



Deep Learning for Healthcare

Deep Generative
Models

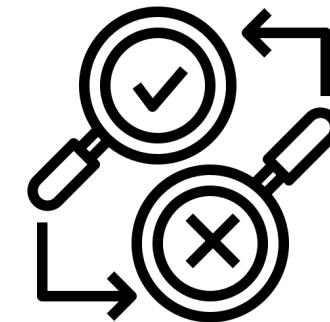
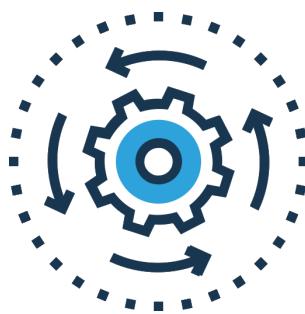
Jimeng Sun

Outline

- Generative Adversarial Network (GAN)
 - Method
 - Application: Synthetic patient record generation
- Variational Autoencoder (VAE)
 - Method
 - Application: Molecule generation

Generative Adversarial Network (GAN)

- **Goal:** to create new samples that resemble training data
- **Idea:** train two neural networks:



Generator for creating
realistic but synthetic samples

Discriminator for differentiating
synthetic and real samples

Healthcare applications of GAN



Generate synthetic electronic health records

Choi et al. 2017. "Generating Multi-Label Discrete Patient Records Using Generative Adversarial Networks." In *Proceedings of the 2nd Machine Learning for Healthcare Conference*

Zhang et al. 2020. "Ensuring Electronic Medical Record Simulation through Better Training, Modeling, and Evaluation." *Journal of the American Medical Informatics Association: JAMIA* 27 (1): 99–108.



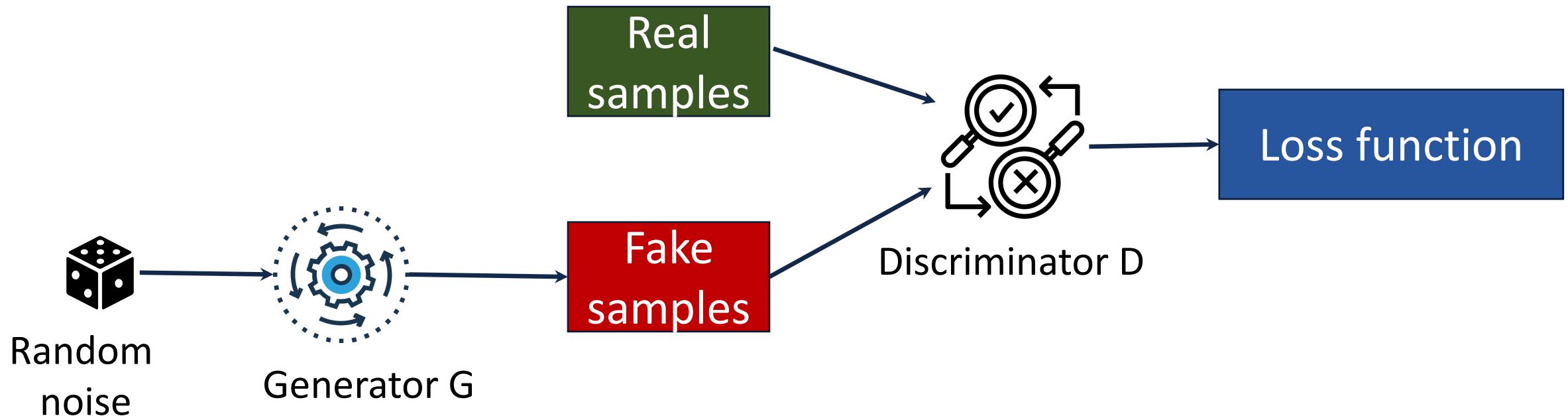
De Novo design: molecule generation

De Cao, Nicola, and Thomas Kipf. 2018. "MolGAN: An Implicit Generative Model for Small Molecular Graphs." *arXiv [stat.ML]*. arXiv. <http://arxiv.org/abs/1805.11973>.

Jin, Wengong, Kevin Yang, Regina Barzilay, and Tommi Jaakkola. 2018. "Learning Multimodal Graph-to-Graph Translation for Molecular Optimization." *arXiv [cs.LG]*. arXiv. <http://arxiv.org/abs/1812.01070>.

Maziarka, Łukasz, Agnieszka Pocha, Jan Kaczmarczyk, Krzysztof Rataj, Tomasz Danel, and Michał Warchoł. 2020. "Mol-CycleGAN: A Generative Model for Molecular Optimization." *Journal of Cheminformatics* 12 (1): 2.

GAN Structure



- Generator learns to create realistic data (fake samples)
- Discriminator learns to separate fake samples with real samples

GAN Loss function

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))]$$

Log-likelihood that real samples are real

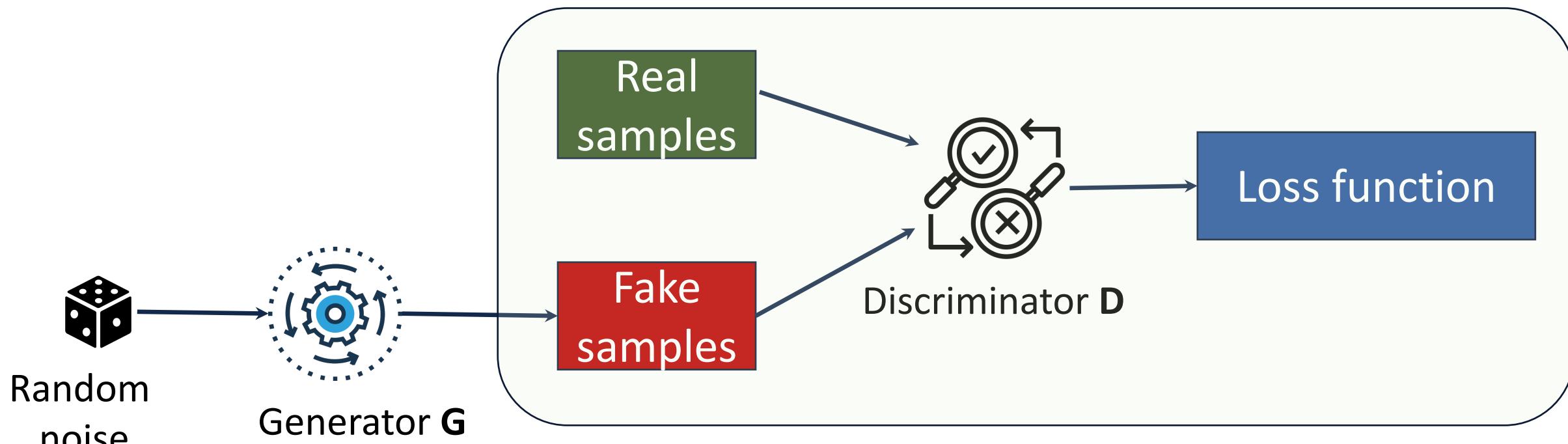
Log-likelihood that fake samples are fake

Sample **real data**

Sample **random noise**

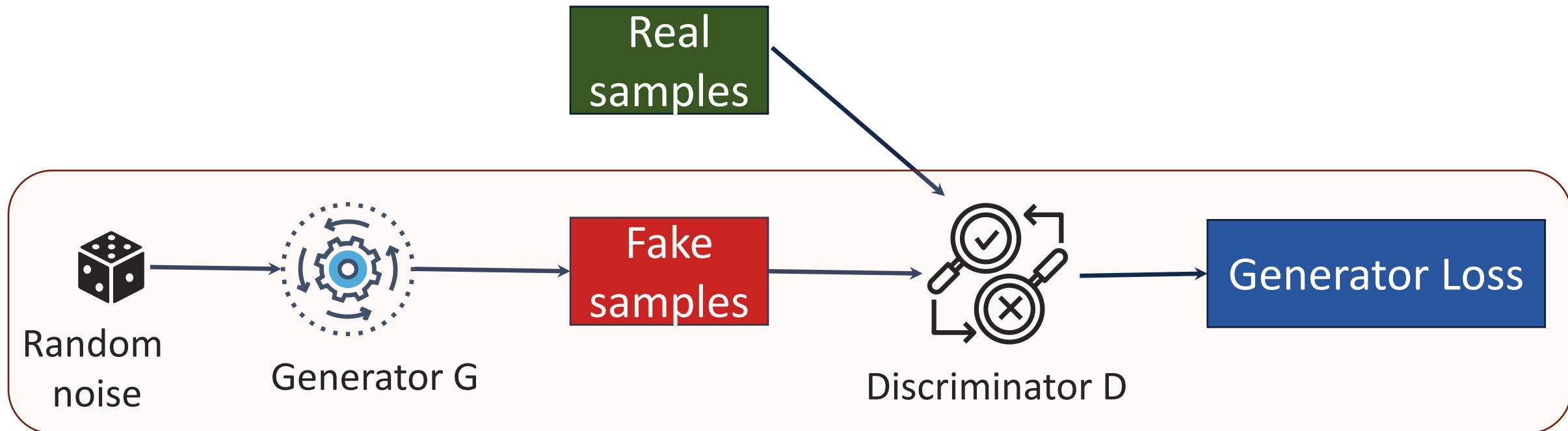
A diagram illustrating the GAN Loss function. The equation is $\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))]$. Above the equation, there are two curly braces. The left brace covers the first term $\mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})]$ and is labeled "Log-likelihood that real samples are real". It points to the text "Sample **real data**". The right brace covers the second term $\mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))]$ and is labeled "Log-likelihood that fake samples are fake". It points to the text "Sample **random noise**".

