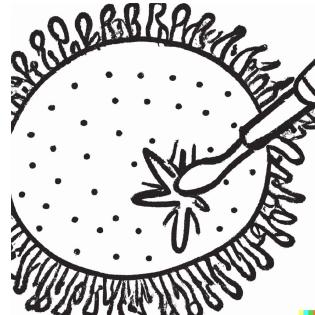
Relapse



Enrollment



Apheresis



Manufacturing

... removal of healthy lymphocytes and manufacture of a targeted dose of CAR-T cells ...



Relapse



Enrollment



Apheresis



Manufacturing



Treatment



Management

... treatment with
investigational CAR-T
therapy and careful
management of patients
through the side effects
that follow ...



... finally, a new
response to treatment,
disease control ...



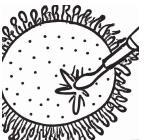
Relapse



Enrollment



Apheresis



Manufacturing



Treatment



Management



Response



Control



Home.

... and a return Home.



Relapse
Enrollment
Apheresis
Manufacturing
Treatment
Management
Response
Control

Home

Now, what if we
understood this journey
for 100 patients?



For 3,000?

1 patient. 1 journey. 1 experience.



Medidata AI — 3,000 CAR-T patient journeys today ... and growing



If we could learn from the combined experience of the whole world, rather than just relying on the experience of a single doctor or hospital or sponsor, what?

How well could we do?



* visuals created using DALL-E

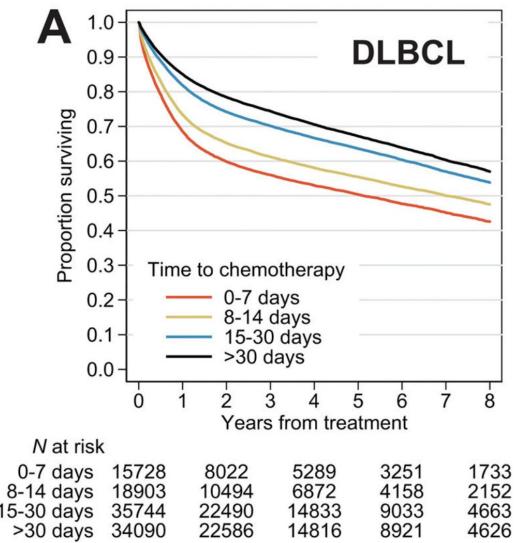


Every step in the patient's
journey provides
opportunities to design
better trials to prove the
benefit of new treatments



Relapse

Example



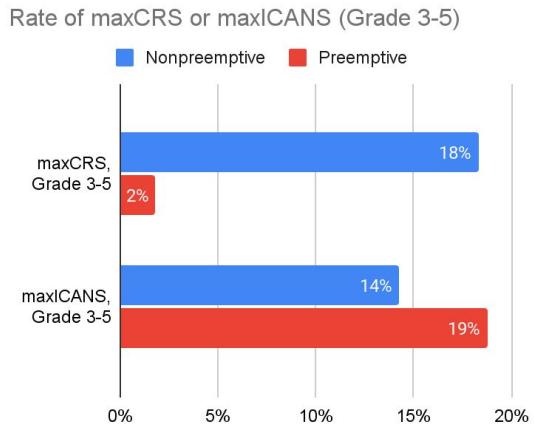
The length of **time between diagnosis and treatment** may be the **single strongest predictor of the prognosis** for trial participants with high risk lymphoma

Alshaibani, A. et al (2019). *High risk patients with diffuse large B cell lymphoma are not enrolled on clinical trials.* JCO 37, e19058–e19058.



Management

Example



Pre-emptive treatment of CAR-T recipients with tocilizumab & dexamethasone reduces severe CRS rates but does not decrease rates of ICANS (neurotoxicity)

EHA 2023

Authors:

Esther Nie, MD, PhD^{2†}; Penelope Lafeuille, MS^{1†}; Sheila Diamond, MS, CGC¹; Jacob Aptekar, MD, PhD¹; Vibhu Agarwal, PhD, MBA¹

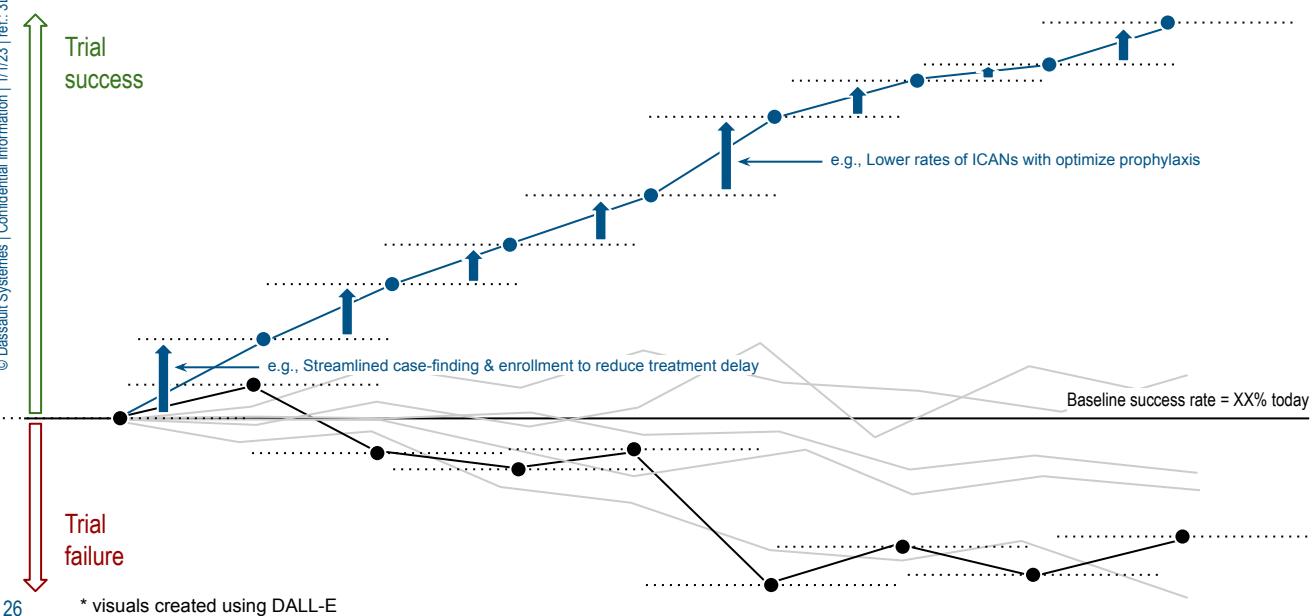
Affiliations:

1. Medidata, a Dassault Systèmes company, New York, NY
2. Stanford University, Stanford, CA



Probability of Technical & Regulatory Success

Schematic



If we designed trials so that we made every decision informed by the experience of treating thousands of similar patients

How well could we do?



Select case examples and scientific publications

CASES

Top 20 BioPharma #1

Safer trials, better outcomes

Identified patients at high risk of developing severe CRS based on pre-infusion characteristics. Recommended optimal dosing regimen to minimize CRS by 20-30%

Top 20 BioPharma #2

Broader trial population, greater unmet need

Recommended eligibility criteria for 1L DLBCL trial to expand label and accelerate enrolment. Reduced trial timeline by 6-12 months

EU Biotech

Data-driven trial design, optimization

Developed optimal treatment strategies to minimize CRS and ICANs. 30% fewer cases of grade 3+ AEs. Identified patient likely to respond to therapy based on prior treatments

Top 20 BioPharma #3

Safer treatments in the clinic=bigger market

Predicted which patients are at higher risk of developing CRS based on prior treatment history and patient characteristics. ~20% higher accuracy in identifying high patients

PUBLICATIONS



[Predictors of severe CRS in longitudinal CAR T-cell clinical trial data](#)
Poster presentation, 2022



[Deriving Predictive Features of Severe CRS from Pre-Infusion Clinical Data in CAR T-Cell Therapies](#)
Poster presentation, 2022

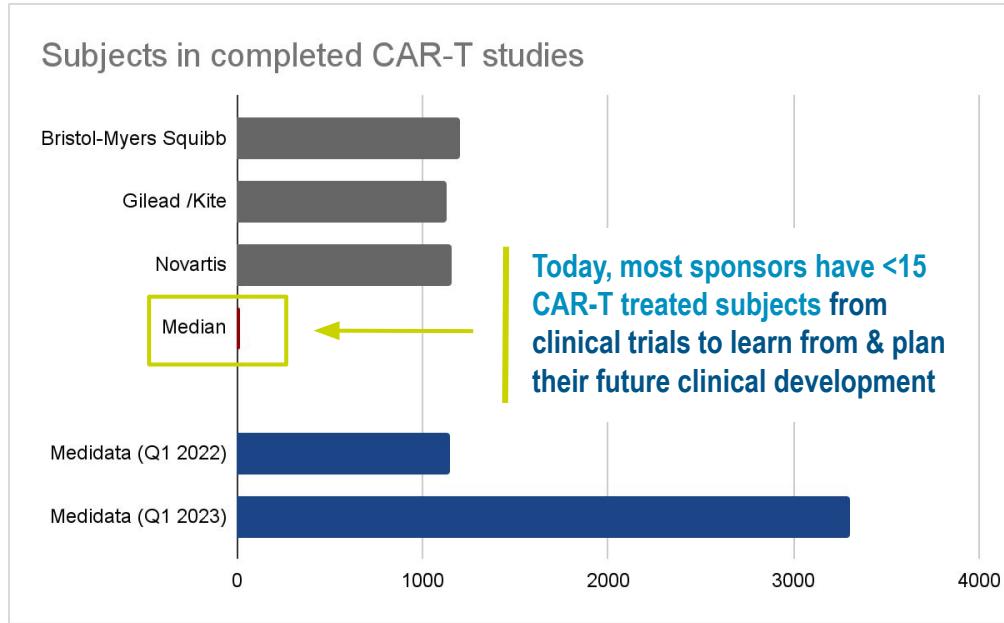


[Evaluating Early Risk of Cytokine Release Syndrome \(CRS\) Induced by CAR T-Cell Therapy Using Pooled Clinical Trial Data](#)
Poster presentation, 2022



Cooccurrence patterns of CRS and ICANS in patients undergoing autologous CD19-targeted CAR T-cell treatments
In submission, 2023

60% of CAR-T data has been generated by 4 sponsors, and most sponsors working in this space have data on <15 subjects



