



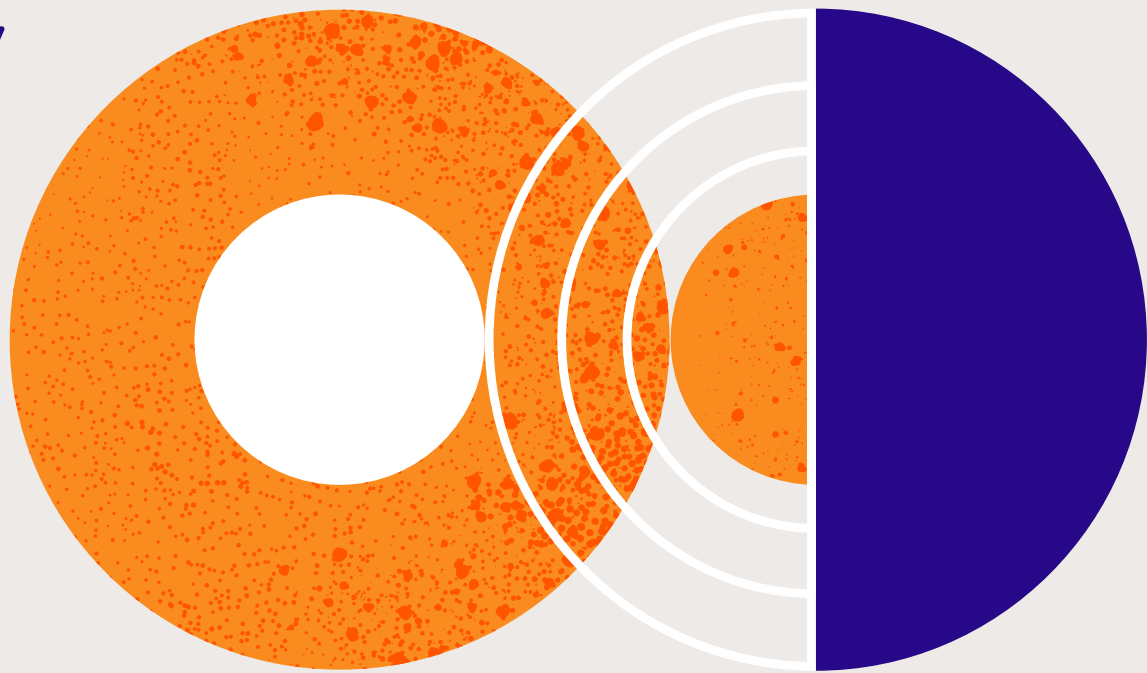
**UGM** Boston  
Oct. 4th 23  
 Chemaxon

# Generative AI Begins to Dominate the AI Conversation in Early Drug Discovery

**Joe Michel**  
Director of Informatics,  
Cytokinetics

# Generative AI Begins to Dominate the AI Conversation in Early Drug Discovery

Joe Michel,  
Head of Informatics,  
Cytokinetics Inc.  
South San Francisco, CA

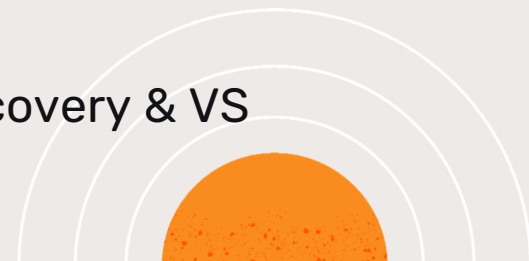


# New Focus of AI: Generative AI

How Generative AI shifts our conversation from predictors to generating medicines

## Agenda:

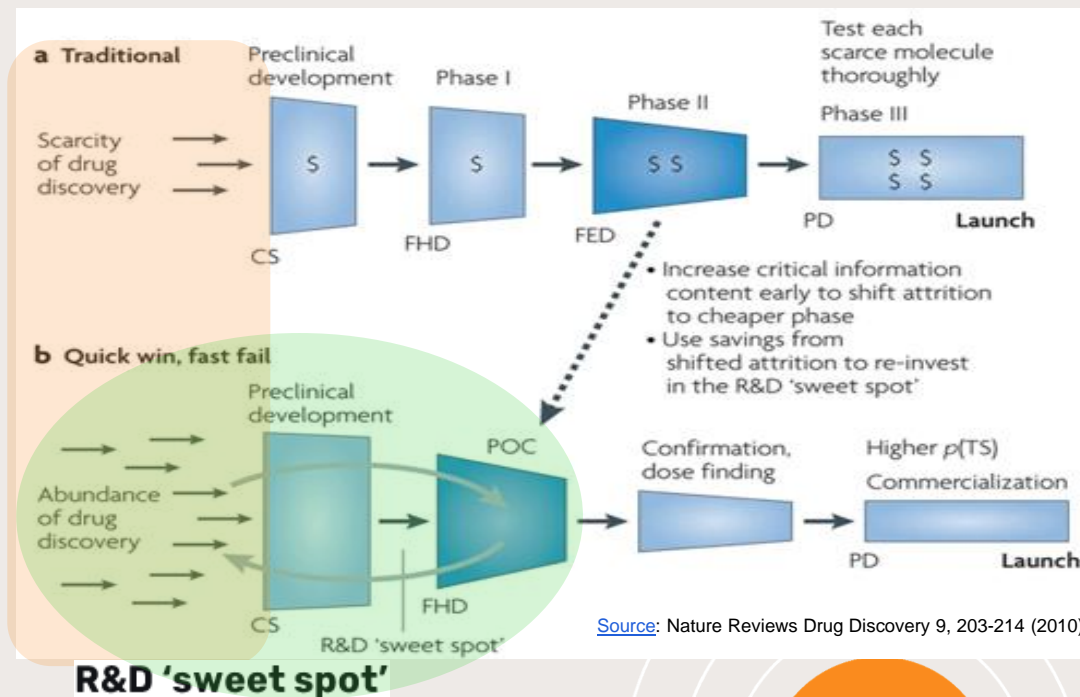
- Need for advancements in Research Drug Discovery
- Generative AI has made impact
- Recent Headlines - Generative AI (Gen AI) continues contributes
- What has changed? What is Generative AI / Foundation models?
- Generative AI Framework
- Generative AI Examples in Drug Discovery Ligand Discovery & VS



# New Focus of AI: Generative AI

What is the need? Cost, Time, and Quality...

Quick win, fast fail  
drug development  
paradigm

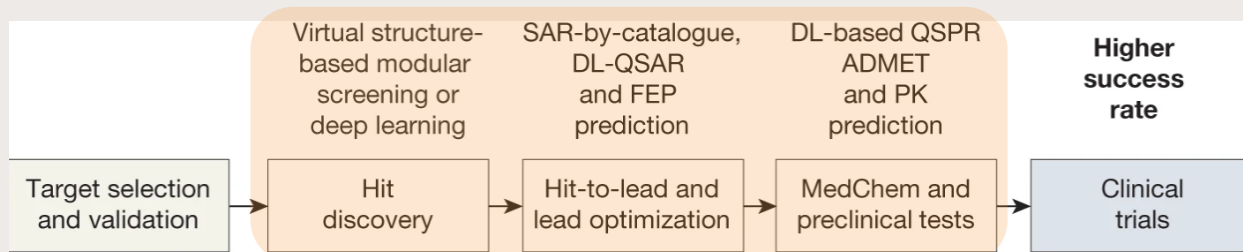


# New Focus of AI: Generative AI

One of the areas AI technology could make an impact

- Computer-aided drug discovery (CADD) been around for decades, as well as Ligand Discovery Technologies.
- How can Ligand Discovery Technologies improve “quick win, fast fail”?