Discriminator D



- D classifies real and fake samples
- Discriminator loss penalizes any mistake for misclassifying fake from real
- Update parameters via backprop from Discriminator loss

Generator G



- Random noise will be converted to realistic samples (fake samples)
- D classifies the fake samples
- Generator loss penalizes for failing to fool D
- Backprop through D and G

Generator Loss

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))]$$

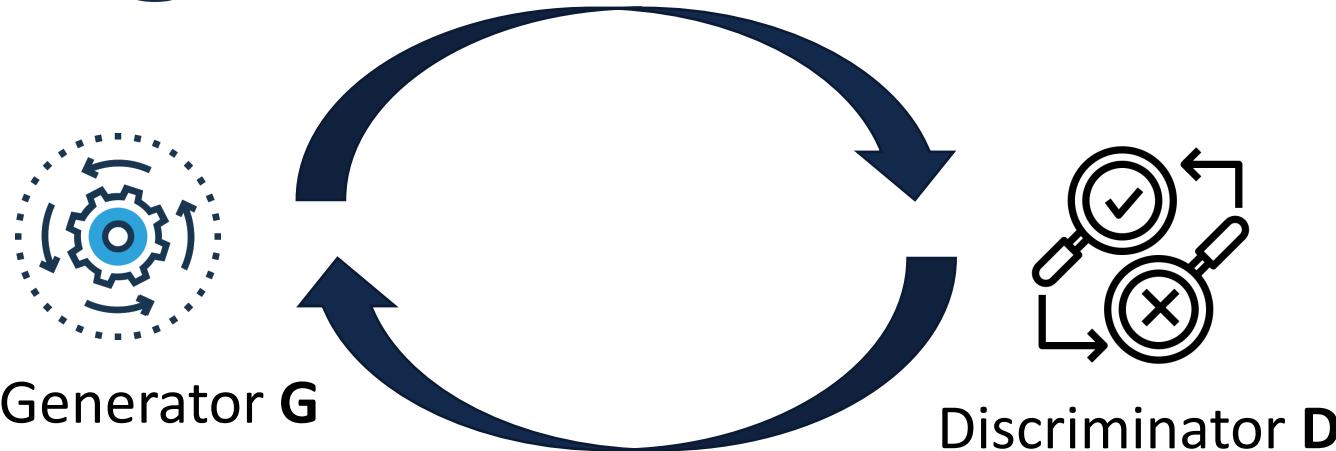
Generator Loss

- The generator loss is equivalent to

$$L_G = -\mathbb{E}_{\mathbf{z}} \log D(G(\mathbf{z}))$$

which is easier to train

Training GAN



- Iteratively, train D for one or few epochs, then train G for one or few epochs
- Convergence of GAN training is still an open problem
 - No reliable metrics to indicate the convergence
 - Overtraining can lead to performance degradation

Generating Multi-Label Discrete Patient Records using Generative Adversarial Networks

How to generate realistic EHR data

Choi, Edward, et al. 2017. “Generating Multi-Label Discrete Electronic Health Records Using Generative Adversarial Networks.” *Machine Learning for Health Care 2017*.

Healthcare Research and Data

- Massive data are electronically collected by healthcare organizations
- Privacy concerns often hinder data science research in healthcare



Approaches to handle privacy concerns

- De-identification by perturbation methods
 - Derived from the real data; not impregnable to attacks
- Generating synthetic datasets
 - Has not been realistic enough for machine learning tasks



Objective of MedGAN



Generate realistic EHR data using Generative Adversarial Network (GAN)



- High-quality: Synthetic datasets can train ML models that perform similarly to the ones trained on the real datasets

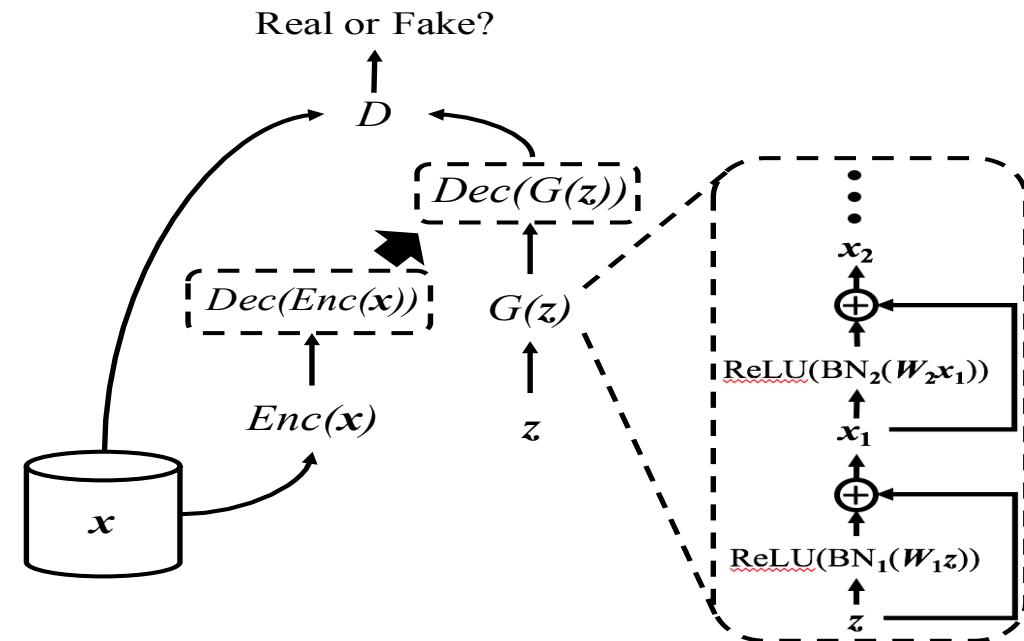


- Privacy-preserving: Practically impossible to gain knowledge on real patients from the synthetic dataset

MedGAN Architecture



- x : real-data
- z : random noise
- Enc : encoder of autoencoder
- Dec : decoder of autoencoder
- G : Generator
- D : Discriminator
- BN : Batch-norm



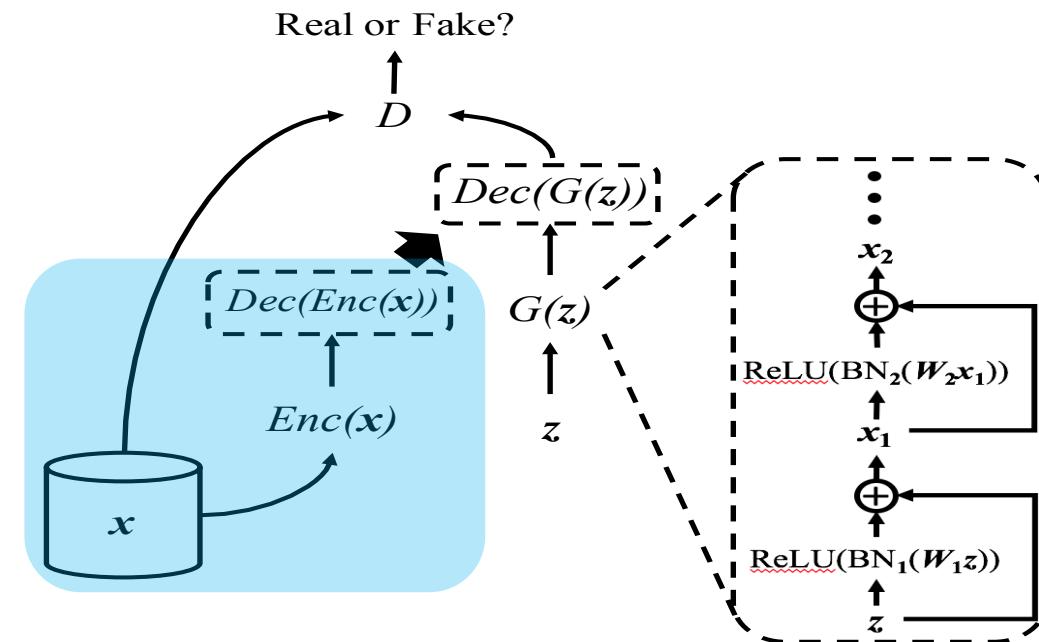
Generating completely synthetic data using GAN