50+ trials

CAR-T or T-Cell Engager (TCE) clinical trials run on Medidata's platform



>5k+ subjects

available to Medidata AI in unique CAR-T and TCE synthetic dataset

Synthetic data supports a healthy, resilient biopharma ecosystem

Forests rely on **nutrient sharing** by fungi between plants (*mycorrhiza*) to support a healthy, local ecosystem

Nutrients move between plants via a network of fungal hyphae (tubes) that digests, mobilizes & delivers water & essential minerals to neighboring plants & trees to enhance growth & disease resistance within a small, geographical area

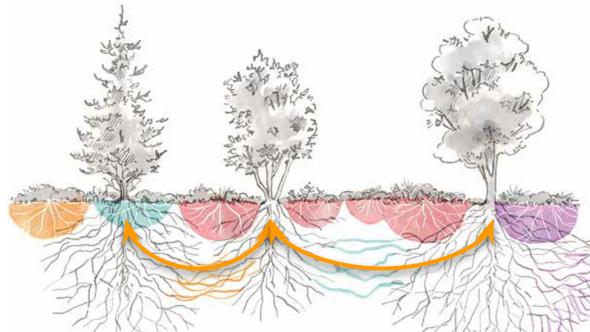
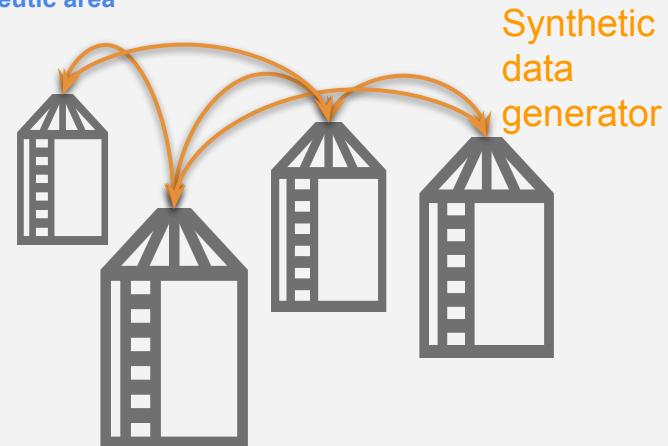


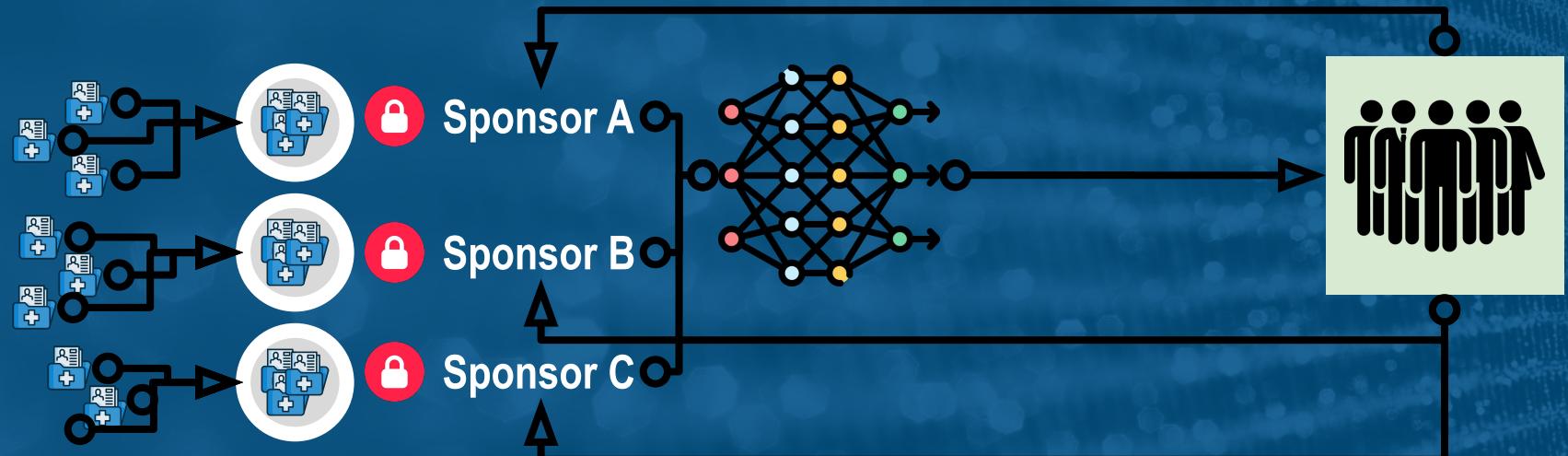
Image from book: <https://www.amazon.com/Mycorrhizal-Planet-Symbiotic-Support-Fertility/dp/160358658X?asin=160358658X&revisionId=4&format=4&depth=1>

Biopharma relies on **knowledge sharing** to support a healthy, productive ecosystem for innovation

Insights move between companies via Synthetic data network that digests, mobilizes & delivers information & insights to PharmaCos & Biotechs to develop safe, effective medicines within a therapeutic area



Synthetic data is a solution to encourage collaboration and continue innovation in pharma



Clinical Trials

Unlearn.AI, a startup developing a 'digital twin' service for clinical trials, raises \$50M

Kyle

COMMENT | June 1, 2022

Can AI-generated prognostic forecasts substitute patients in oncology trials?

Menk's latest bet is to generate clinical evidence using digitally simulated 'predicted outcomes' rather than actual patients.

Unlearn Receives Draft Qualification Opinion from European Medicines Agency for Using the PROCOVATM Framework to Implement TwinRCTs™

The three-step PROCOVATM procedure provides a clear framework for implementing Unlearn's TwinRCT™ solution to accelerate Phase 2 and Phase 3 clinical trials

May 12, 2022 08:00 AM Eastern Daylight Time

Real World Data

The People in This Medical Research Are Fake. The Innovations Are Real.

synthetic-data technology, by creating artificial patient

populations, has the potential to speed up innovations without

compromising privacy

Anthem Looks to Fuel AI Efforts With Petabytes of Synthetic Data

Health insurance company is working with Google Cloud to generate 1.5 to 2 petabytes of synthetic data

aimed at detecting fraud, delivering personalized

service to members

Action acquires synthetic data trailblazer Replica Analytics

Acquisition will give Action customers ability to tap into previously inaccessible, high utility, global health data to conduct transformational healthcare research

NEWS PROVIDED BY
Action Inc.
Action on
Jan 04, 2022, 07:00 ET

Industry at Large

The market for synthetic data is bigger than you think

By 2024, 60% of the data used for the development of AI and analytics projects will be synthetically generated

By Andrew White | July 24, 2021 |

0 Comments | Search all blog posts

Fake It to Make It: Companies Beef Up AI Models With Synthetic Data

American Express experiments with AI-generated fake fraud patterns to sharpen its models' ability to detect rare or uncommon swindles

By Sara Castellanos
July 23, 2021 5:30 am ET WSJ PRO

Synthetic Data landscape

Available open-source solutions



 SmartNoise

 MEDIDATA
Simulants



synthcity

gretel-synthetics by
 gretelTM

SDV

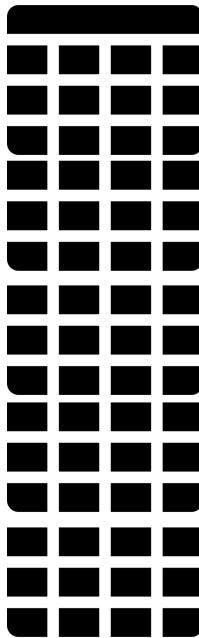
The Synthetic Data Vault

 datacebo

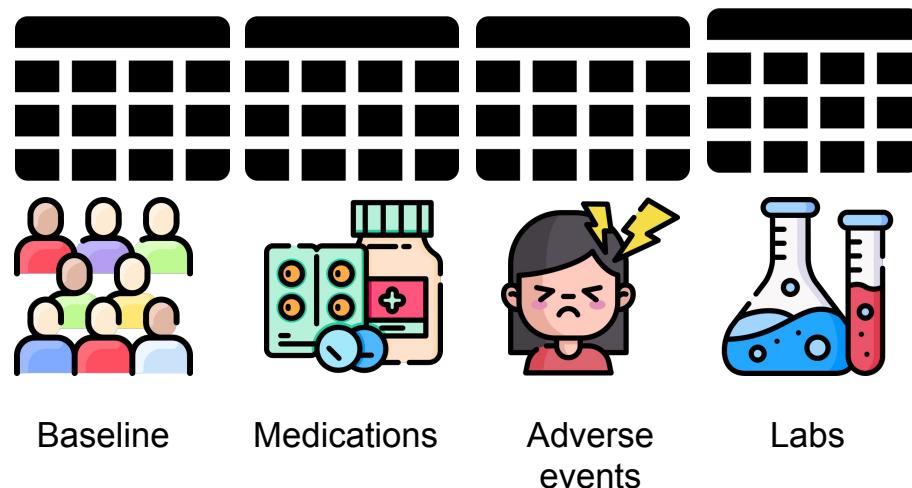
 DS MEDIDATA

 DASSAULT
SYSTEMES

Limitations in current solutions and the uniqueness of clinical trial data

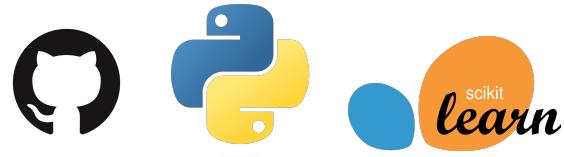


VS.