## Ligand Discovery Technologies



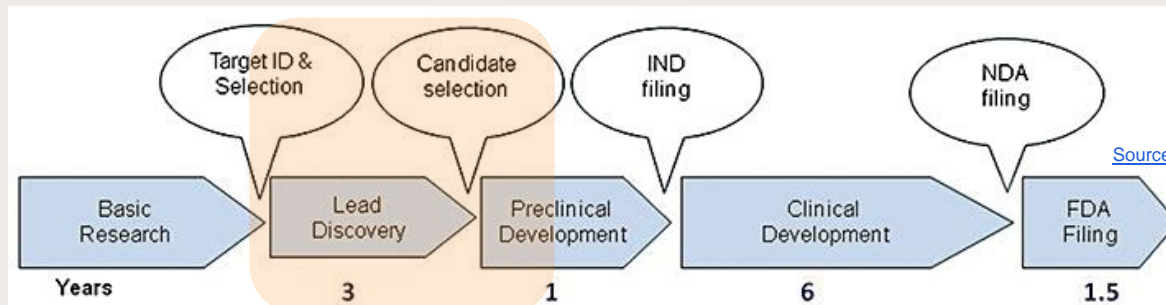
[Source](#): Nature 616, 673-685 (2023) , ISSN 1476-4687



# New Focus of AI: Generative AI

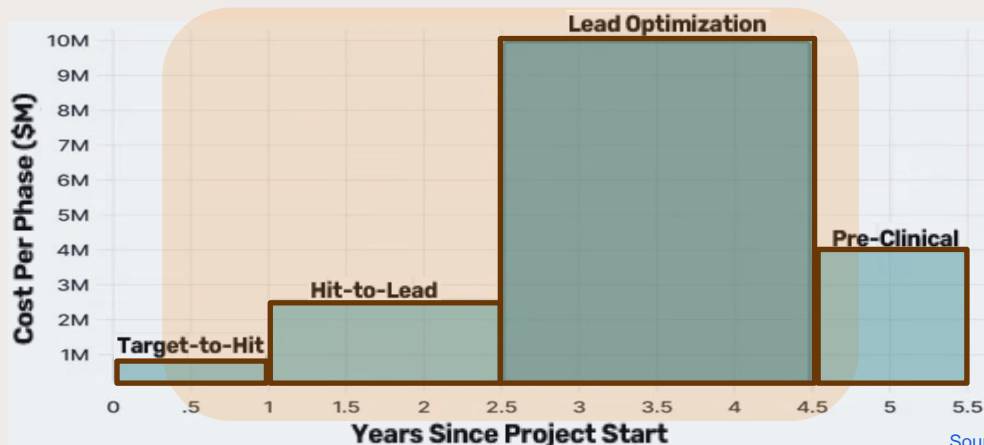
One of the areas AI technology could make an impact

**Typical Years of Drug Discovery Phases:**



Source: Br J Pharmacol 2011; 162(6): 1239-1249

**Within Discovery Phases, Typically, Lead Optimization Phase is the most costly:**



Source: Nature Reviews Drug Discovery 9, 203-214 (2010)

# Generative AI has made “a breakthrough”

## It's here, it's real and starting to impact Drug Discovery

- In 2023, ChatGPT (a chatbot developed by Microsoft's OpenAI) was “born”
- DeepMind (AlphaFold) CEO, Demis Hassabis, [spoke on this topic](#), **AlphaFold IS A major break-through** Understanding and controlling protein folding is arguably the most important challenge in structural biology ([Martnez, 2014](#); [Jankovic and Polovic, 2017](#))

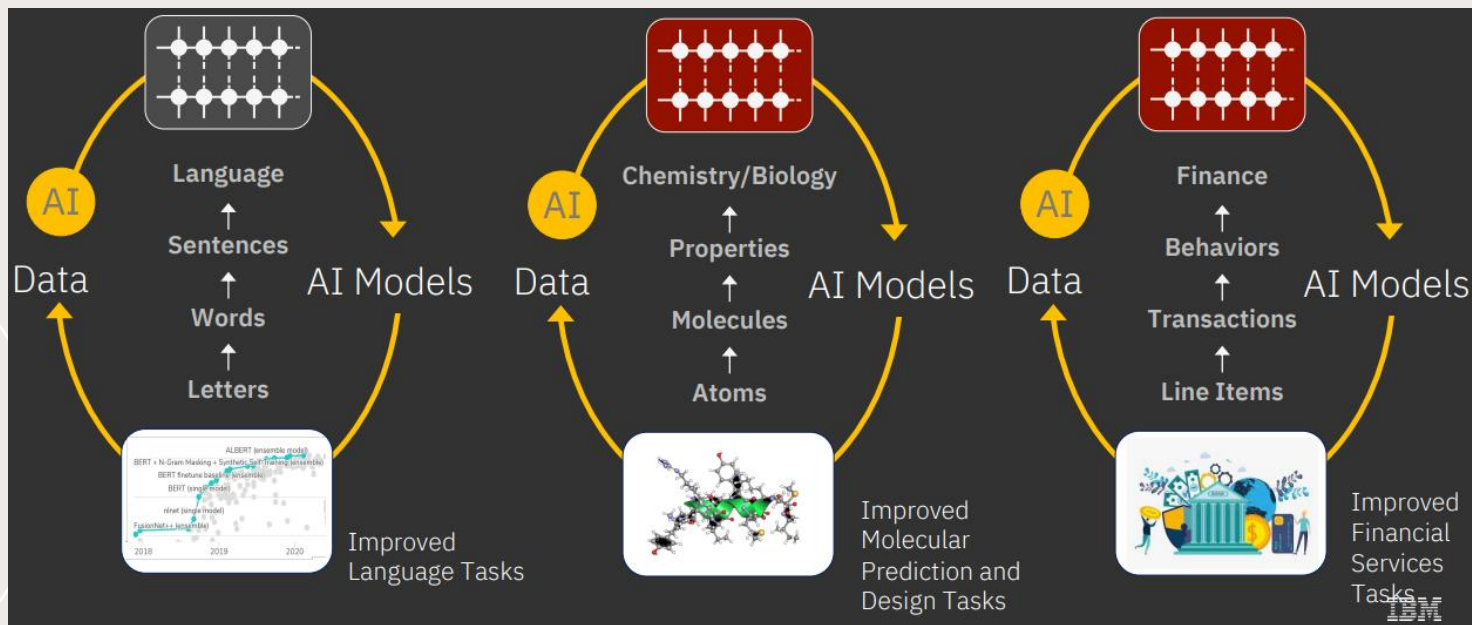


- Nvidia's CEO Jensen Huang stated, “**We are in the 'iPhone moment' of AI**.” “2023 is expected to be the most exciting year in the field of AI” ([source](#): keynote speech 2023)