MedGAN Architecture

Pre-train Autoencoder

$$\frac{1}{m} \sum_{i=0}^m \mathbf{x}_i \log \mathbf{x}'_i + (1 - \mathbf{x}_i) \log(1 - \mathbf{x}'_i)$$

where $\mathbf{x}'_i = Dec(Enc(\mathbf{x}_i))$



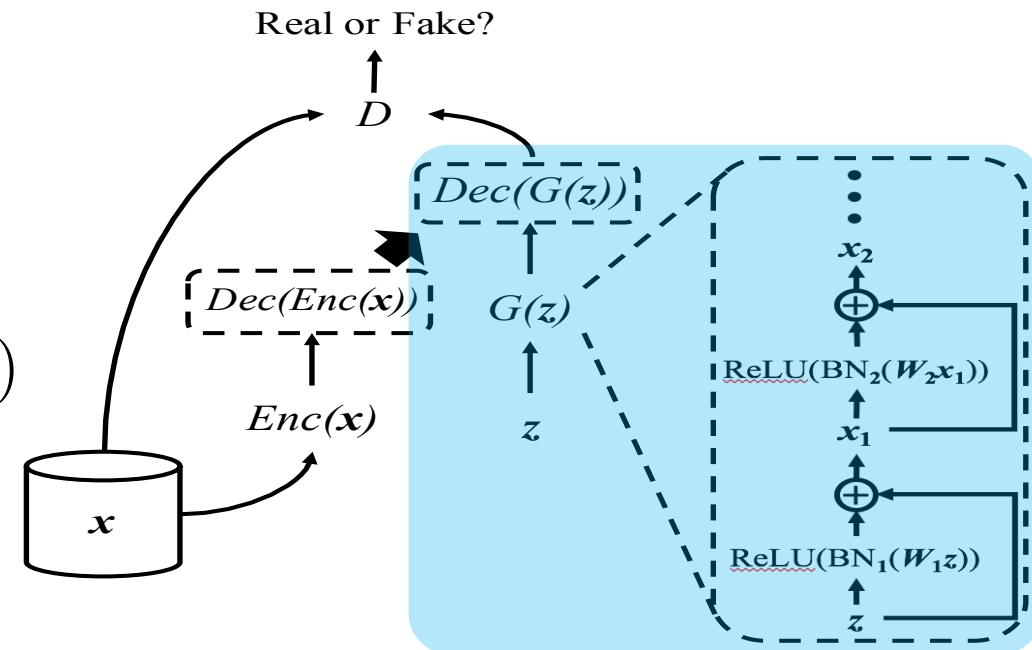
MedGAN Architecture



Update Generator parameters

$$\theta_{g,dec} \leftarrow \theta_{g,dec} + \alpha \nabla_{\theta_{g,dec}} \frac{1}{m} \sum_{i=1}^m \log D(\mathbf{x}_{\mathbf{z}_i})$$

where $\mathbf{x}_{\mathbf{z}_i} = Dec(G(\mathbf{z}_i))$

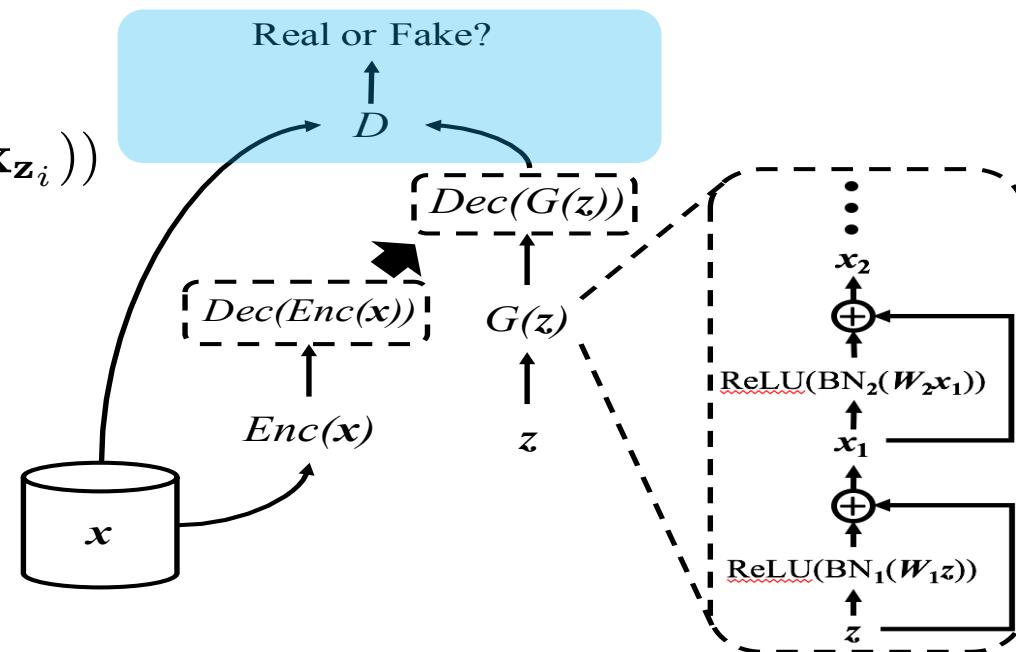


MedGAN Architecture

Update Discriminator parameters

$$\theta_d \leftarrow \theta_d + \alpha \nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m \log D(\mathbf{x}_i) + \log(1 - D(\mathbf{x}_{z_i}))$$

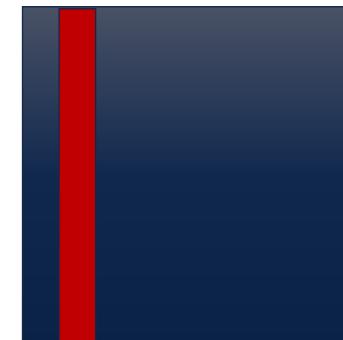
where $\mathbf{x}_{z_i} = Dec(G(z_i))$



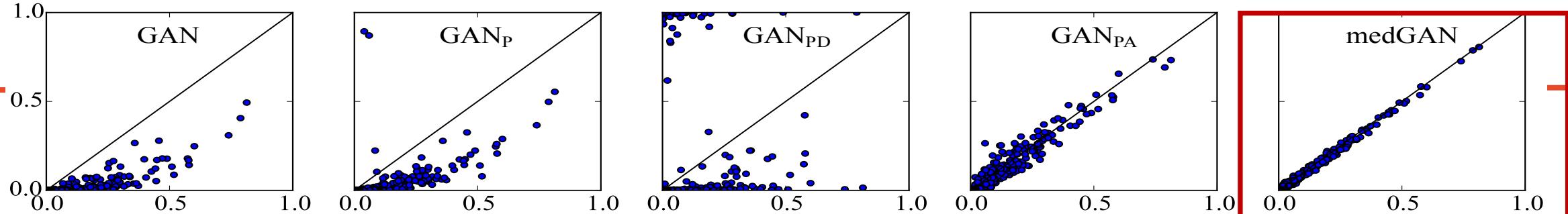
Quantitative Evaluation



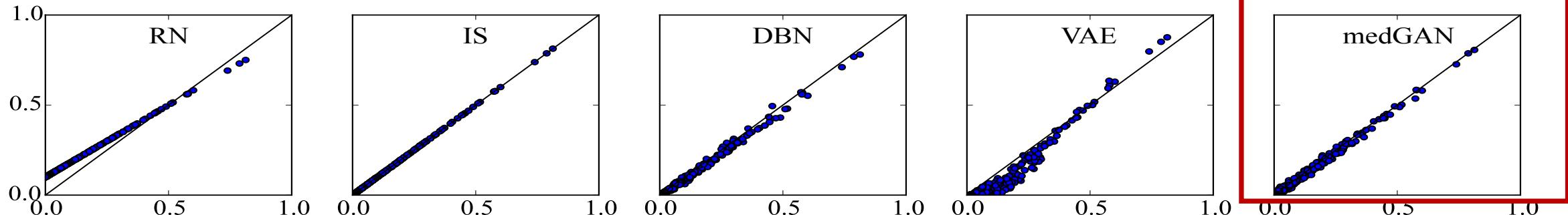
- Dimension-wise probability
 - Compare the dimension-wise probability of the training data and the synthetic data
- Dimension-wise prediction
 - Assess how well medGAN captures the inter-dimensional relationships
 - Select one dimension k as the label, use the remaining dimensions as features for a logistic regression (LR) classifier.
- Experiments conducted on Sutter PAMF



Dimension-wise Probability



(a) Dimension-wise probability performance of various versions of medGAN.