Icons Source: Flaticon; Full credits in the last slide

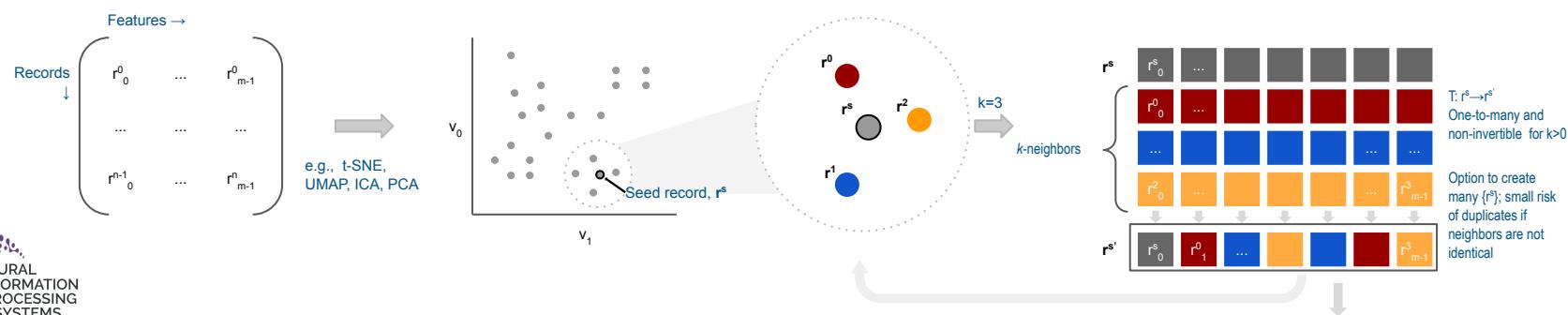
Simulants: Open-source Methodology



Privacy, fidelity & numerosity can be easily tuned and adjusted through choice of embedding, number of neighbors & feature linking

Simplified Process Diagram

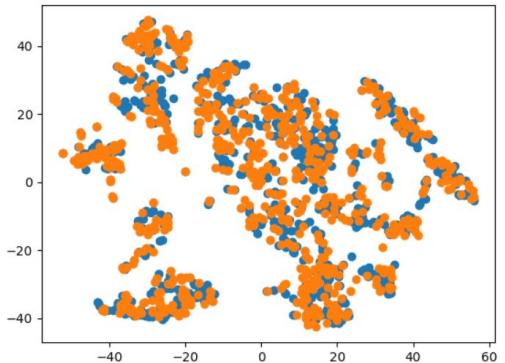
Load source dataset, R → Embed records, r^i , in metric space → Select seed, r^s , and identify nearest neighbors → Generate new record(s), $r^{s'}$, by randomly selecting features from k -neighbors



1. M Beigi, A Shafquat, J Mezey, JW Aptekar [Synthetic Clinical Trial Data while Preserving Subject-Level Privacy](#) - NeurIPS 2022 Workshop on Synthetic Data for Empowering ML Research
2. M Beigi, A Shafquat, J Mezey, JW Aptekar **Simulants: Synthetic Clinical Trial Data via Subject-Level Privacy-Preserving Synthesis** - AMIA 2022
3. Open-source Simulants GitHub: <https://github.com/mdsol/Simulants>
4. Simulants is implemented in Python and uses packages from open-source libraries including scikit-learn and numpy

Qualitative: t-SNE Embedding

t-SNE Embedding & Overlap

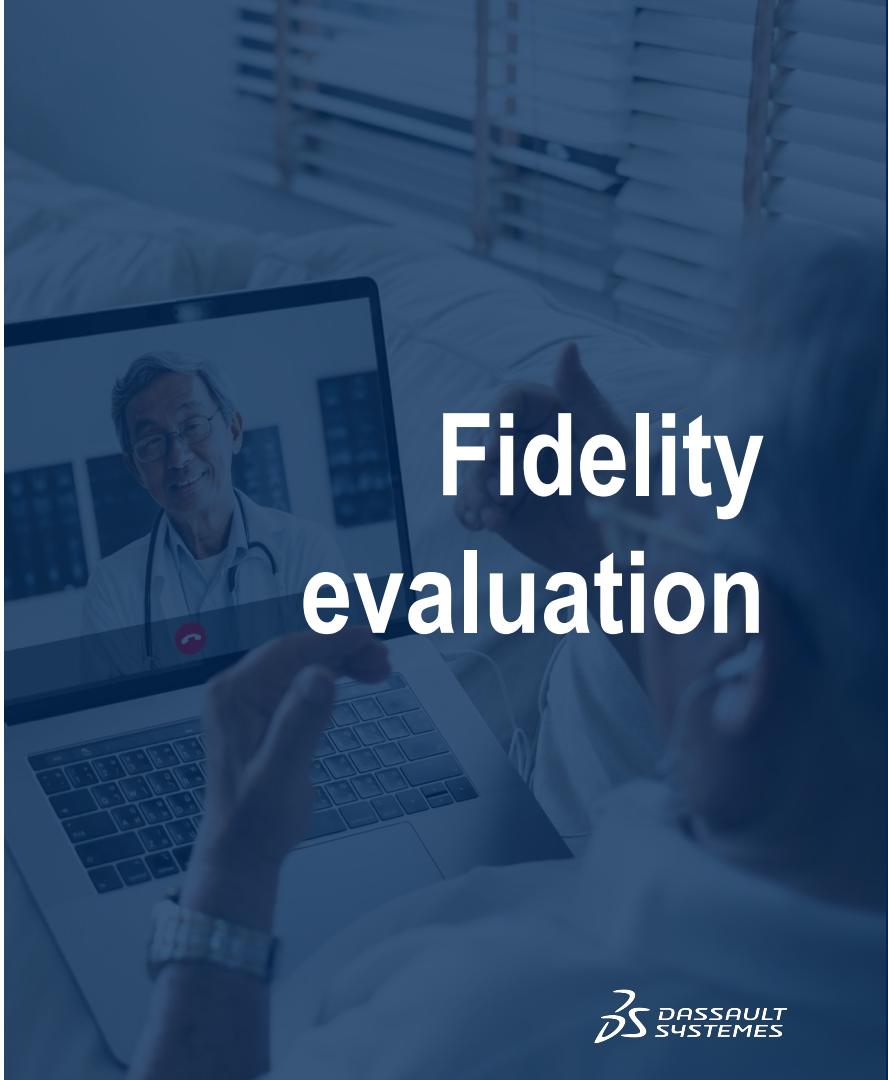


Alignment between source/template and Simulants generated synthetic dataset indicates synthesis reflects the source population well

Other open-source fidelity metrics available at the MIT/Data Cebo SDV project

SDV
The Synthetic Data Vault


1. M Beigi, A Shafquat, J Mezey, JW Aptekar [Synthetic Clinical Trial Data while Preserving Subject-Level Privacy](#) - NeurIPS 2022 Workshop on Synthetic Data for Empowering ML Research
2. M Beigi, A Shafquat, J Mezey, JW Aptekar [Simulants: Synthetic Clinical Trial Data via Subject-Level Privacy-Preserving Synthesis](#) - AMIA 2022



Fidelity evaluation

Quantitative: Univariate measures

Comparison of Mean and other statistical similarity metrics

	Template	Simulants
NUM_CYC	7.030086	7.379656
TTE_SD1	56.193410	54.720630
TTE_DEATH	249.376791	254.849570
TTE_BOR	50.110315	49.786533
HEIGHT	161.403023	160.781633
TTE_PFS	179.895415	187.706304
TTE_GLYC_AE	85.014327	85.451289
PRIOR_RADIOIX_TIME	64.816619	68.999284
AGE	65.110315	64.769341
PRIOR_CANC_SURG_TIME	265.462751	230.846705
SLD_BASELINE	61.173352	60.402579
ECOG	0.638968	0.651862
TTE_PR1	62.312321	61.693410
DEPTH_RESP	-24.353725	-24.694126
TTE_SKIN_AE	24.965616	28.411175
PRIOR_CHEMO_TIME	659.383954	659.811605
TTE_CR1	52.563037	53.174785
TTE_PD1	179.385387	179.308023
TTE_2L	179.388252	182.510029

Comparison of mean across numerical features shows the mean of source/template and Simulants generated synthetic data are close.

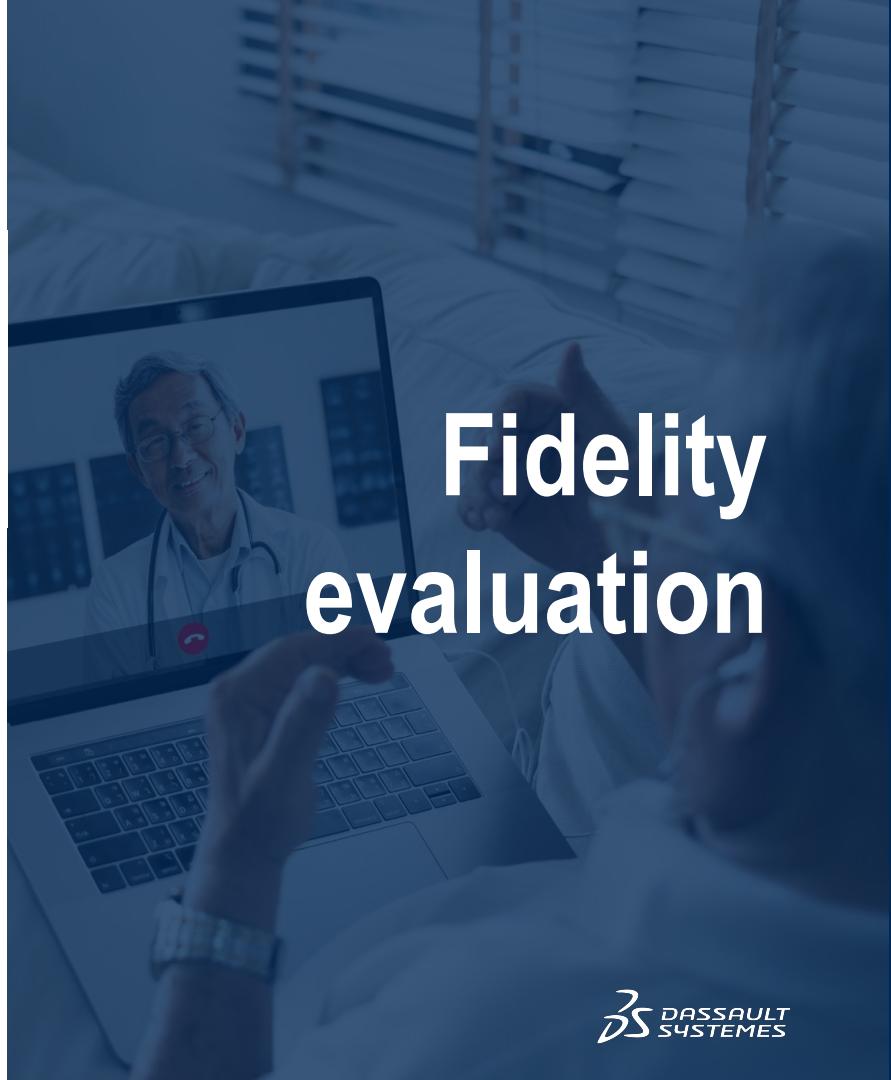
Other open-source fidelity metrics available at the MIT/Data Cebo SDV project