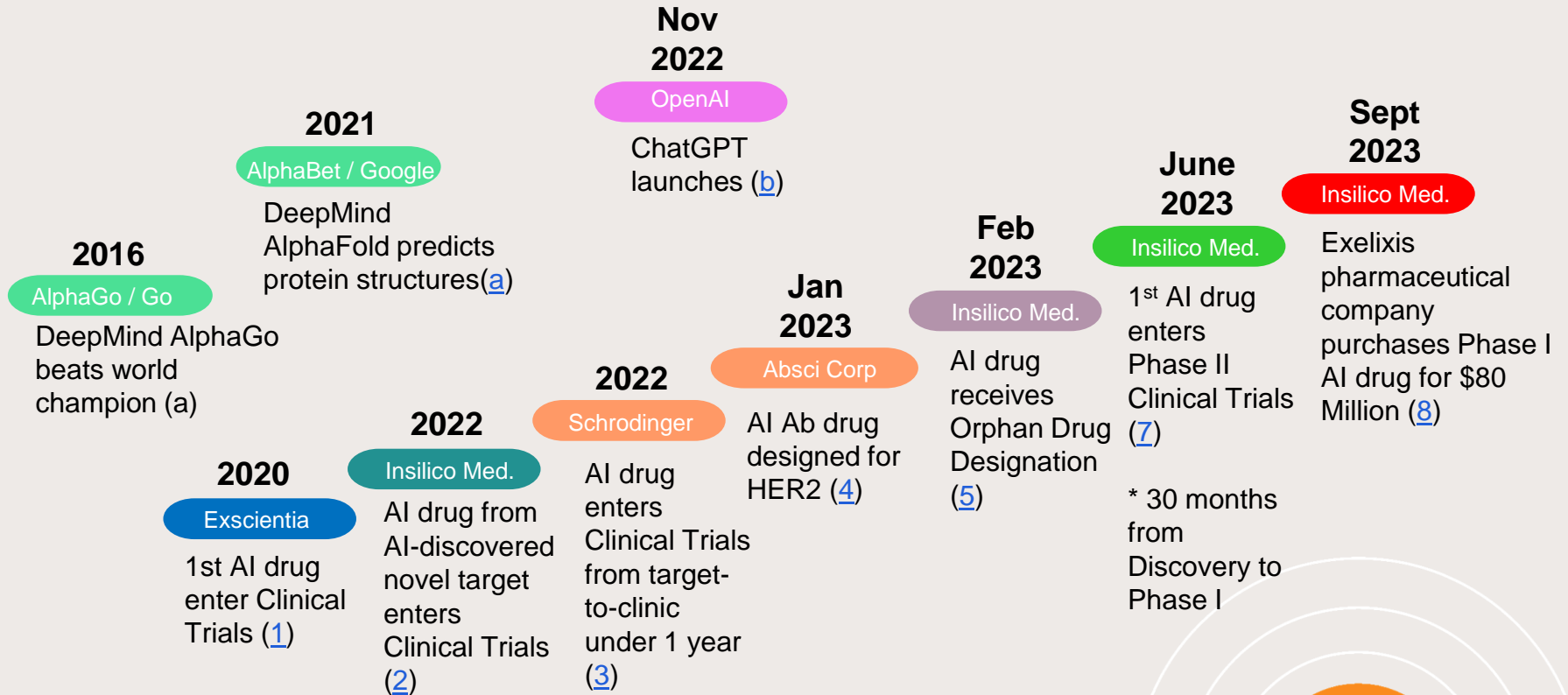
# Generative AI continues to “breakthrough”

The same AI breakthroughs happening in language are impacting other scientific and enterprise applications





# Generative AI impact in recent years



# Gen AI – Recent Headlines

## AI Year in Review (2023)

**Jan 2023:** Absci - deploys zero-shot generative AI in antibody design

**March 2023:** MIT reveals DiffDock, which could support faster, safer drug development

**March 2023:** Nvidia launches BioNeMo Cloud to boost drug discovery with generative AI

**June 2023:** Insilico Medicine's AI drug enters phase 2 study

**June 2023:** Sanofi reveals plan to put AI at the center of its operations

**July 2023:** AI-aided drug ulotaront fails phase 3 studies

**July 2023:** Nvidia invests \$50 million in Recursion to boost AI-driven drug discovery

**July 2023:** AI startup Causaly raises \$60 million to expedite drug development

**August 2023:** Pharos iBio develops anticancer drug with AI platform

**August 2023:** Generative AI tool boasts 79% accuracy in predicting clinical trial outcomes

**August 2023:** UCF researchers reveal AI-assisted drug screening technology

**August 2023:** Recursion uses AI to bridge protein/chemical universe, predicted targets for 36 bill. Compounds  
Recursion is using Cyclica's MatchMaker technology, NVIDIA DGX Cloud supercomputing and DeepMind's [AlphaFold2 database](#) to screen Enamine REAL Space