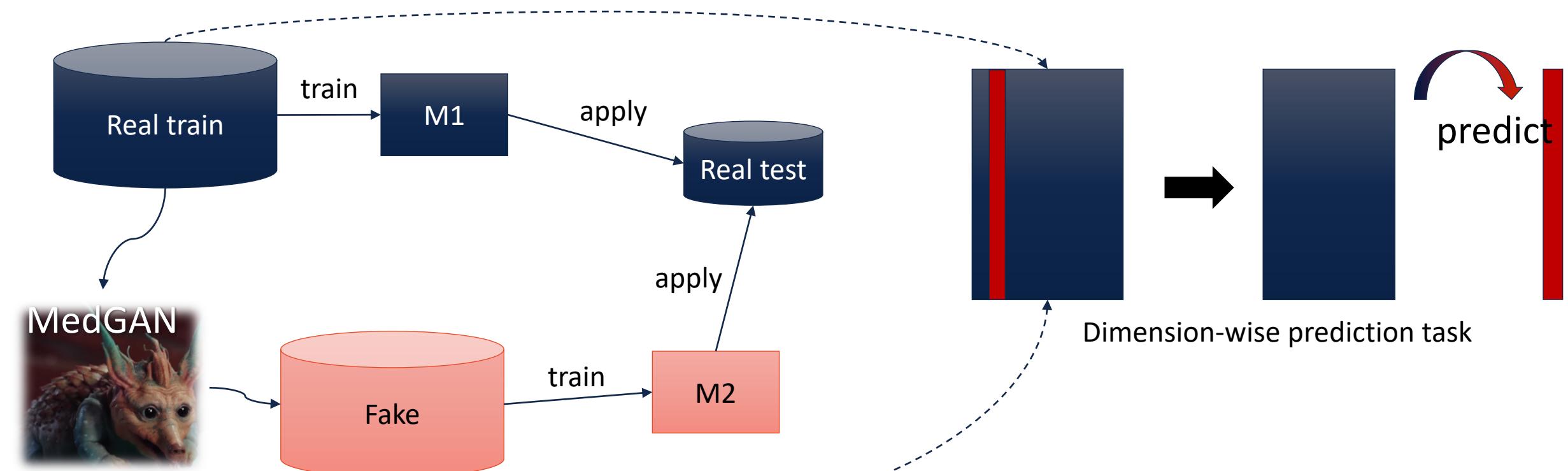


(b) Dimension-wise probability performance of baseline models and medGAN.

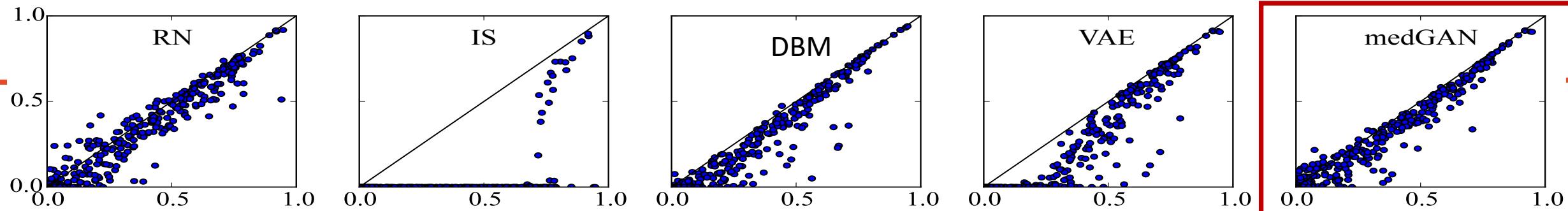
- X-axis: dimension-wise probability of real data
- Y-axis: dimension-wise probability of synthetic data
- A dot represents a single dimension (i.e. medical code)
- Diagonal line indicates the ideal performance

Quantitative Evaluation: Dimension-wise prediction

Select one dimension as the target, use the remaining dimensions as features for a logistic regression (LR) classifier.



Dimension-wise Prediction



Real data +
random noise

Independent
dimensional wise
sampling

Deep Boltzmann
machine

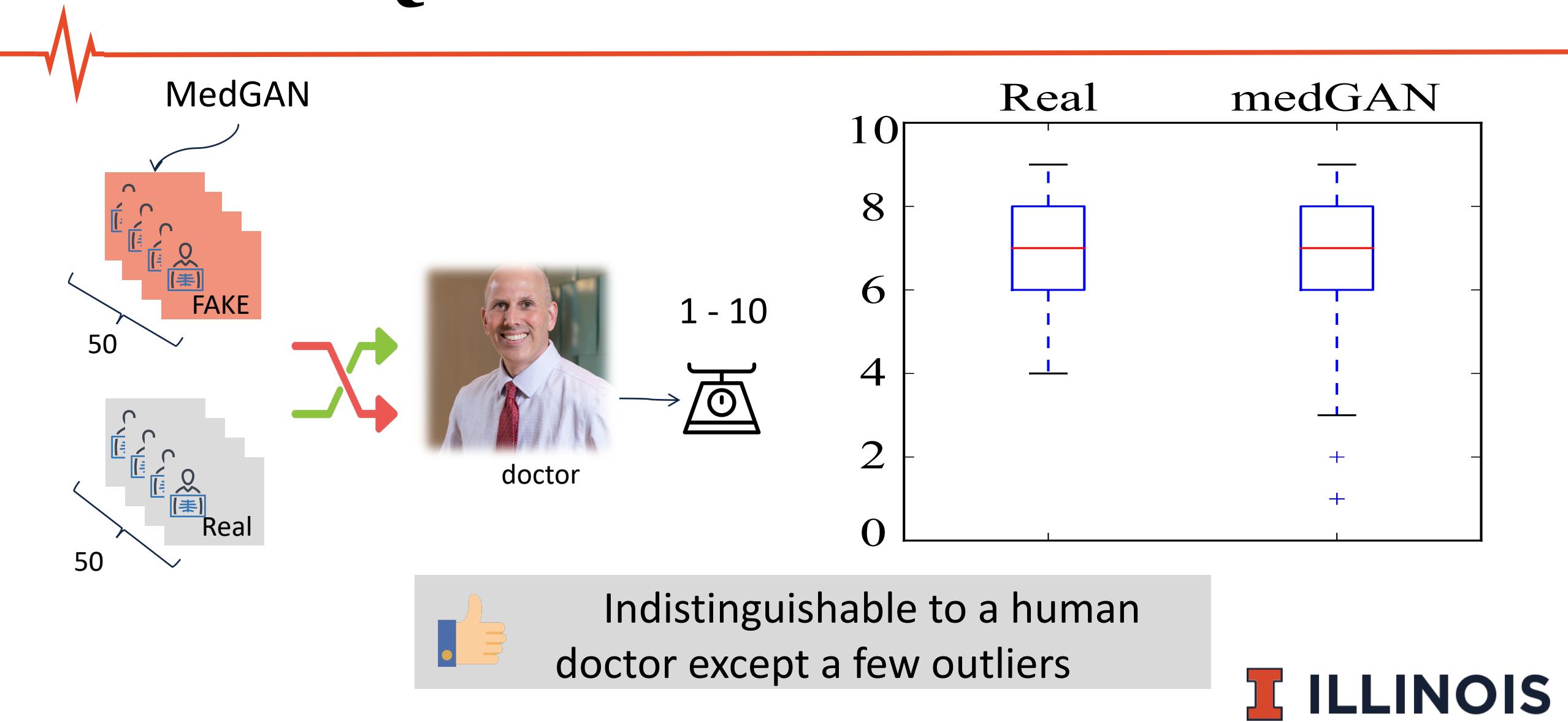
Variational
autoencoder

- X-axis: F1-score of LR trained by real data
- Y-axis: F1-score of LR trained by synthetic data
- A dot represents a single dimension selected as the label
- Diagonal line indicates the equal performance with real data



Data generated by medGAN performs great for ML tasks

Qualitative Evaluation



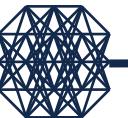
Conclusion



- MedGAN generates high-quality synthetic data for ML tasks
 - Dimension-wise prediction
 - Qualitative evaluation by a human doctor
- MedGAN preserves privacy
 - Attacker can confirm the presence of very small number of patients
 - Attacker cannot gain useful knowledge on unknown attributes