SDV

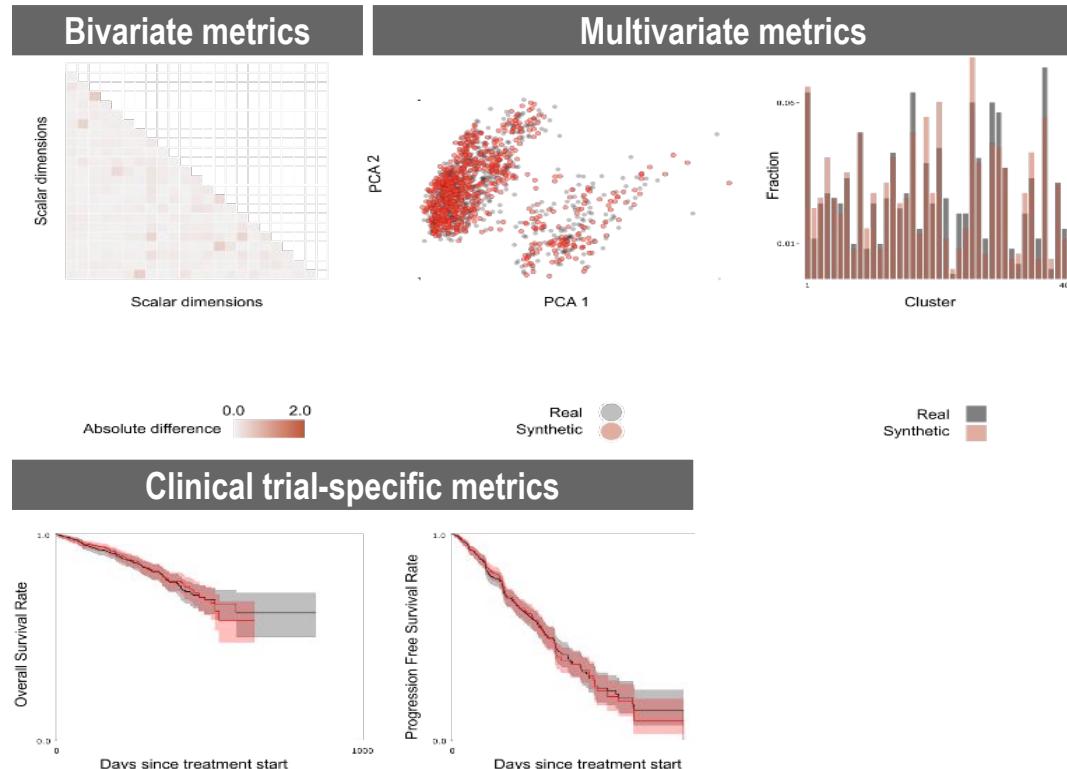
The Synthetic Data Vault



1. M Beigi, A Shafquat, J Mezey, JW Aptekar [Synthetic Clinical Trial Data while Preserving Subject-Level Privacy - NeurIPS 2022 Workshop on Synthetic Data for Empowering ML Research](#)
2. M Beigi, A Shafquat, J Mezey, JW Aptekar [Simulants: Synthetic Clinical Trial Data via Subject-Level Privacy-Preserving Synthesis - AMIA 2022](#)



Quantitative metrics



Bivariate metrics like comparison of pairwise correlations allow assessment of preservation of correlations across features*

* All fidelity metrics were computed using open-source packages including scipy, tsne, scikit-learn, numpy etc.

Clinical trial specific metrics like comparison of survival probability as computed using Kaplan-Meier curves allow assessment of preservation of clinical insights in synthetic data

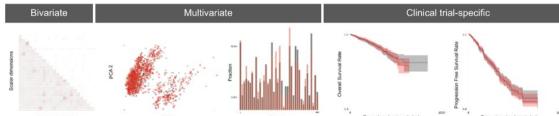
1. M Beigi, A Shafquat, J Mezey, JW Aptekar [Synthetic Clinical Trial Data while Preserving Subject-Level Privacy](#) - NeurIPS 2022 Workshop on Synthetic Data for Empowering ML Research
2. M Beigi, A Shafquat, J Mezey, JW Aptekar [Simulants: Synthetic Clinical Trial Data via Subject-Level Privacy-Preserving Synthesis](#) - AMIA 2022

Results – Fidelity & cross validation

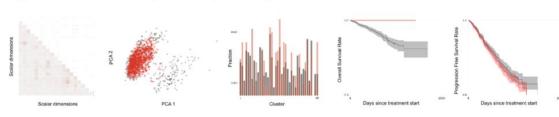
Benchmarks – Simulants outperforms all methods in Clinical Trial data

A. Algorithm = Simulants

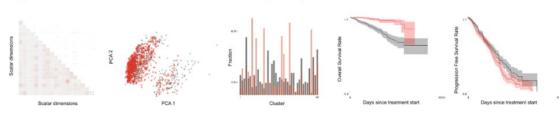
Dataset = Non Small Cell Lung Cancer



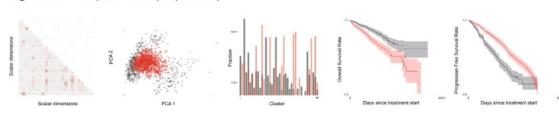
B. Algorithm = Gaussian Copula GAN (GC GAN)



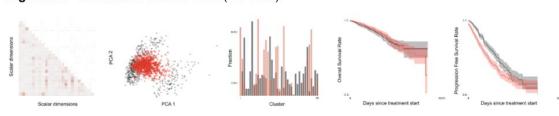
C. Algorithm = Tabular Variational AutoEncoder (TVAE)



D. Algorithm = Copula GAN (Cop. GAN)



E. Algorithm = Conditional Tabular GAN (CT GAN)



F. Fidelity quantification summary

Category	Univariate	Bivariate	Multivariate	Clinical trial - specific			
	Dimensions where K-S Chi-Square p<0.05	Correlation, mean absolute difference	Silhouette score	"Bag of Words" distance	log-rank Overall Survival K-M		
Best	0/N	0	0	1	1		
Worst	N/N	1	1	0	0		
Dataset	Algorithm						
Non small cell lung cancer	Simulants	0/171	0.11	0.000	0.042	0.47	0.89
	GC GAN	56/171	0.08	0.004	0.106	8.75E-26	0.46
	TVAE	98/171	0.10	0.014	0.215	5.62E-16	0.39
	Cop. GAN	64/171	0.22	0.040	0.236	9.63E-08	9.03E-09
	CTGAN	58/171	0.20	0.027	0.242	0.56	0.6
Diffuse Large B Cell Lymphoma	Simulants	2/174	0.07	0.000	0.024	0.99	NA
	GC GAN	88/174	0.09	0.009	0.210	3.31E-35	NA
	TVAE	94/174	0.12	0.026	0.159	3.82E-09	NA
	Cop. GAN	82/174	0.22	0.014	0.330	2.49E-04	NA
	CTGAN	83/174	0.22	0.018	0.320	3.02E-11	NA
Acute Lymphoblastic Leukemia	Simulants	11/142	0.08	0.000	0.013	NA	NA
	GC GAN	72/142	0.20	0.032	0.261	NA	NA
	TVAE	73/142	0.48	0.001	0.079	NA	NA
	Cop. GAN	101/142	1.14	0.017	0.172	NA	NA
	CTGAN	115/142	1.10	0.027	0.151	NA	NA
Acute Myeloid Leukemia	Simulants	0/108	0.19	0.000	0.081	0.06	NA
	GC GAN	36/108	0.29	0.007	0.192	0.68	NA
	TVAE	58/108	0.30	0.066	0.387	5.86E-35	NA
	Cop. GAN	44/108	0.42	0.032	0.304	0.11	NA
	CTGAN	36/108	0.42	0.021	0.311	0.20	NA

Takeaways

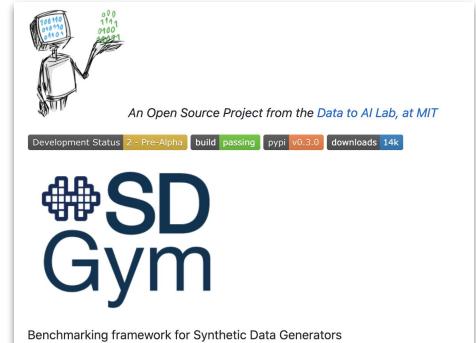
- Simulants outperforms all other approaches in fidelity benchmarks across multiple different clinical trial datasets (note: each dataset composed of 3-20 clinical trials)**

Results – Fidelity & cross validation

Benchmarks – Simulants outperforms all methods in canonical data too

Outperforms state-of-the-art GANs, VAEs and RNNs in canonical datasets at all scales in open benchmarks (SD-GYM / MIT AI Lab)

Group	Synthesizer	Description	Name	mnist-12	mnist-28 news	Real Datasets				Legend	
			Shape	60000x145	60000x785	31644x59	adult	covtype	credit		
Embedding	Simulants	Our method which uses the high dimensional embeddings.		0.9	0.9	0.1	0.8	0.8	1	1	0.9
Trivial	Identity	The synthetic data is the same as training data.		0.9	0.9	0.1	0.8	0.8	1	1	0.9
	Uniform	Each column in the synthetic data is sampled independently and uniformly.		0.1	0.1	-4	0.5	0.1	0.6	0.1	0.5
	Independent	Each column sampled independently. Continuous columns use Gaussian Mixture Model and discrete columns use the PMF of training data.		0.1	0.2	-0.06	0.6	0.4	0.9	0.7	0.7
GANs	Medgan	Minibatch averaging to efficiently avoid mode collapse.		0.4	0.1	-6	0.6	0.4	0.9	0.9	0.6
	VEEGAN	The method features a reconstructor network, reversing the action of the generator by mapping from data to noise.		0.4	0.2	-3.00E+08	0.7	0.2	0.9	0.5	0.8
	CTGAN	Models discrete and continuous columns.		0.1	0.1	-0.07	0.8	0.6	1	1	0.9
	CopulaGAN	Uses GaussianCopules to make the underlying CTGAN model task of learning the data.		0.2	0.2	-0.06	0.8	0.6	1	1	0.9
	TableGAN	Generates synthetic data using a convolutional neural network which optimizes the label column's quality.		0.1	0.1	-6	0.8	0	1	nan	0.9
VAEs	TVAE	Based on the VAE-based Deep Learning data synthesizer.		0.8	0.8	-0.02	0.8	0.7	1	1	0.9
Others	CLBN	Uses Bayesian networks.		0.7	0.2	-7	0.8	0.6	1	0.9	0.9
	PrvBN	A differential privacy method which uses a Bayesian network to model the correlation among the attributes.		nan	nan	nan	0.8	0.5	1	0.9	0.9
	GaussianCopulaCategorical	Based on copula functions and uses a CategoricalTransformer.		0.1	nan	-5	0.5	0.4	1	1	0.9
	GaussianCopulaCategoricalFuzzy	Based on copula functions and uses a CategoricalTransformer with fuzzy=True.		0.2	0.2	-9	0.8	0.4	1	0.8	0.8
	GaussianCopulaOneHot	Based on copula functions and uses a One-HotEncodingTransformer.		0.5	0.5	-4	0.8	0.5	1	0.9	0.9