# AI (DL) Breakthroughs over the decades

**Backprop**  
Rise of  
Neural Networks

**1960's .  
1990's**

**ReLU  
(rectified linear unit)**  
Activation Function  
breakthrough

**2000's**

**CNN Improvements**  
e.g. ReLu w/ Dropout

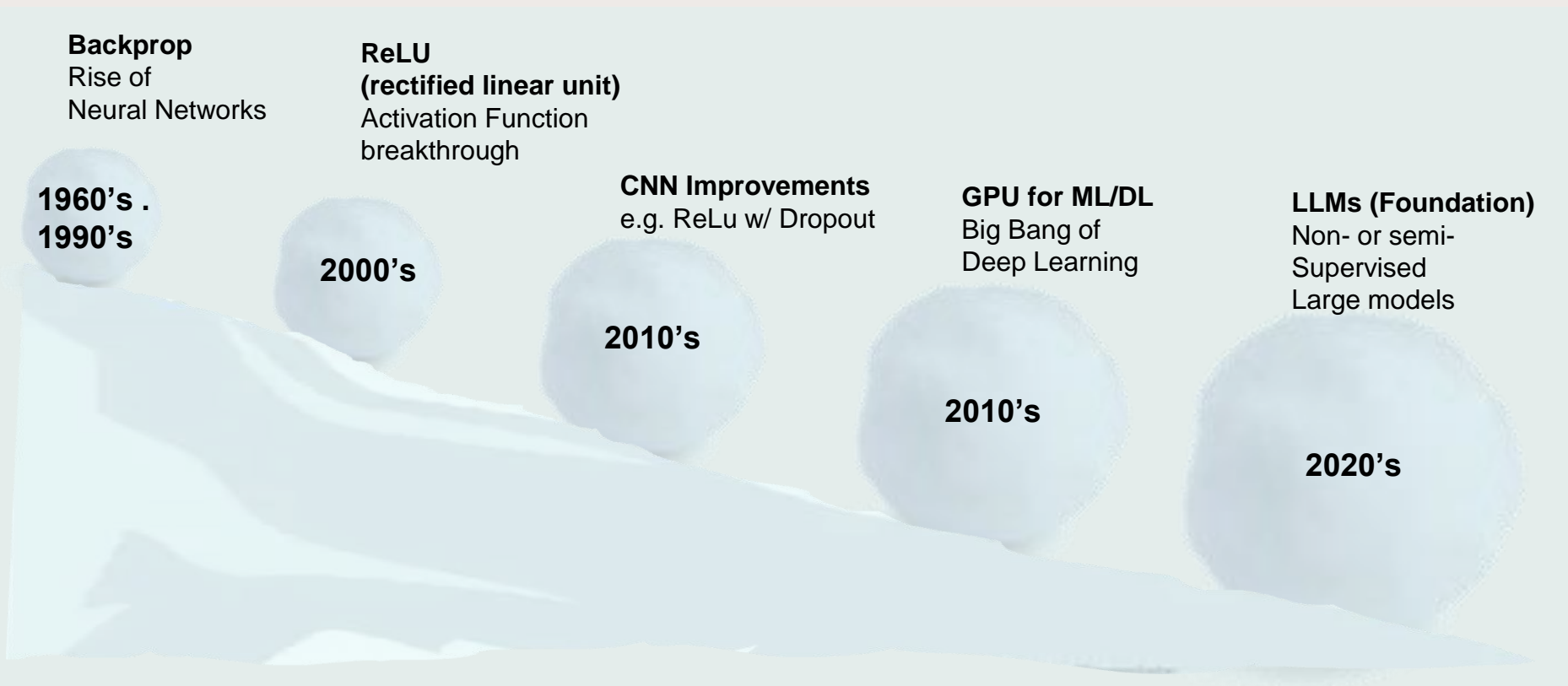
**2010's**

**GPU for ML/DL**  
Big Bang of  
Deep Learning

**2010's**

**LLMs (Foundation)**  
Non- or semi-  
Supervised  
Large models

**2020's**



# What's required for Gen AI?

LLMs or Foundation models is a fundamental AI “breakthrough”

- ✓ **Pre-trained** on unlabeled **large** datasets of different modalities (e.g. molecules, proteins)
- ✓ Leverage **self-supervised learning**
- ✓ Learn **generalizable & adaptable data representations** which can be effectively used in multiple downstream tasks (e.g. molecule generation)

e.g.

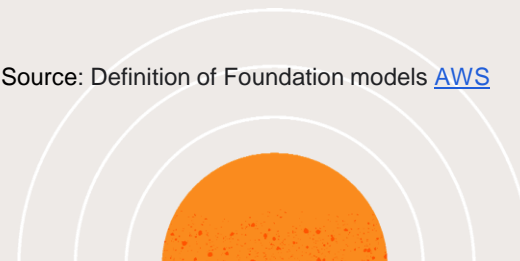
GPT-3 (ChatGPT) - 175 billion parameters  
IBM MolFormer-XL – 1.1 billion molecules

[Source](#): Nature Machine Intelligence 4, 1256-1264



Source: Definition of Foundation models [AWS](#)

More input records, more potential input parameters...



# What's definition of Generative AI?

## It makes something, not just a predictor...

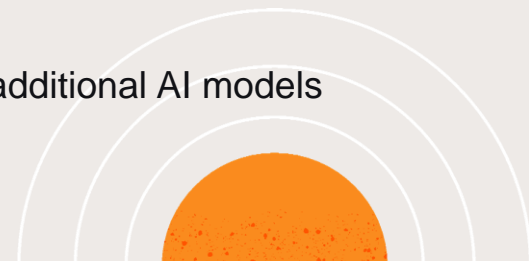
Generally Speaking:

- Generative AI is artificial intelligence capable of generating text, images, other media, using generative models. It produces beyond just a mathematical predictor value.

For Ligand Discovery Technologies:

- Generative AI is trained to understand molecules and can generate new molecules randomly or can be controlled by specific criteria of interest (**molecule generator**).
  - Represented Molecular Space – (VAEs or GANs) Gen AI that can produce a representation of molecules based on an input molecule or information (approx. Gaussian distribution)
  - Reinforcement Learning - This control mechanism usual requires additional AI models separate from the molecule generator

[Source](#): What is Generative AI and How does it work?

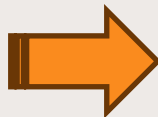


# Another “breakthrough” in DD?

## Producing beyond predictions or classifications

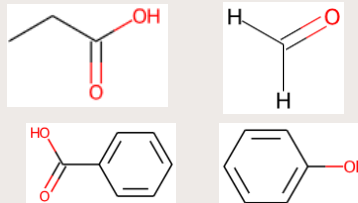
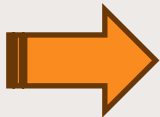
DALL-E (GPT-3)

**Text Description:**  
Chemist synthesizing molecule  
in Biotech lab  
in Boston Massachusetts



Molecule Gen AI

**Text Description:**  
SMILE String  
[CH3][CH2]C(O)=O  
&  
Potentially other ligand info



Example Tools:

Application: Chemistry42 by Insilico Med

Software Tool Kit: IBM GT4SD by IBM (open source [release March 2023](#))

# Nuts & Bolts - 3 popular Gen AI model types

- Breakthrough – Foundation models -> Ability to leverage different learning approaches, including unsupervised or semi-supervised learning
- 3 Main Deep Learning Models used for Generative AI:
  - Denoising **Diffusion** Probabilistic **models** (DDPMs) – (also considered Foundation models)  
Two-step process during training. The two steps are forward diffusion and reverse diffusion
  - Variational autoencoders (VAEs)  
Collaboration of two neural networks (encoder and decoder), with a subsampling middle layer
  - Generative adversarial networks (GANs)  
Pitting two neural networks against each other: a generator that generates new examples and a discriminator



# How can Gen AI help discover Ligands

## Guided molecule generation is required

Power comes in combination of Deep Learning Techniques (supervised & unsupervised):

- BASIS: VAEs & GANs capable of generating molecule set based on representative input (e.g. IBM)
- Predictors – use or train needed predictors (bread and butter)
- Additionally apply Reinforcement Learning (RL) – Can be combined with VAE or GAN to further optimize generative molecular structure & improve chances of success

[Source:](#) What is Generative AI and How does it work?