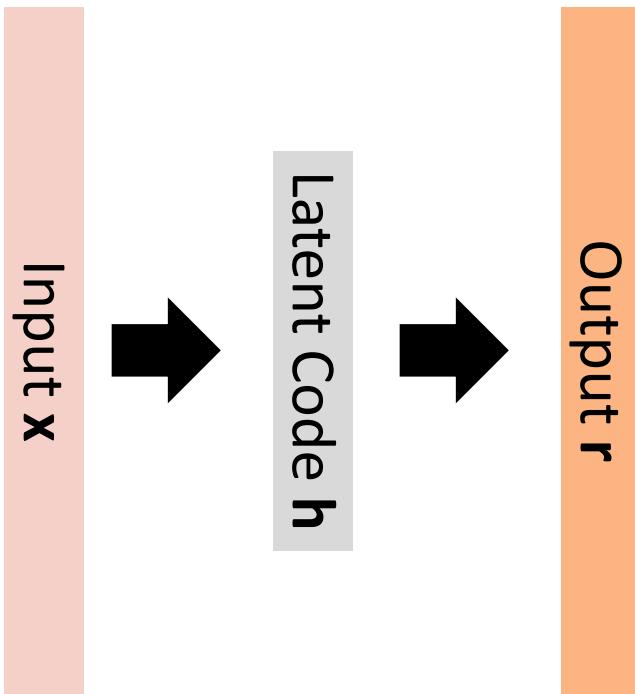
Choi, Edward, et al. 2017. “Generating Multi-Label Discrete Electronic Health Records Using Generative Adversarial Networks.” *Machine Learning for Health Care 2017*.

Variational Autoencoder (VAE)



- VAE is a generative model for creating realistic data samples
- VAE is in the intersection of **deep learning** and **probabilistic graphical model**
- VAE has many healthcare applications in
 - Molecule generation
 - Medical imaging analysis

Review: Autoencoders

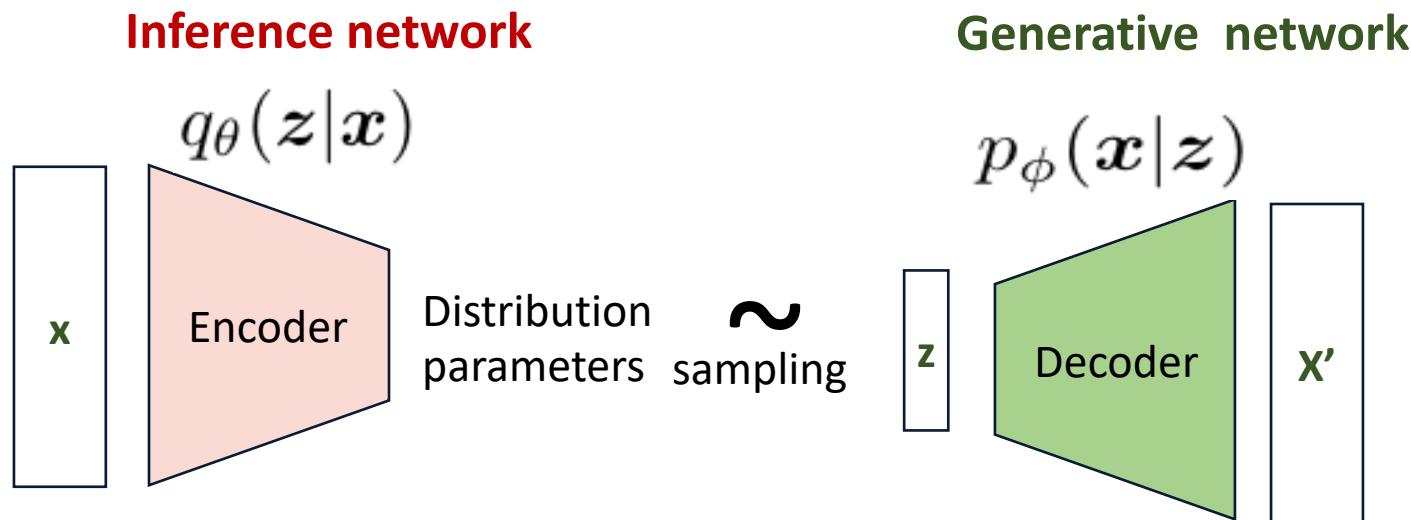
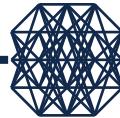


Minimize reconstruction error

$$L(\mathbf{x}, \mathbf{r}) = \|\mathbf{x} - \mathbf{r}\|^2 \quad \text{for Gaussian input}$$

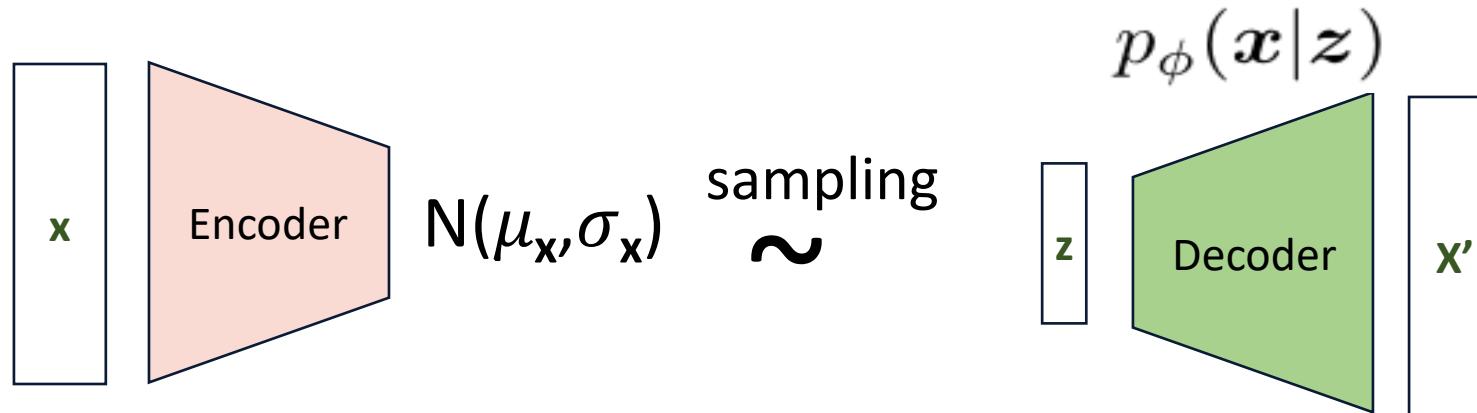
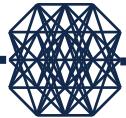
$$L(\mathbf{x}, \mathbf{r}) = - \sum_i [\mathbf{x}_i \log \mathbf{r}_i + (1 - \mathbf{x}_i) \log(1 - \mathbf{r}_i)] \quad \text{for binary input}$$

VAE: Deep learning view (high-level)



- Encode input as a distribution
- Sample a point from the encoding distribution
- Decode from the sample to generate a realistic data point

VAE: Deep learning view



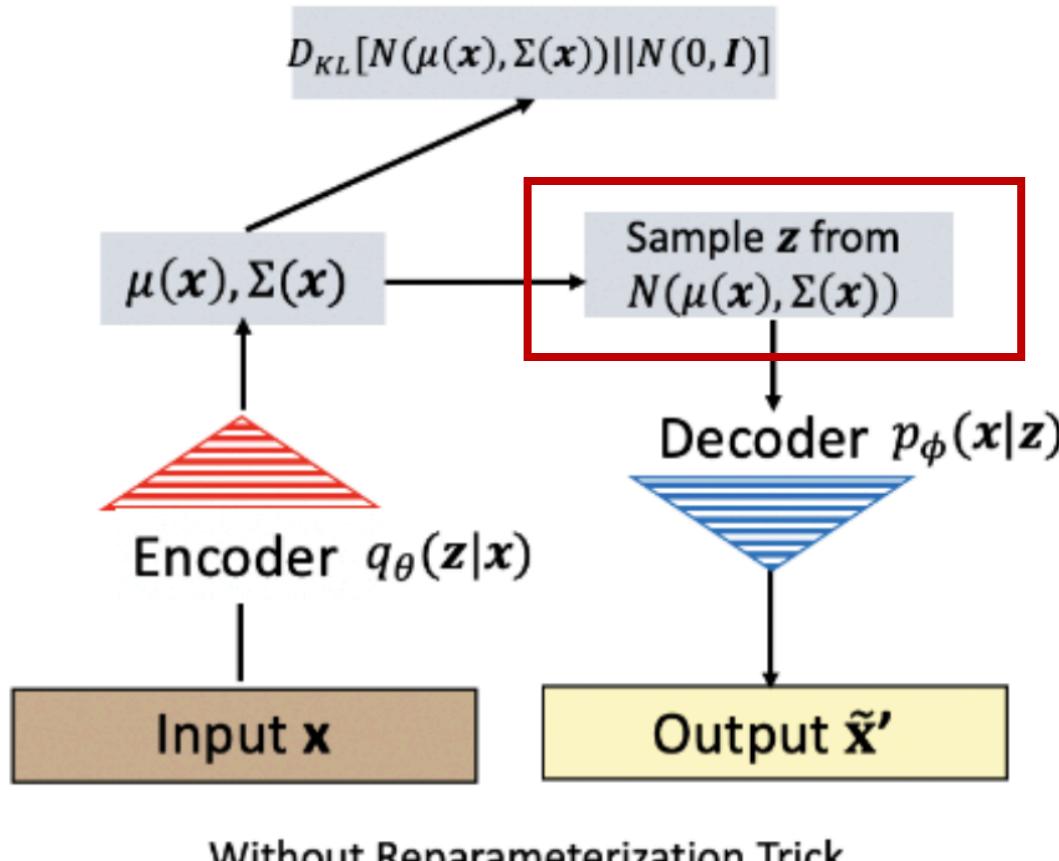
- Encoder

$$q_{\theta}(\mathbf{z}|\mathbf{x}) \sim \mathcal{N}(\mu(\mathbf{x}), \Sigma(\mathbf{x})) \quad p_{\phi}(\mathbf{z}) \sim \mathcal{N}(0, \mathbf{I})$$