Designing a privacy framework for synthetic clinical trial generation



Synthetic data is a solution to encourage collaboration and continue innovation in pharma



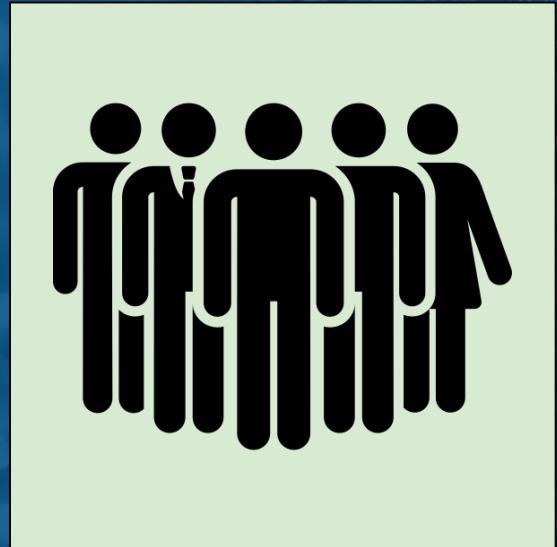
Sponsor A

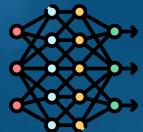


Sponsor B



Sponsor C





Personas

Data contributor

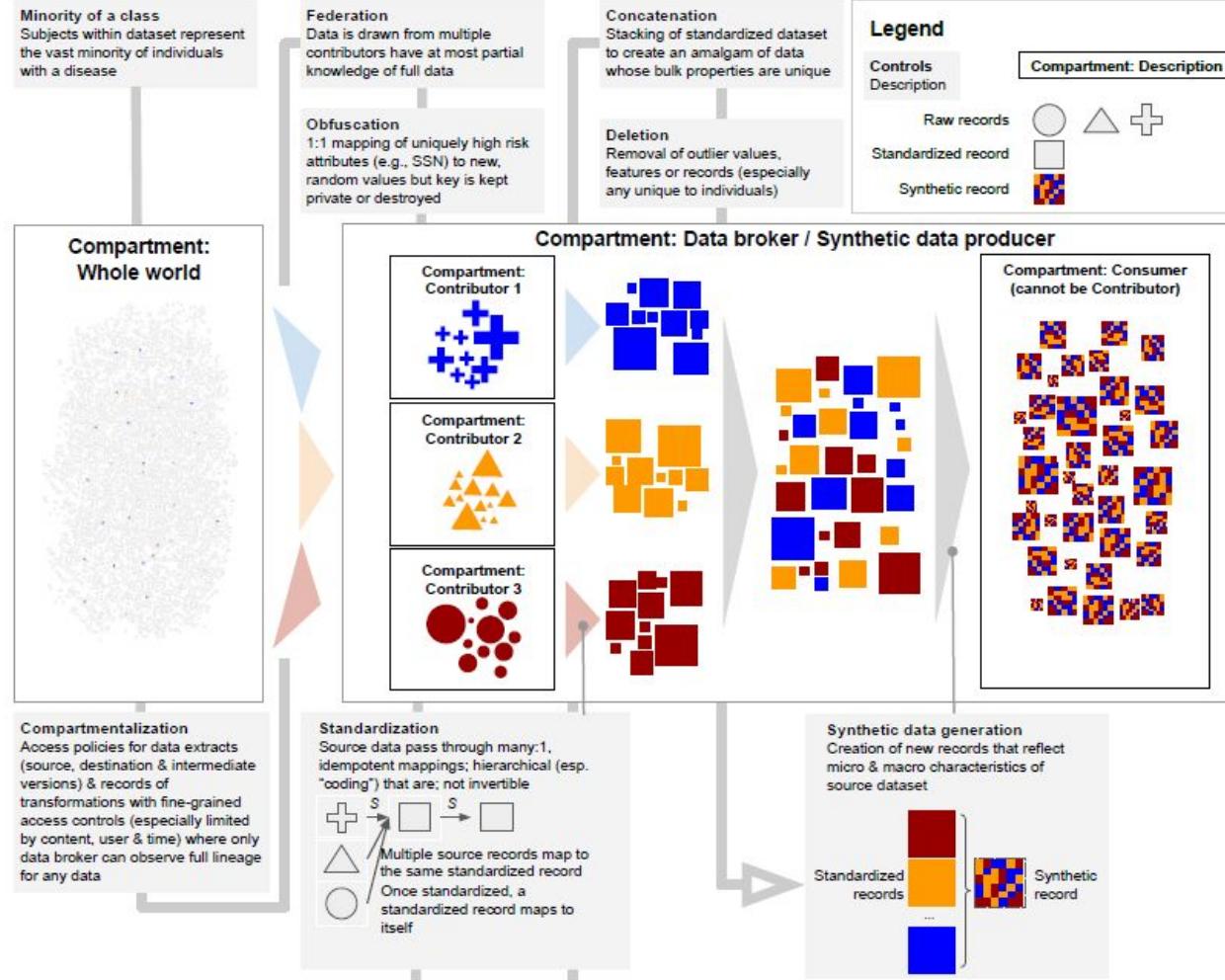
Synthetic Data broker

Data consumer

Competing interests

	Data Contributor privacy	Patient privacy	Privacy of proprietary information	Synthetic data fidelity
Data contributor	✓	✓	✓	
Synthetic Data broker	✓	✓	✓	✓
Data consumer				✓

Overview of privacy system design

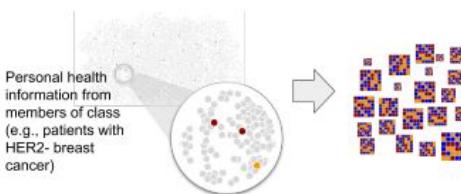
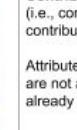
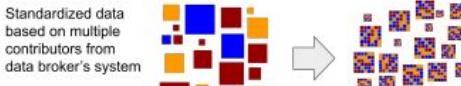
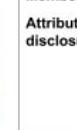


1. A Shafquat, J Mezey, M Beigi, J Sun, JW Aptekar **A source data privacy framework for synthetic clinical trial data-** NeurIPS 2022 Workshop on Synthetic Data for Empowering ML Research



Adversarial Scenarios

Attack scenarios

	Adversary	Defender	Key risk	Key safeguard	Feasibility / Privacy risk	Examples illustrating risk
A.	External attack 		Membership disclosure	Defender's data represents minority of class (i.e., only 1% of breast cancer patients are in dataset) 	High / Low	1% (4 of 300) Eligible studies in first line, EGFR mutant lung cancer compose source dataset shared with customers – representing the vast minority of possible studies; actual number of patients in dataset <0.01% patients with condition diagnosed in a single year
B.	Contributor attack 		Contributor disclosure (i.e., contributor reveals contribution to dataset) Attribute & Membership disclosure are not at risk because contributor already has this information	Compartmentalization by policy (i.e., contributors may not access datasets to which they contribute) 	Low / Low	4% (3,000 of 80,000) Original strings from a study in raw format appear in the synthetic dataset for a first line advanced stage lung cancer dataset; no tables or columns in common post standardization; all strings pre-processed for k-anonymity at k=2, 1-2% overlap between unrelated studies at baseline
C.	Omniscient attack 		Membership disclosure Attribute disclosure	Data synthesis algorithm (i.e., mapping from source→synthetic data formally limits disclosure risk, ϵ) 	Low / Low	7-13% Maximum improvement in attribute disclosure attacks for held-in vs held-out tranches of subjects in Lymphoma, Lung Cancer and Leukemia datasets (detailed privacy method in companion manuscript for this conference)

Legend

- Minority of class
- Compartmentalization
- Federation
- Obfuscation
- Standardization
- Deletion
- Synthetic data generation
- Active safeguard
- Inactive safeguard

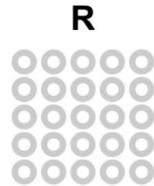
1. A Shafquat, J Mezey, M Beigi, J Sun, JW Aptekar A source data privacy framework for synthetic clinical trial data- NeurIPS 2022 Workshop on Synthetic Data for Empowering ML Research



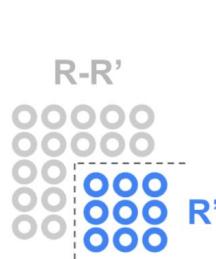
Overall privacy preservation score

A. Load source, R

\bullet = Real record



B. Partition into R' and $R-R'$

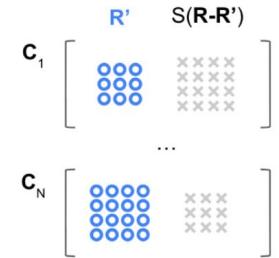


C. Produce synthetic datasets, $S(R-R')$ and $S(R)$

\times = Synthetic record



D. Repeat to create many mutually exclusive source-synthetic complement pairs, C



E. Create prediction scenarios, P

$P: \{x_1, \dots, x_j\} \rightarrow x_k$,
where x_k not in $\{x_1, \dots, x_j\}$

$$\left[P_1: \{x_1, \dots, x_j\} \rightarrow x_k \right]$$

...

$$\left[P_m: \{x_1, \dots, x_j\} \rightarrow x_k \right]$$

For each P , a random subset of features, $\{x_1, \dots, x_j\}$, is used to predict another random feature, x_k