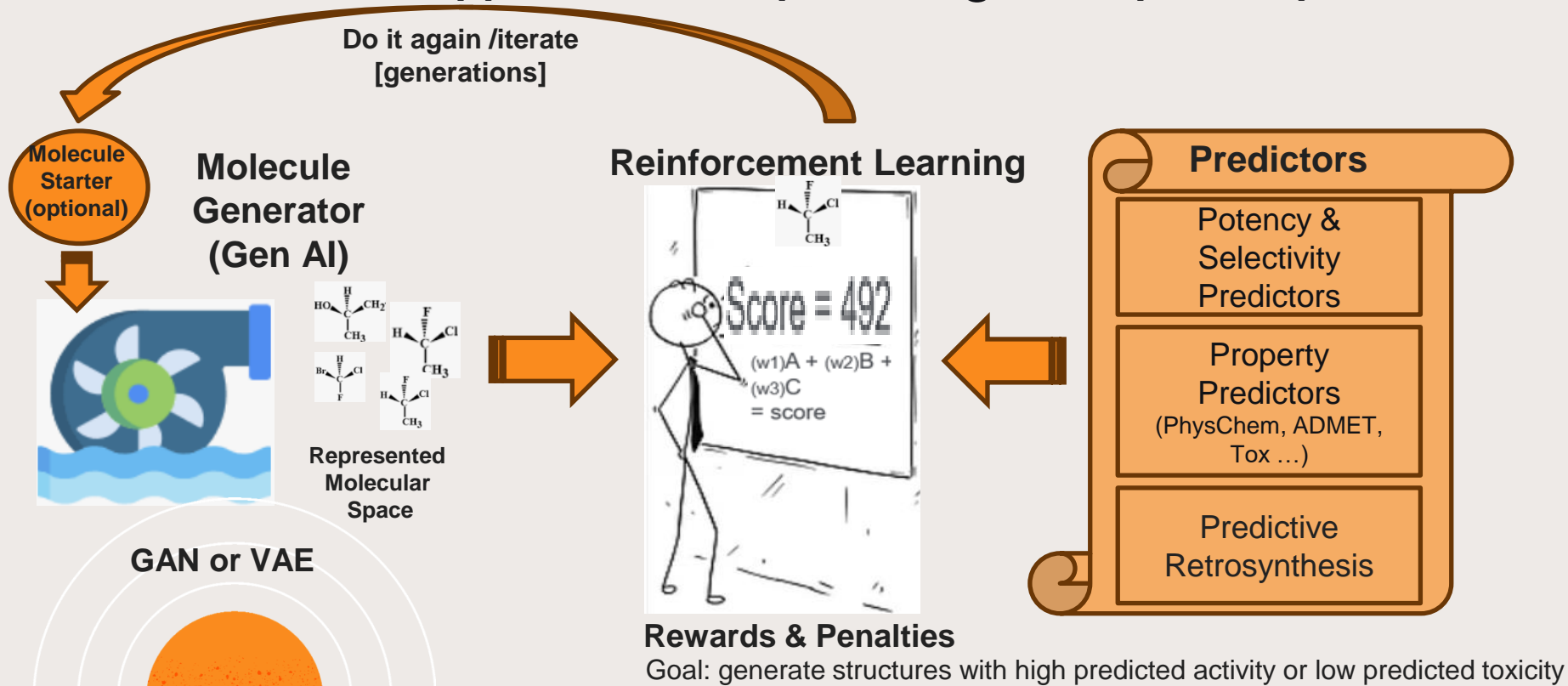
[Source:](#) Saturn Cloud Generative AI explanation





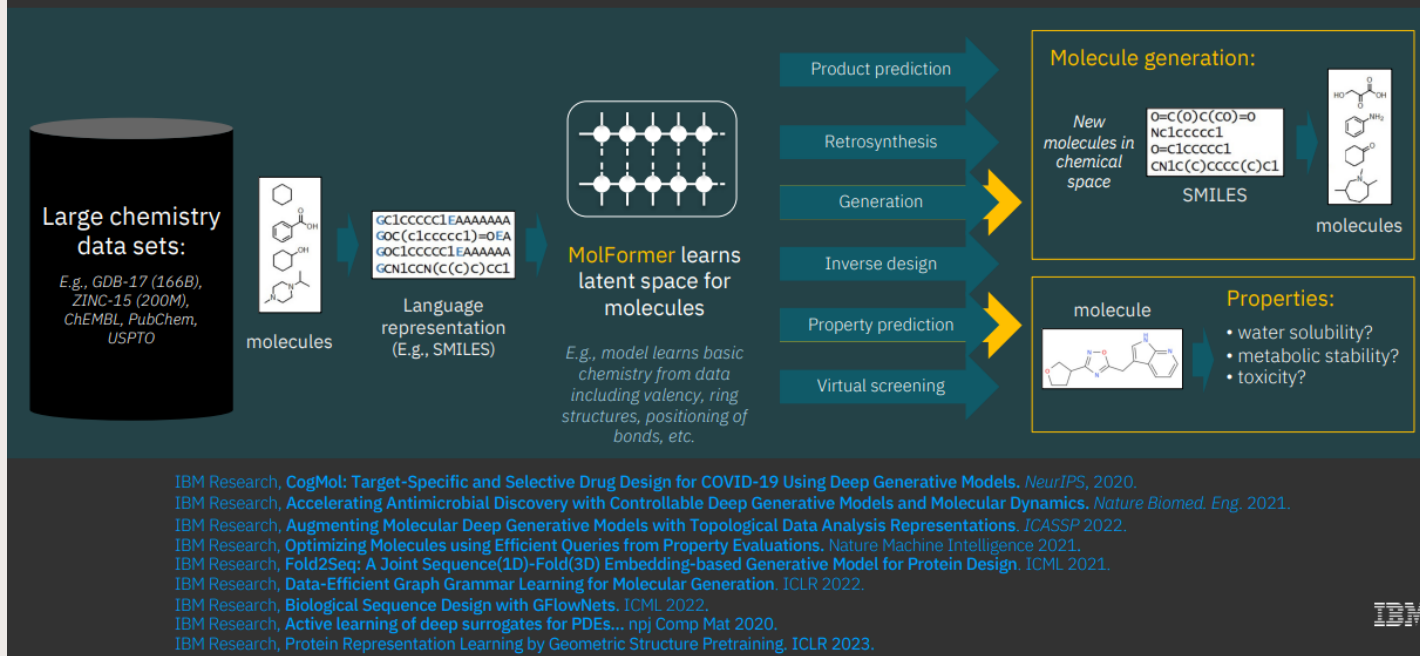
# Gen AI Framework

Several model types can help strengthen your system



# IBM Generative AI / Foundation AI

Foundation Models learn the language of chemistry/biology from data and can power up a multitude of discovery tasks – We call them MolFormer



IBM Research, **CogMol: Target-Specific and Selective Drug Design for COVID-19 Using Deep Generative Models**. *NeurIPS*, 2020.

IBM Research, **Accelerating Antimicrobial Discovery with Controllable Deep Generative Models and Molecular Dynamics**. *Nature Biomed. Eng.* 2021.

IBM Research, **Augmenting Molecular Deep Generative Models with Topological Data Analysis Representations**. *ICASSP 2022*.

IBM Research, **Optimizing Molecules using Efficient Queries from Property Evaluations**. *Nature Machine Intelligence* 2021.

IBM Research, **Fold2Seq: A Joint Sequence(1D)-Fold(3D) Embedding-based Generative Model for Protein Design**. *ICML 2021*.