- Loss function for data point \mathbf{x}

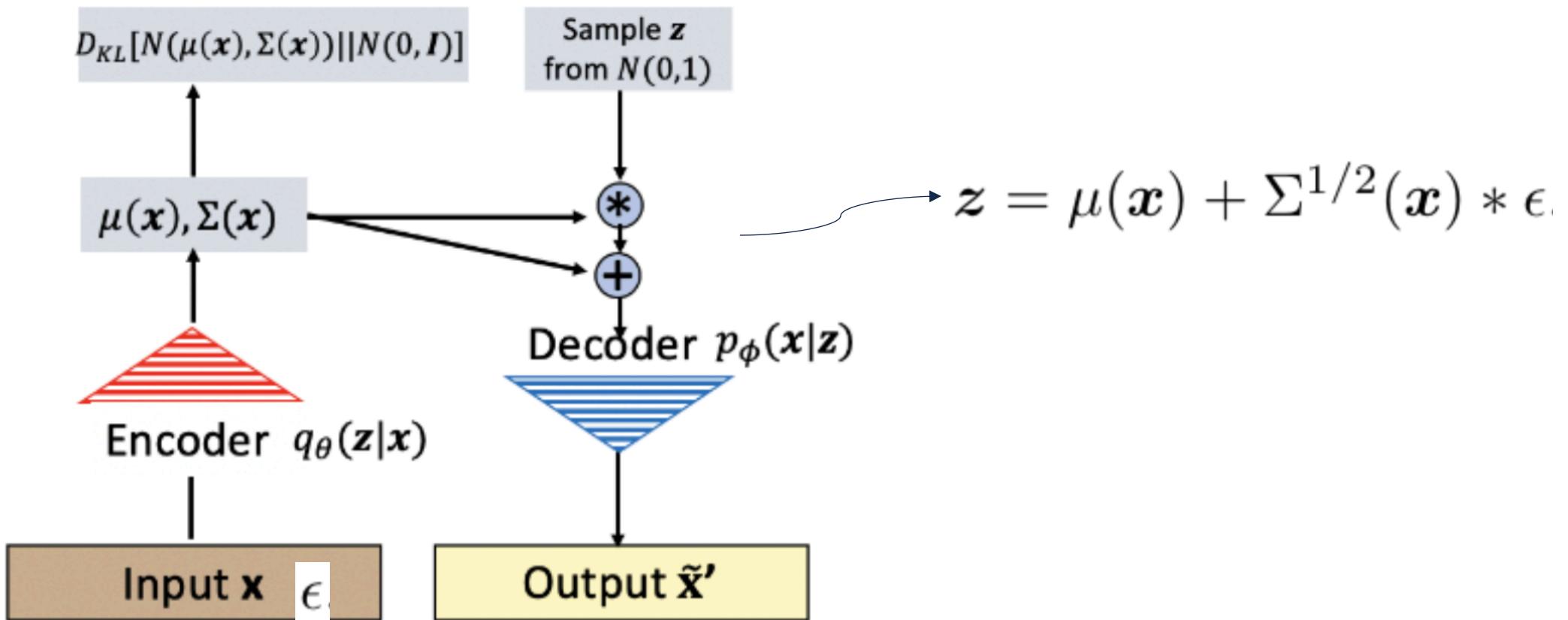
$$l_{\mathbf{x}} = -\mathbb{E}_{\mathbf{x} \sim D} [\log p_{\phi}(\mathbf{x}|\mathbf{z}) + \mathbb{D}_{\text{KL}}(q_{\theta}(\mathbf{z}|\mathbf{x}) \| p_{\phi}(\mathbf{z}))]$$

VAE: Reparameterization trick



- Sampling step makes backprop difficult

VAE: Reparameterization trick (cont.)



With Reparameterization Trick

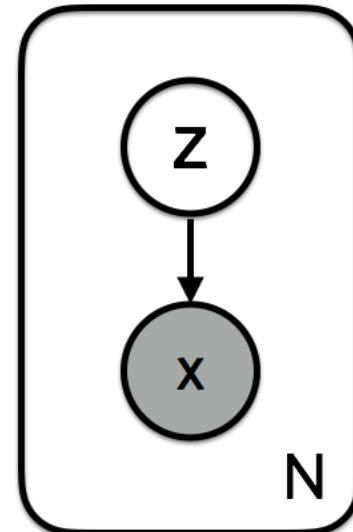
VAE: Probability view



- Posterior probability (encoder)

$$p_{\phi}(z|x) = \frac{p_{\phi}(x|z)p(z)}{p_{\phi}(x)}$$

where $p_{\phi}(x) = \underbrace{\int p_{\phi}(x, z) dz}_{\text{Difficult to compute}}$



$$p_{\phi}(z|x) \longrightarrow q_{\theta}(z|x) \text{ Such that } \mathbb{D}_{\text{KL}}(q_{\theta}(z|x) \| p_{\phi}(z|x)) \text{ is minimized}$$

Approximated with

VAE: Variational approximation



$$\log p_{\phi}(\mathbf{x}) = \underbrace{\mathbb{E}_{q_{\theta}(z|x)} \left[\log \frac{p_{\phi}(\mathbf{x}, z)}{q_{\theta}(z|x)} \right]}_{\text{ELBO}} + \mathbb{D}_{\text{KL}}(q_{\theta}(z|x) \| p_{\phi}(z|x)).$$

- Maximize ELBO == Minimize DL

VAE: Variational approximation (cont.)



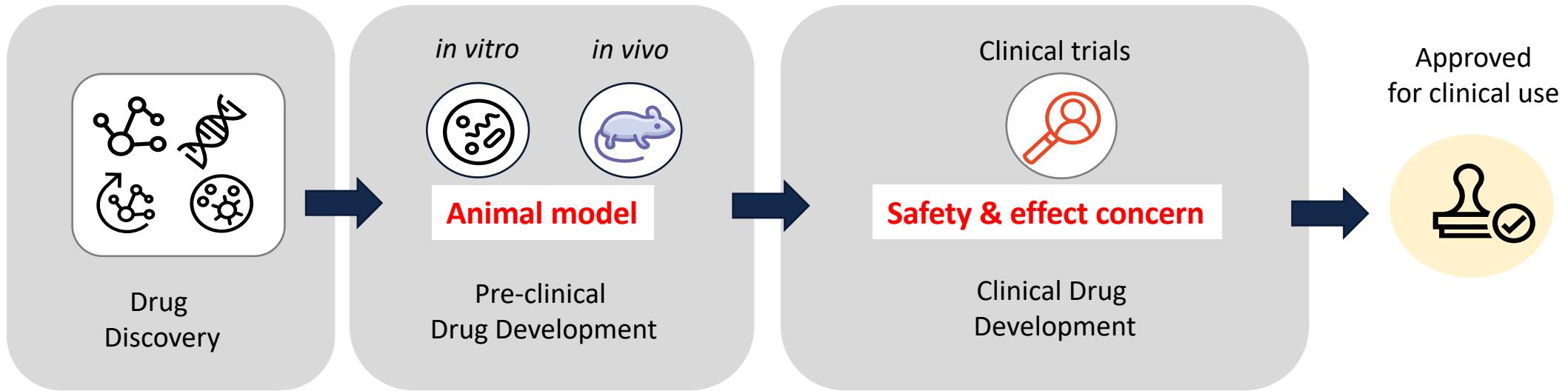
$$\text{ELBO}(\mathbf{x}; \boldsymbol{\phi}, \boldsymbol{\theta}) = \mathbb{E}_{q_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x})}[\log p_{\boldsymbol{\phi}}(\mathbf{x}|\mathbf{z})] - \mathbb{D}_{\text{KL}}(q_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x}) \| p_{\boldsymbol{\phi}}(\mathbf{z}))$$



$$l_{\mathbf{x}} = -\mathbb{E}_{\mathbf{x} \sim D} [\log p(\mathbf{x}|\mathbf{z}) + \mathbb{D}_{\text{KL}}(q_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x}) \| p_{\boldsymbol{\phi}}(\mathbf{z}))]$$