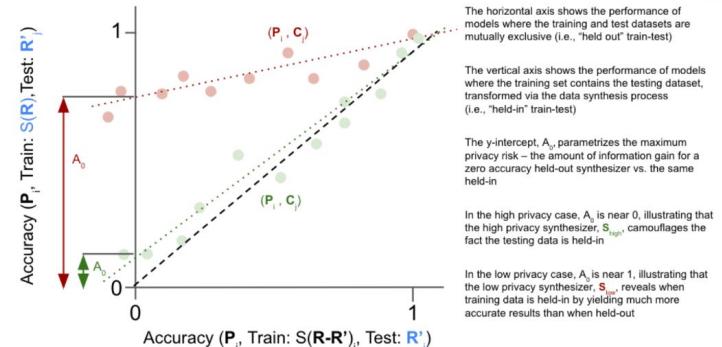
F. For each P_i and C_j ,

Assess held-out accuracy:
Train P_i on $S(R-R')_j$ and test accuracy in R'_j
(Refer to horizontal axis in G)

and compare to

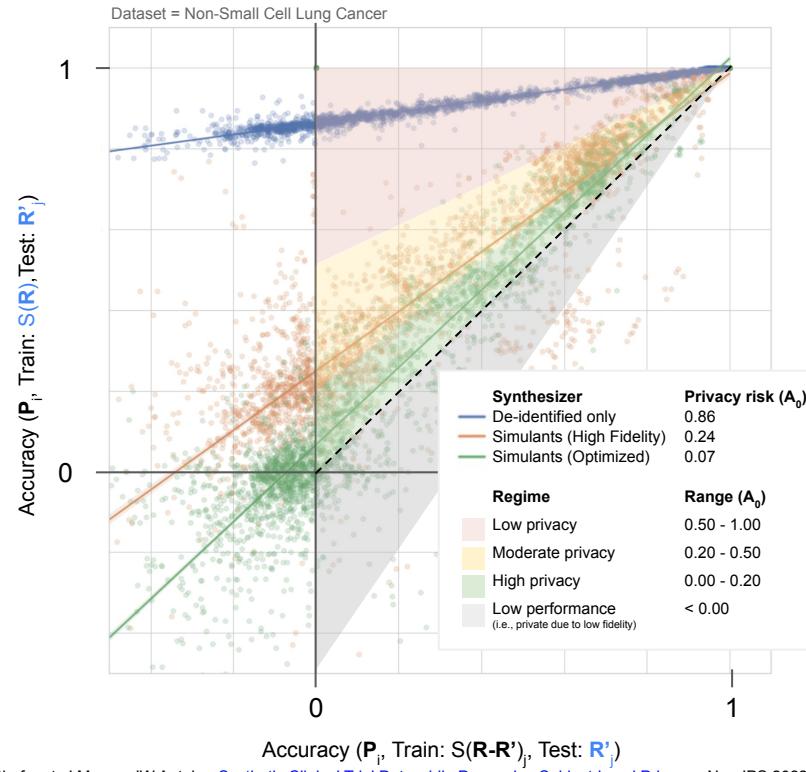
held-in accuracy:
Train P_i on $S(R)$ and test accuracy in R'_j
(Refer to vertical axis in G)

G. Illustration of results for low and high privacy synthesizers, S_{low} and S_{high}



Trial simulation- Standard tests, privacy

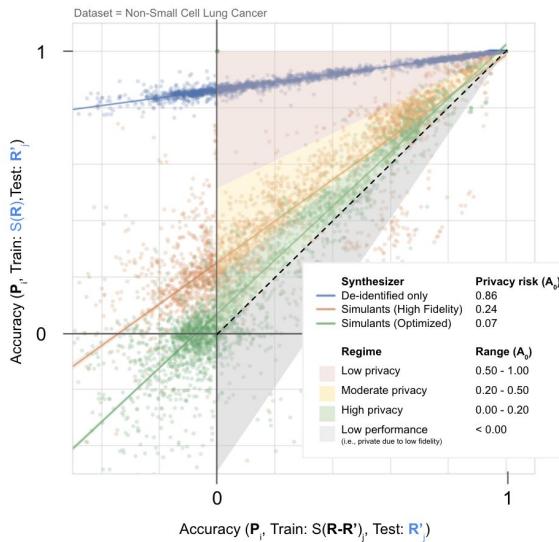
Representative privacy plots for high-fidelity & optimized Simulants



Results – Privacy

Benchmarks

H. Representative privacy plots for high-fidelity & optimized Simulants



I. Quantification of privacy risk (A_0) by synthesizer over datasets

Synthesizer	Dataset			
	NSCLC	DLBCL	AML	ALL
De-identified only	0.86	0.86	0.77	0.86
Simulants (Optimized)	0.07	0.09	0.13	0.08
Gaussian Copula GAN	0.05	0.05	0.04	0.02
TVAE	0.08	0.26	0.05	0.10
Copula GAN	-0.14	-0.01	-0.06	-0.14
CTGAN	-0.18	-0.12	-0.10	-0.10

Takeaways

- **Simulants delivers high privacy** at par with all other synthesizers tested (including Generative Adversarial Networks [GANs] and Variational Autoencoders [VAEs])
- Results for Novartis project are consistent with these benchmarks

Data Augmentation and applications in machine learning



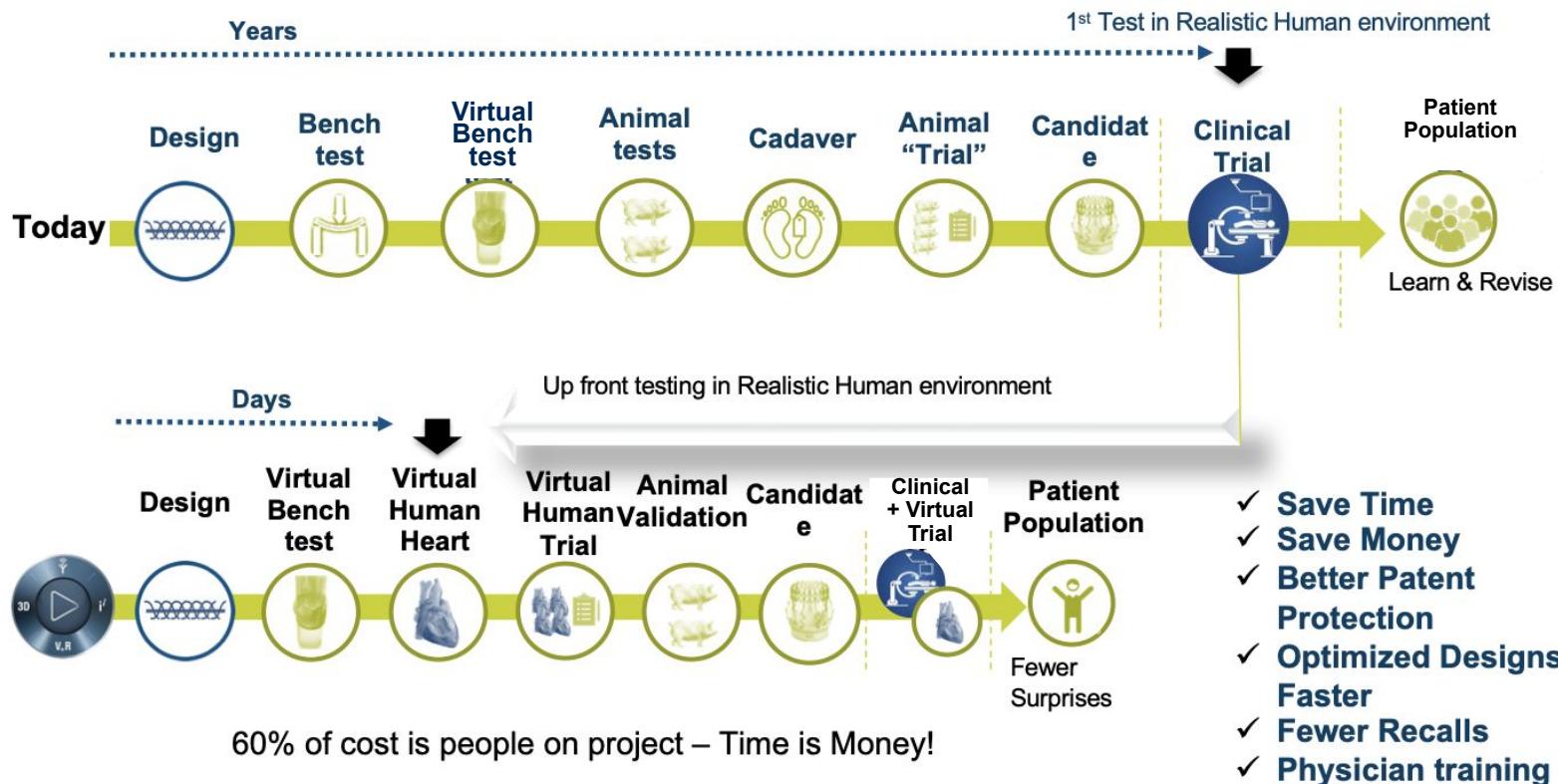
Accelerating *in silico* clinical trials



Full presentation available here:

<https://events.3ds.com/living-heart-and-virtual-twin-for-humans-symposium>

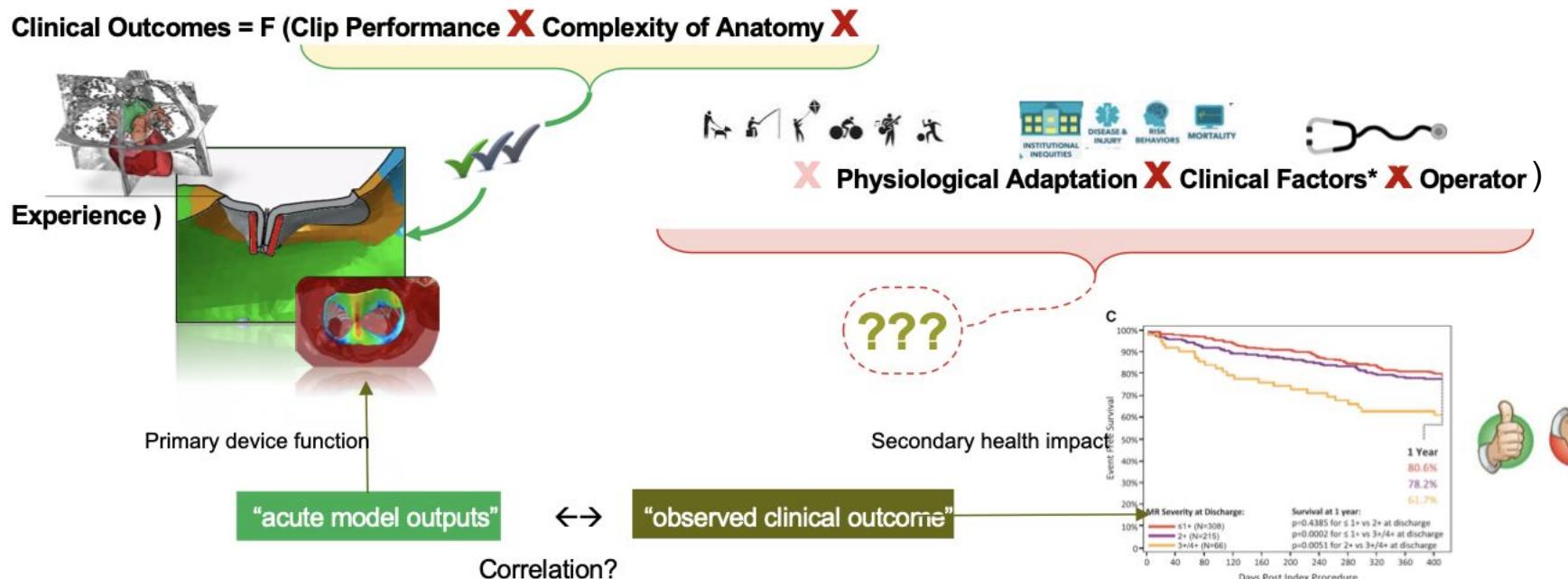
Value Proposition for *In Silico* Clinical Trial – “Years to Days”