IBM Research, **Data-Efficient Graph Grammar Learning for Molecular Generation**. *ICLR 2022*.

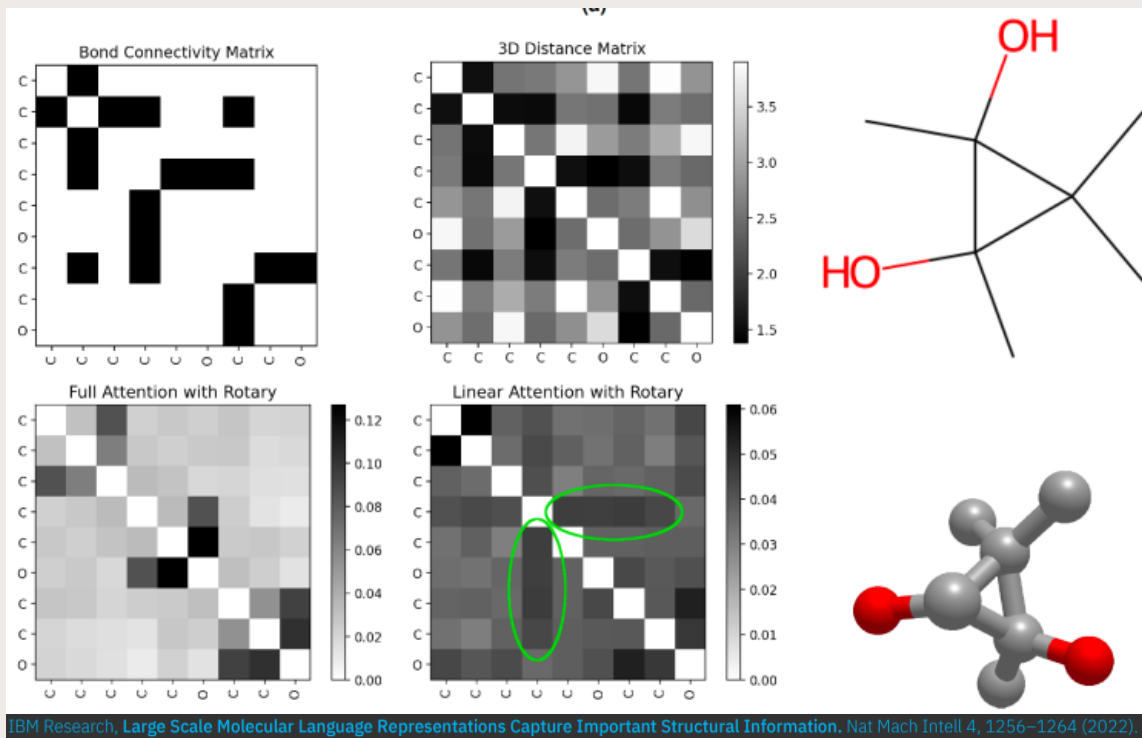
IBM Research, **Biological Sequence Design with GFlowNets**. *ICML 2022*.

IBM Research, **Active learning of deep surrogates for PDEs...** *npj Comp Mat* 2020.

IBM Research, **Protein Representation Learning by Geometric Structure Pretraining**. *ICLR 2023*.

IBM

# IBM Molformer captures structural info



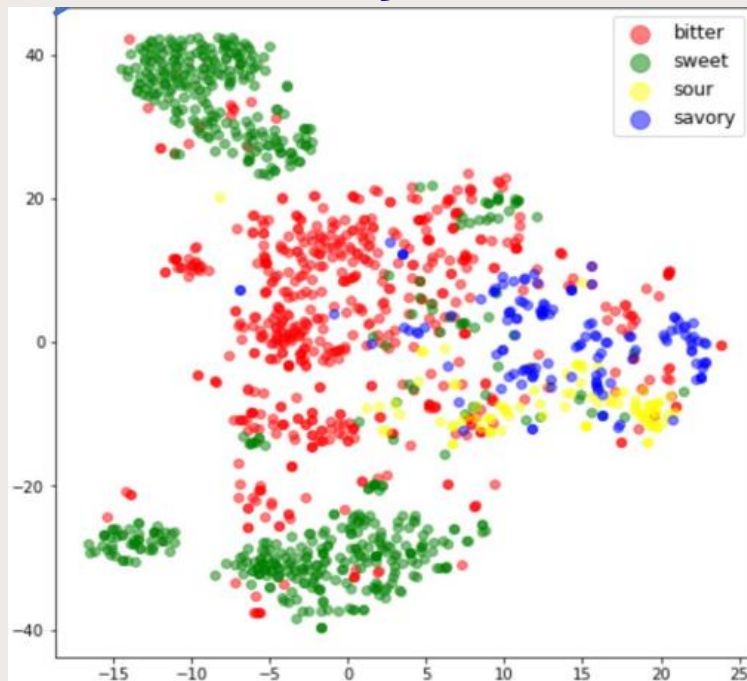
[Source:](#) The Promise of Foundation Models and Generative AI (2023), Payel Das, Principal Research & Master Investor

# Molformer Learns molecular taste w/o labels (unsupervised)

Visualization of unsupervised MolFormer Embeddings in t-SNE space and separation of flavor molecules in that space.

[source](#)

t-SNE (t-distributed Stochastic Neighbor Embedding) is an unsupervised non-linear dimensionality reduction technique for data exploration and visualizing high-dimensional data



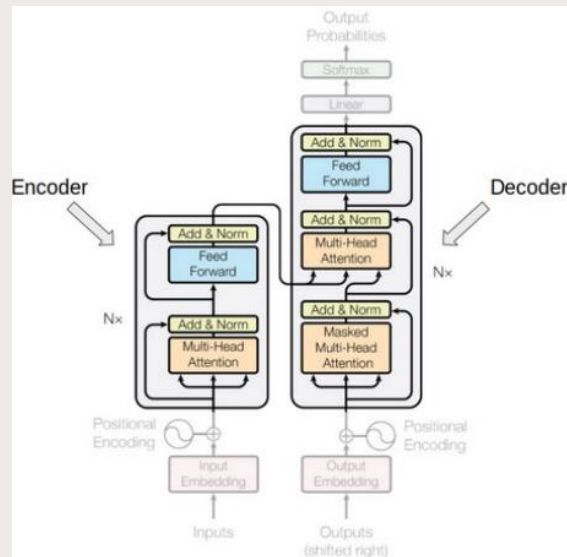
IBM Research, [Cloud-Based Real-Time Molecular Screening Platform with MolFormer](https://2022.ecmlpkdd.org/wp-content/uploads/2022/09/sub_1455.pdf), ECML PKDD (2022). [https://2022.ecmlpkdd.org/wp-content/uploads/2022/09/sub\\_1455.pdf](https://2022.ecmlpkdd.org/wp-content/uploads/2022/09/sub_1455.pdf)

[Source](#): The Promise of Foundation Models and Generative AI (2023), Payel Das, Principal Research & Master Investor