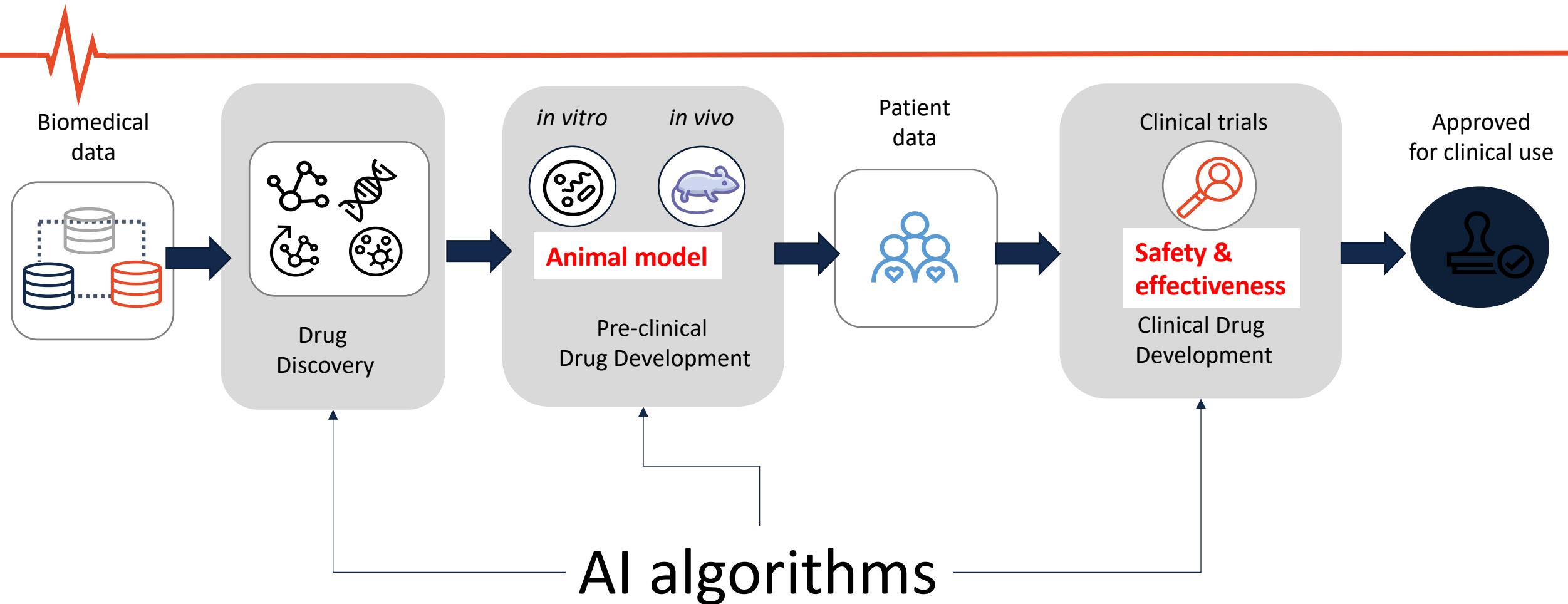
Deep learning for Drug discovery

Traditional Drug Discovery & Development Process



	Drug discovery	Pre-clinical	Phase 1	Phase 2	Phase 3
Time spent	4-5 years	1-2 years	1-2 years	1-2 years	2-3 years
\$ spent	\$550M	\$125M	\$225M	\$250M	\$250M
Output	5,000 - 10,000 compounds	10-20 candidates	5-10 candidates	2-5 candidates	1-2 candidates

AI/ML to the Rescue: Why and How?



Molecule generation with VAE

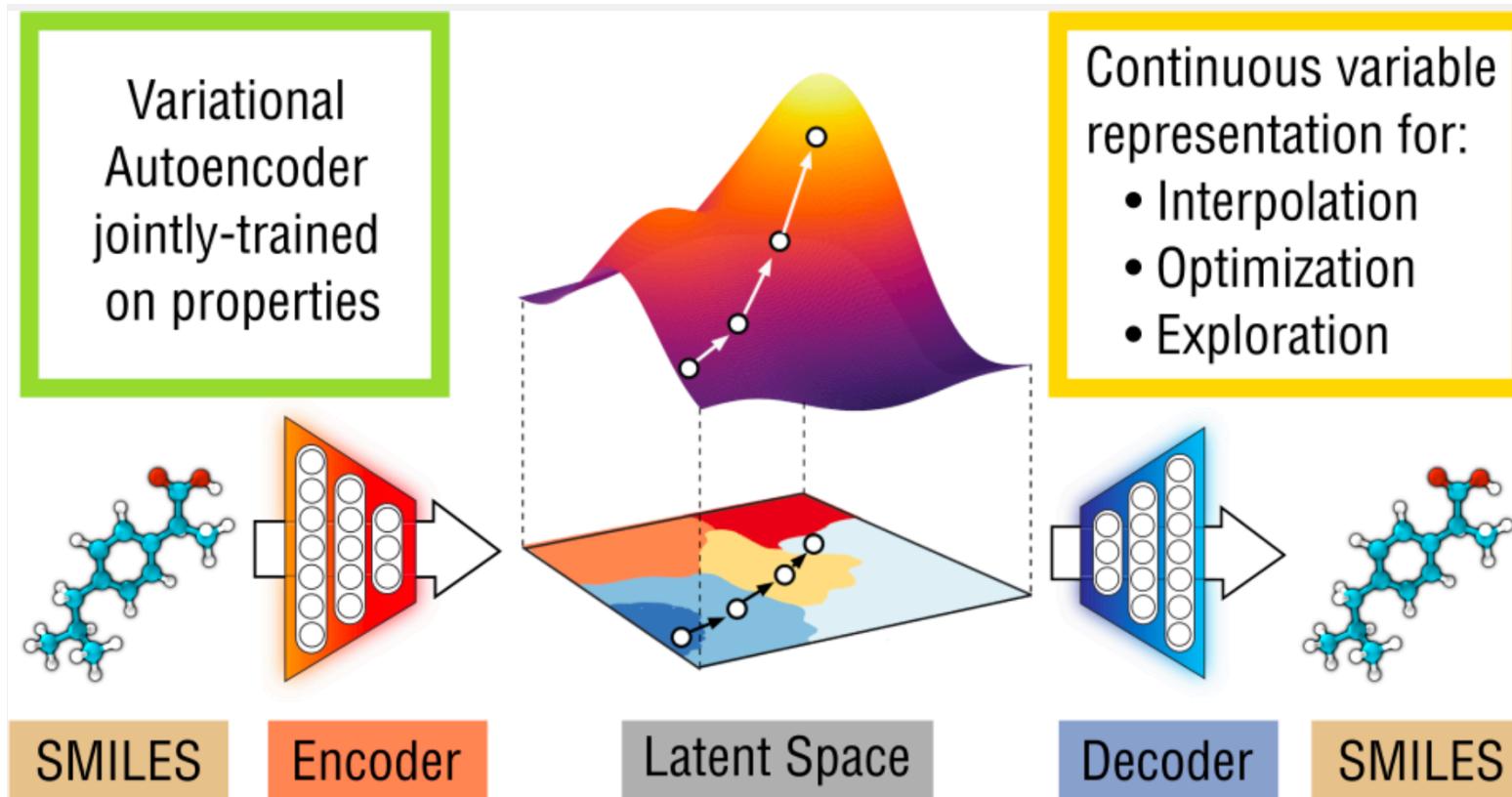
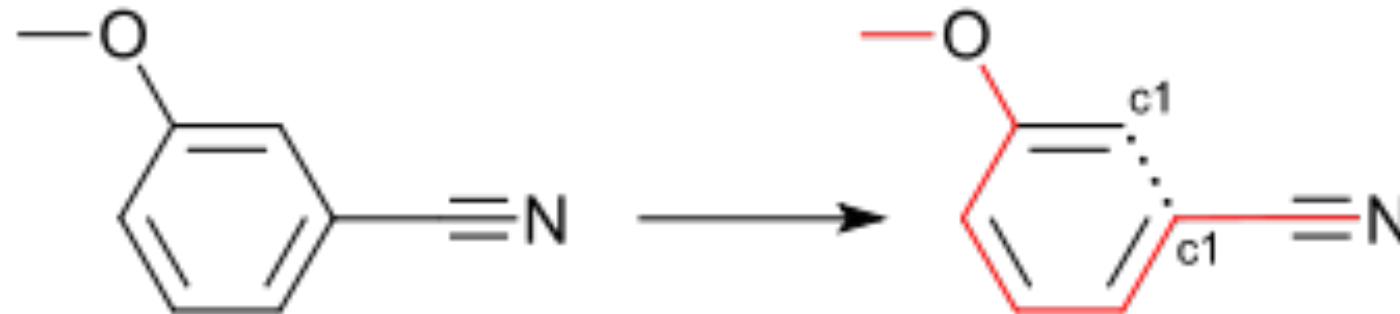


Image credit and citation: Gómez-Bombarelli et al. 2018. "Automatic Chemical Design Using a Data-Driven Continuous Representation of Molecules." *ACS Central Science* 4 (2): 268–76.