- ✓ Save Time
- ✓ Save Money
- ✓ Better Patent Protection
- ✓ Optimized Designs Faster
- ✓ Fewer Recalls
- ✓ Physician training

iSCT Key Challenge – Defining & Proving the Hypothesis

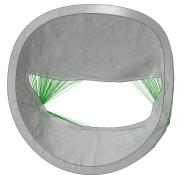
Identify strong correlation between acute model outputs (surrogate endpoints) & clinical endpoints



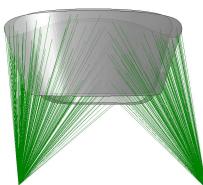
Virtual Patient Engine Schematic

Input:

Initial VPC - A collection of (physics-based) patient model definitions & pre-operative simulation results



Top view

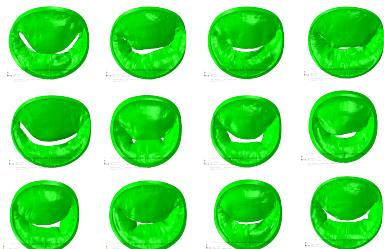


Front view



Output:

iSCT VPC - A physics-based VPC with targeted pre-operative characteristics to be treated in an iSCT

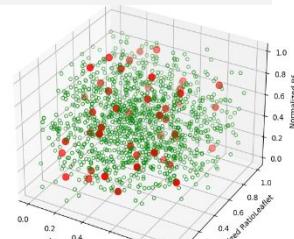


Machine Learning Powered VPE

Create surrogate model to accelerate VPC creation

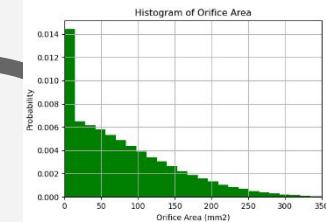


Obtain targeted representative VPs

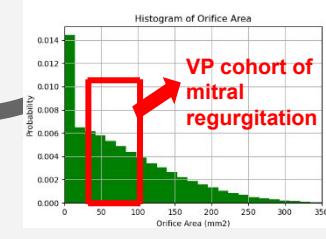


Build surrogate-based VPC

Build surrogate-based VPC



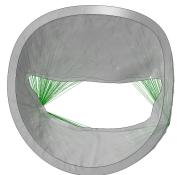
Down sample to get VPs with user defined characteristics



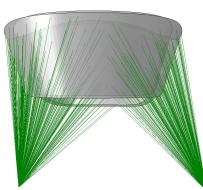
Virtual Patient Engine Schematic

Input:

Initial VPC - A collection of (physics-based) patient model definitions & pre-operative simulation results



Top view



Front view



Output:

iSCT VPC - A physics-based VPC with targeted pre-operative characteristics to be treated in an iSCT

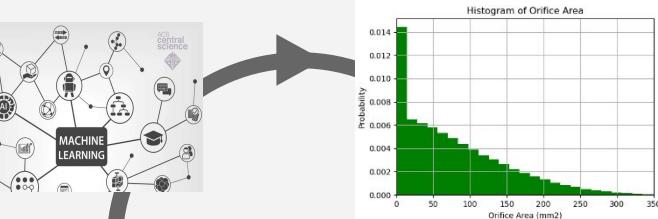
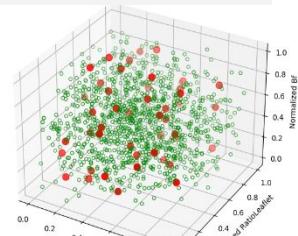


Machine Learning Powered VPE

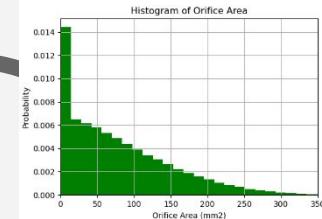
Create surrogate model to accelerate VPC creation



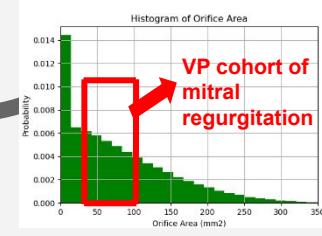
Obtain targeted representative VPs



Build surrogate-based VPC



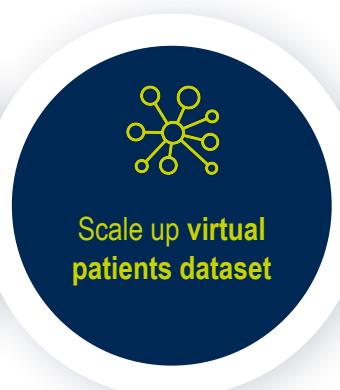
Down sample to get VPs with user defined characteristics



Virtual Patient Engine Schematic



Trained on virtual patients



Scale up virtual patients dataset



Reduces dataset to patient population of interest



Select from remaining ~100 virtual patients

Create surrogate model to accelerate VPC creation

Computationally expensive and impractical to scale

Build surrogate-based VPC

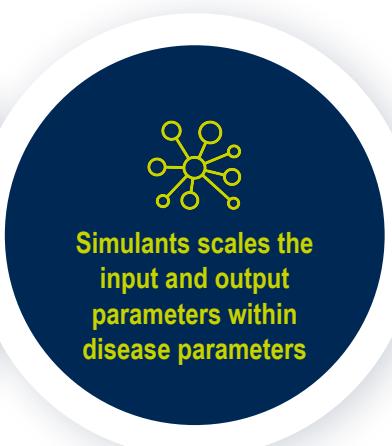
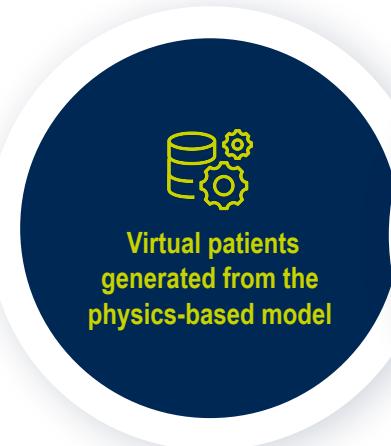
Generates redundant samples that don't pass the quality control criteria

Down sample to get VPs with user defined characteristics

Obtain targeted representative VPs

14x reduction in total sample size for selection

Virtual Patient Engine Schematic

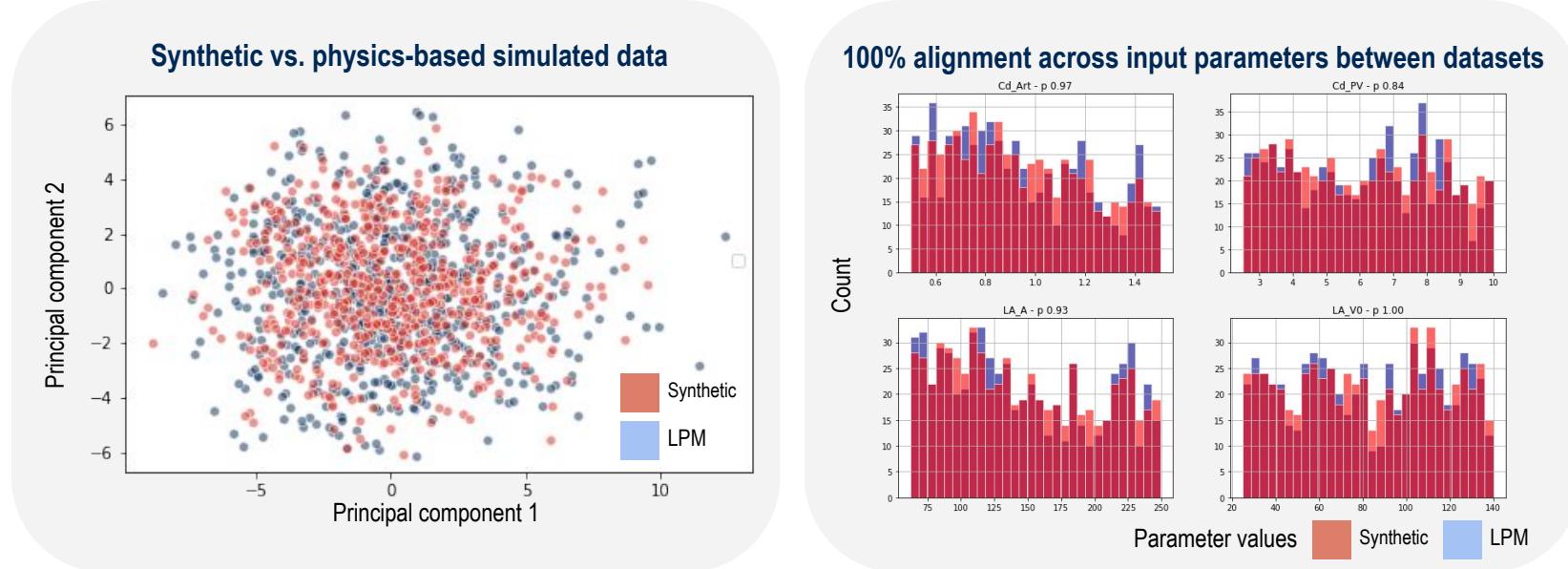


Training dataset generated using the initial physics-based population

Simulants pipeline generates synthetic data comprising of parameters of interest

Obtain targeted representative VPs

Synthetic data aligns with physics-based model simulated patients

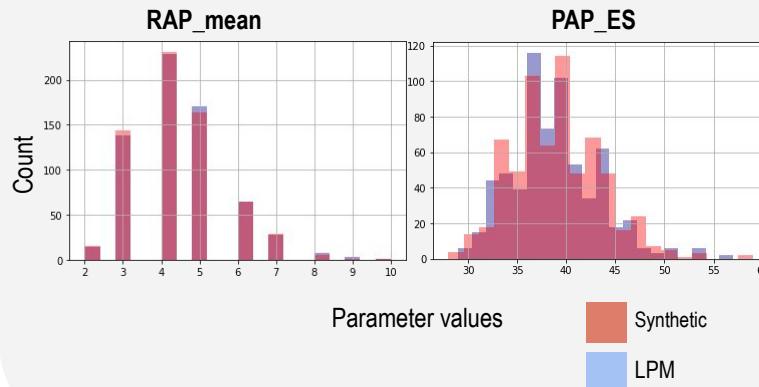


Alignment between source/template and Simulants generated synthetic dataset indicates synthesis reflects the source population well