# MolFormer: Foundational transformer

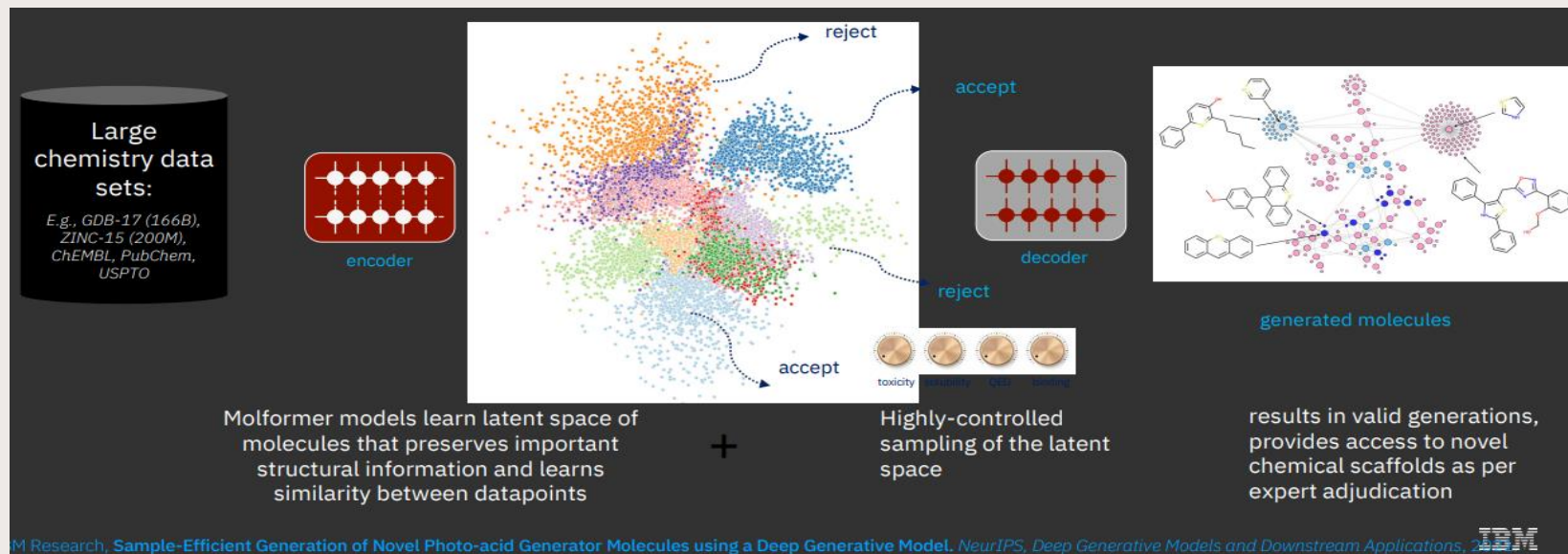
MolFormer-XL – a specific example from MolFormer family

- Trained on up to over **a billion** molecular text strings (SMILES), with relatively limited hardware resources (16 V100 GPUs).
- Scalable and fast to train linear time attention transformers as encoders and decoders
- Relative position embeddings facilitate learning on SMILES
- State-of-the-art, universal chemical language model for wide ranges of **70+ molecular property prediction**
- Shows **emergent behavior**, such as geometry, taste, etc.



# Controllable Generation of Novel Molecules

Large-scale unsupervised pretraining, novel sampling, and optimization methods enable controllable generation of novel molecules w/ desired properties



M Research, [Sample-Efficient Generation of Novel Photo-acid Generator Molecules using a Deep Generative Model](#), *NeurIPS, Deep Generative Models and Downstream Applications*, 2023

# Comparison of MolFormer w/ others

Comparison with existing baselines on classification and regression benchmarks

Dataset	BBBP	Tox21	ClinTox	HIV	BACE	SIDER
Tasks	1	12	2	1	1	27
RF	71.4	76.9	71.3	78.1	<b>86.7</b>	<b>68.4</b>
SVM	72.9	<b>81.8</b>	66.9	<b>79.2</b>	86.2	68.2
MGCN [56]	<b>85.0</b>	70.7	63.4	73.8	73.4	55.2
D-MPNN [57]	71.2	68.9	<b>90.5</b>	75.0	85.3	63.2
Hu, et al. [58]	70.8	78.7	78.9	80.2	85.9	65.2
N-Gram [44]	91.2	76.9	85.5	<b>83.0</b>	87.6	63.2
MolCLR [24]	73.6	79.8	93.2	80.6	<b>89.0</b>	68.0
<b>MoLFoRmER-XL</b>	<b>93.7</b>	<b>84.7</b>	<b>94.8</b>	82.2	88.21	<b>69.0</b>



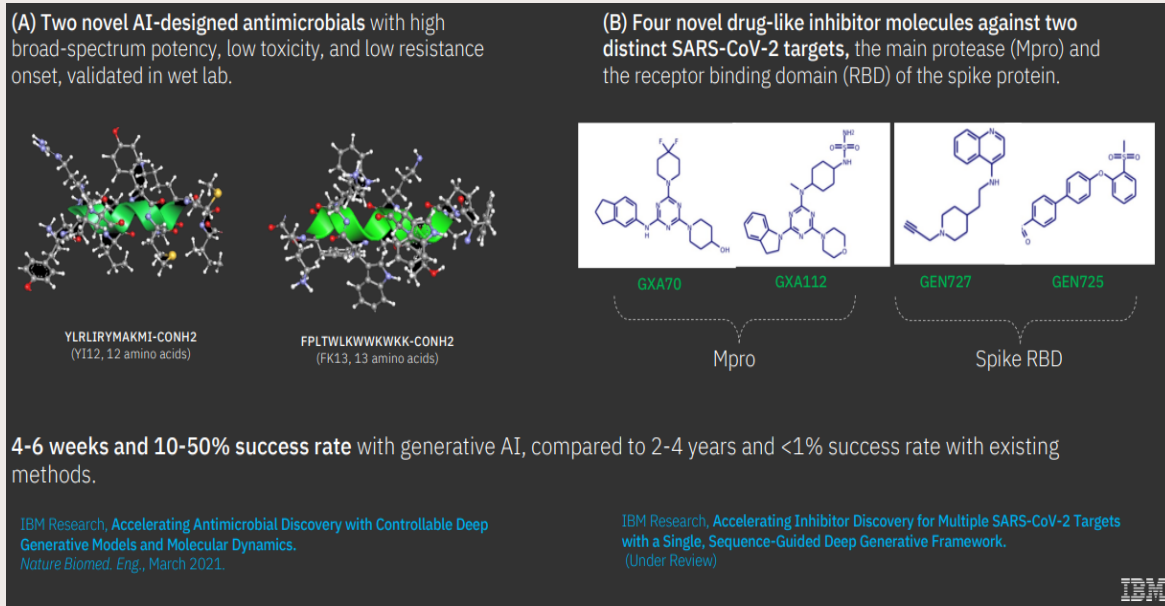
# IBM Utilizes LLMs & Gen AI

## SARS-CoV-2 potential molecules discovered

IBM research created a deep generative framework, CogMol, to design small molecule inhibitors for two different targets –the spike protein receptor binding domain (RBD) and the main protease from SARS-CoV-2