Simplified Molecular-Input Line-Entry System (SMILES)



3-cyanoanisole



() are used to branches

COc(c1)cccc1C#N

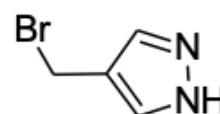
SMILES

Numbers are used to
represent rings

- SMILES string: Traverse the molecular graph in a depth-first manner following the atom with the smallest label at each branch point.

Molecular Representation Learning

Graph



SMILES

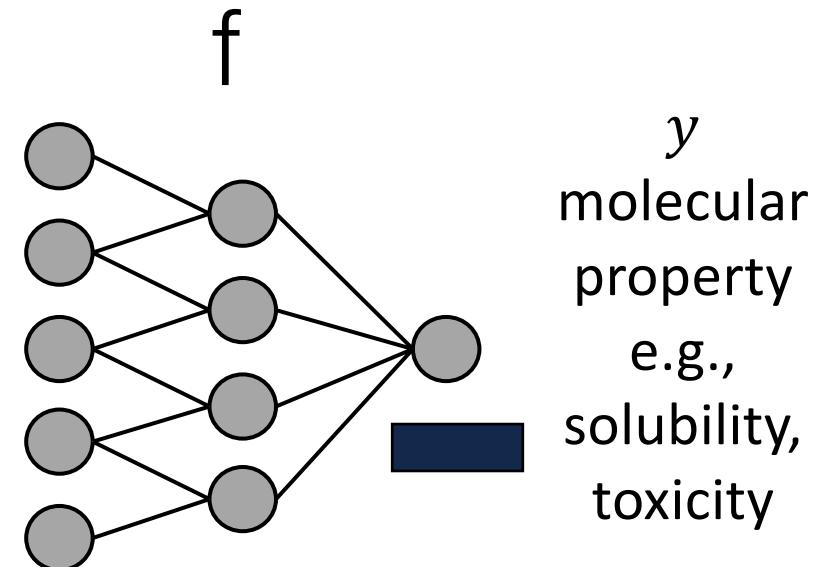
BrCc1c[nH]nc1

One-hot
encoding

Br C c 1 c [n H] n c 1

Br	1	0	0	0	0	0	0	0	0	0	0	0
C	0	1	0	0	0	0	0	0	0	0	0	0
H	0	0	0	0	0	0	0	1	0	0	0	0
c	0	0	1	0	1	0	0	0	0	0	1	0
n	0	0	0	0	0	0	1	0	0	1	0	0
1	0	0	0	1	0	0	0	0	0	0	0	1
[0	0	0	0	0	1	0	0	0	0	0	0
]	0	0	0	0	0	0	0	0	1	0	0	0

vectorize



VAE for molecule generation

- SMILES strings as input
- Seq2Seq model with RNN is used to encode SMILES
- VAE is used to generate new SMILES
- Another property prediction network is also added

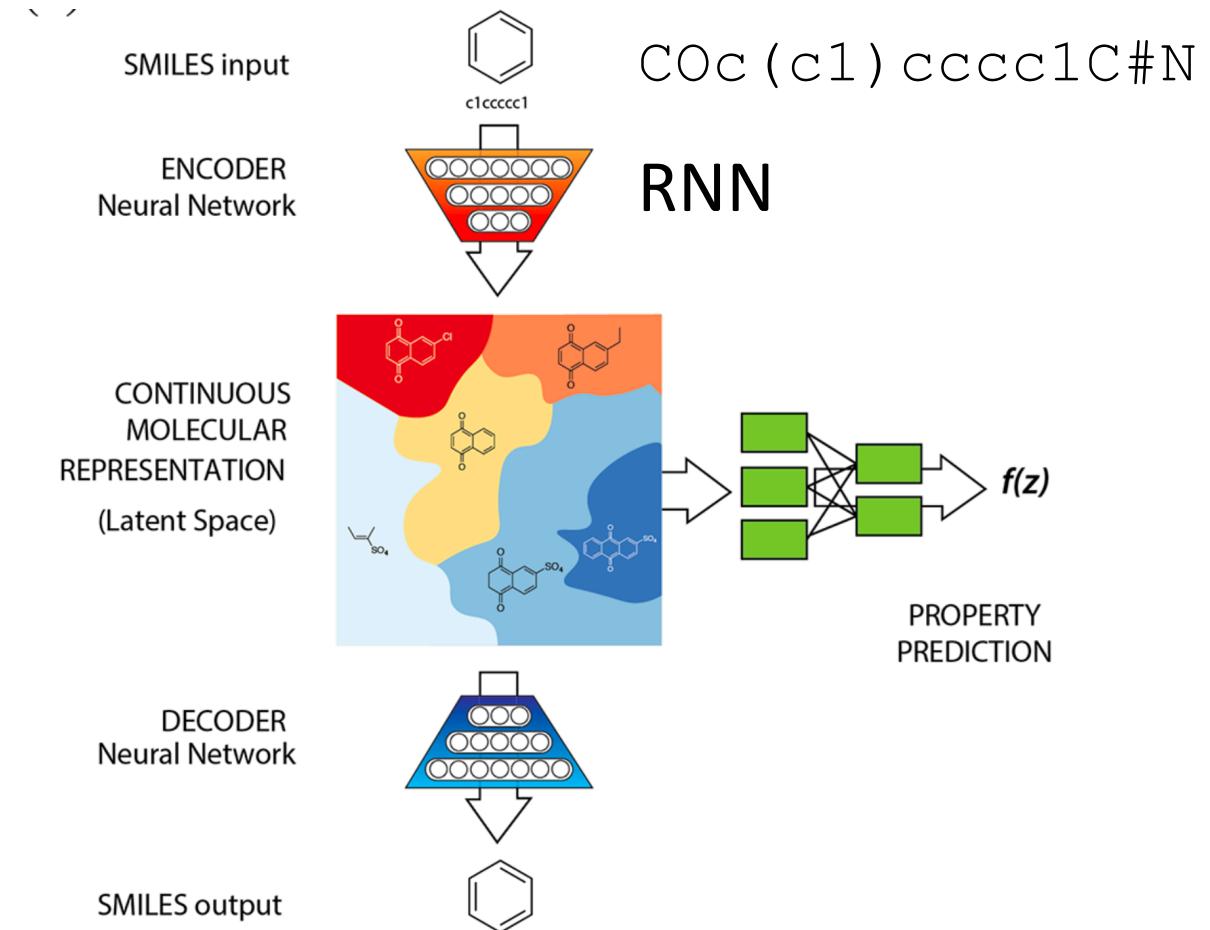


Image credit: Gómez-Bombarelli et al. 2018. "Automatic Chemical Design Using a Data-Driven Continuous Representation of Molecules." *ACS Central Science* 4 (2): 268–76.

Datasets

- 108K molecules from the QM9 data set of molecules
- 250K drug-like molecules in ZINC database

Results



source ^a	data set ^b	samples ^c	logP ^d	SAS ^e	QED ^f
Data	ZINC	249k	2.46 (1.43)	3.05 (0.83)	0.73 (0.14)
GA	ZINC	5303	2.84 (1.86)	3.80 (1.01)	0.57 (0.20)
VAE	ZINC	8728	2.67 (1.46)	3.18 (0.86)	0.70 (0.14)
Data	QM9	134k	0.30 (1.00)	4.25 (0.94)	0.48 (0.07)
GA	QM9	5470	0.96 (1.53)	4.47 (1.01)	0.53 (0.13)
VAE	QM9	2839	0.30 (0.97)	4.34 (0.98)	0.47 (0.08)

water-octanol partition coefficient (logP)
synthetic accessibility score
Drug-likeness

- GA is another method using genetic programming
- Goal is to generate realistic molecules that are similar to the input data