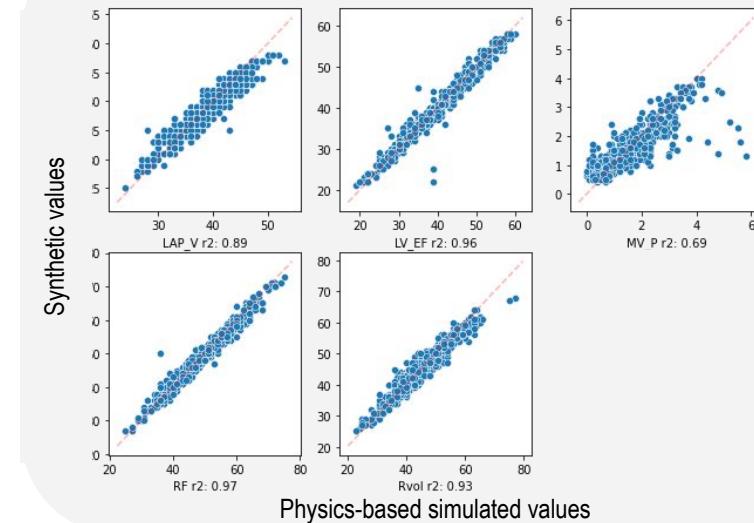
Alignment between input parameters is evaluated by performing the K-S test for input parameters produced by Simulants and the physics-based model. All 14 input parameters had non-statistically significant differences using the K-S test.

Synthetic data aligns with physics-based model simulated patients

97% alignment across output parameters between datasets



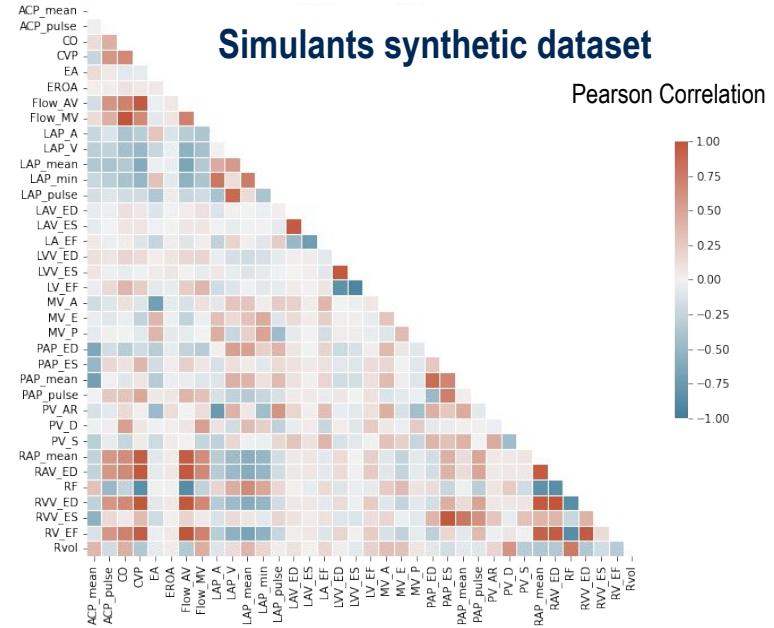
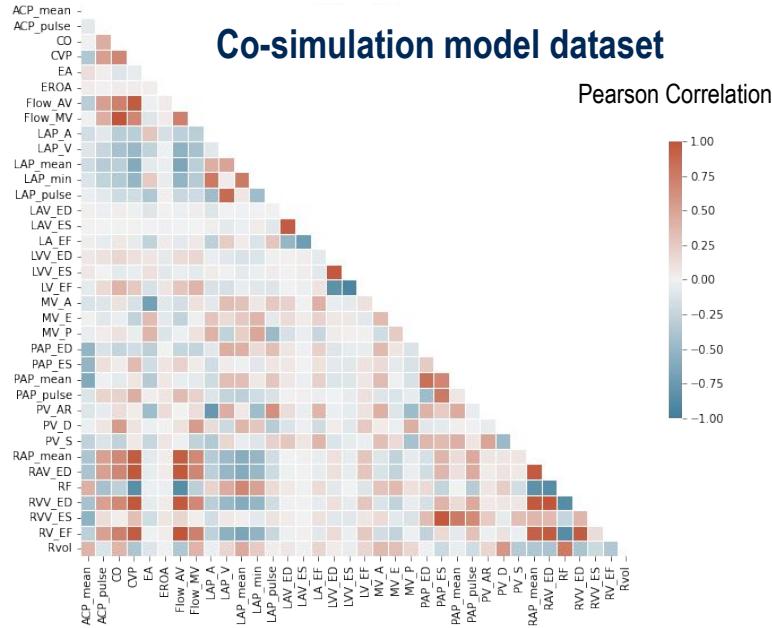
Synthetic vs. LPM simulated data



Alignment between output parameters is evaluated by performing the K-S test for output parameters produced by Simulants and the physics-based model. 34 parameters out of 35 parameters (i.e. MV_P) had non-statistically significant differences.

The diagonal in each plot indicates model performance on par with expectation. The x-axis shows the parameter values from the physics-based model and y-axis shows the Simulants-generated parameter values. High r^2 values (close to 1) indicate Simulants performance and agreement between the two models

Preservation of bidirectional relationships in synthetic data



Heatmaps indicate correlation across features in the source/physics-based simulated dataset and Simulants synthetic dataset. The similarity between the heatmaps indicates the bidirectional relationships and correlations observed in the dataset generated using the physics-based patient model are preserved in the Simulants-generated synthetic dataset.

Acceleration via synthetic data

97%

Agreement between simulated and synthetic output parameter distributions

100%

Agreement between simulated and synthetic input parameter distribution

100%

Preservation of sample size

Simulants is successfully able to mimic the the physics-based co-simulation model

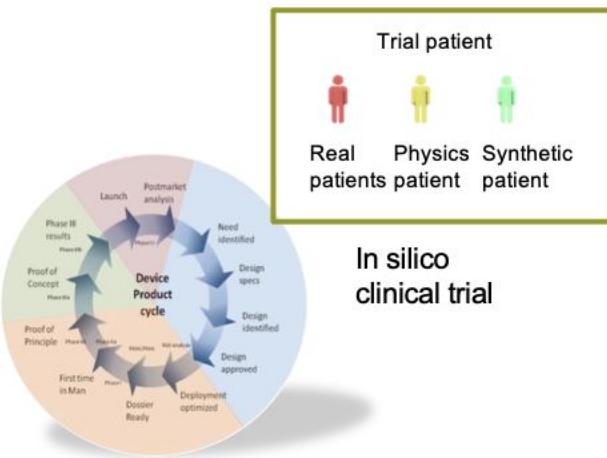
Summary of advantage of Simulants

Though the current VPE framework offers a competitive advantage over the Finite Element Analysis to simulate the living heart, a Simulants powered pipeline offers:

- **Fast and robust** way to generate synthetic parameter sets that reflect the physics-based parameter simulation
- **Reduce the computational cost** of running the Finite Element Analysis where the current process is slow and intensely computationally expensive
- **Doesn't depend on domain knowledge** of parameter boundary values
- **Preservation and scaling of sample size** of the physics-based dataset
- **Removal of redundancy** in the pipeline by only producing parameter sets that are aligned with the parameter distribution
- Potential to **target synthetic generation** of the distribution of interest (e.g. mitral valve regurgitation)

Future directions

Conduction 'in Silico clinical trials' using patient-specific models, physics based population models and synthetic patients to form virtual cohorts for testing the safety and/or efficacy of new drugs and of new medical devices.



Real patient (Patient specific models, acute and long end points, clinical outcome)

Physics patient (population models, detailed physics and physiology correlations, defining hypothesis)

Synthetic patient (Enrichment with VP augmentation, identifying correlation between acute and clinical outcome, proving hypothesis)

Key takeaways



**Data sharing is critical
for continued innovation
in the biotech and
pharma industry**



**Synthetic data provides
a fast, secure and
reliable way to share
private data**



**Synthetic data allows an
innovative way to synthesize and
augment training datasets to
improve AI/ML model performance**

Presentation credit:



Jacob Aptekar, MD PhD
VP, Medidata AI



Mandis Beigi, PhD
Sr. Director, Medidata AI



Jia Chen
Senior Director, Medidata AI

A large background image showing a medical professional in a white coat and stethoscope, and a close-up of a laptop keyboard. A circular portrait of a woman with long dark hair and a smile is overlaid on the right side.

Afrah Shafquat
Sr. Data Scientist II, Medidata AI

Afrah.Shafquat@3ds.com

afrahshafquat

DASSAULT SYSTEMES

Credits

1. Neural-network icons created by Freepik - Flaticon
2. Corporate icons created by Sumitsaengtong - Flaticon
3. Consumer icons created by Freepik - Flaticon
4. Consumer icons created by Freepik - Flaticon
- 5.
6. Sheet icons created by Andrejs Kirma - Flaticon
7. Neural network icons created by Freepik - Flaticon
8. Patient icons created by SBTs2018 - Flaticon
9. Privacy icons created by Fathema Khanom - Flaticon
10. Clinical icons created by Parzival' 1997 - Flaticon
11. Collaboration icons created by Freepik - Flaticon
12. Collaboration icons created by Freepik - Flaticon
13. Group icons created by Freepik - Flaticon
14. Group icons created by Freepik - Flaticon
15. Denied icons created by Alfredo Creates - Flaticon
16. Patient icons created by SBTs2018 - Flaticonk icons created by Freepik - Flaticon
17. Privacy icons created by Freepik - Flaticon
18. Collaboration icons created by small.smiles - Flaticon
19. Collaboration icons created by small.smiles - Flaticon
20. Medicine icons created by max.icons - Flaticon
21. Simulation icons created by Freepik - Flaticon