IBM synthesized 4 potential inhibitors for Mpro

2 of them had inhibitory activity (43 and 34.2  $\mu\text{M}$ )



[Source](#): Four novel SARS-CoV-2 molecules



# IBM Utilizes LLMs & Gen AI

## SARS-CoV-2 potential molecules discovered

VAE (encoder / decoder) - molecules

Pre-trained embeddings of protein

UC San Diego  
SKAGGS SCHOOL OF PHARMACY  
AND PHARMACEUTICAL SCIENCES

BindingDB database

Gen AI  
Used Conditional Latent Space Sampling (CLaSS)  
for attribute-controlled

IBM RXN platform

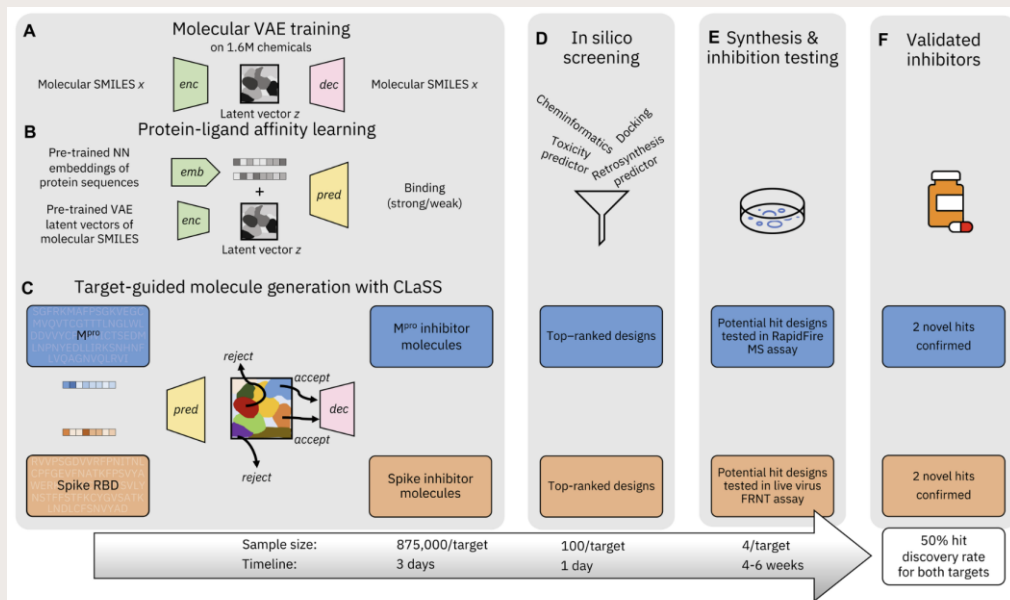


Fig. 1. Overview of our inhibitor discovery workflow driven by CogMol, a sequence-guided deep generative

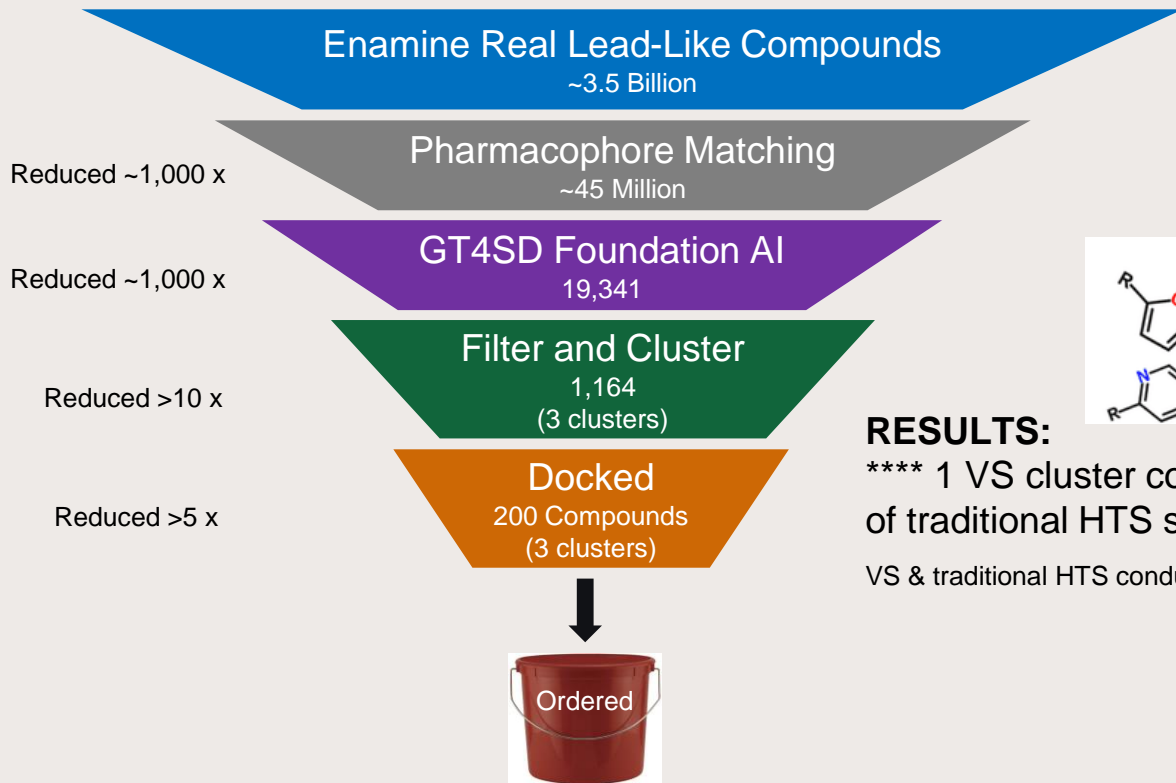
the CogMol Generative Framework relies on

- a chemical VAE,
- a protein sequence encoder,
- and a set of molecular property predictors,

all of which are pretrained on large amount of broad data—i.e., chemical SMILES, protein sequences, and available protein-ligand binding affinities.

# Cytokinetics begins utilization of GT4SD

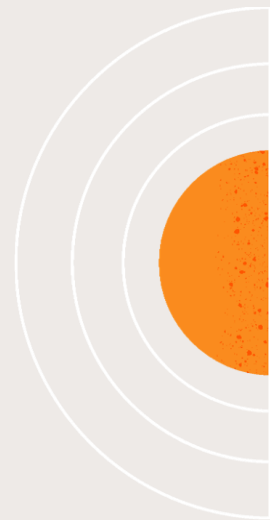
## In-House Virtual Screen



## RESULTS:

\*\*\*\* 1 VS cluster core matched top hit of traditional HTS screen!

VS & traditional HTS conducted independently



# Acknowledgements

## **Cytokinetics Informatics Team**

Marc Garard

Scott Rowland

## **Cytokinetics Research Management**

Bradley Morgan, SVP, Research and Non-Clinical  
Development

Anne Murphy, VP, Biology

Ajay Chawla, VP, Translational Sciences



# Thank you

Joe Michel

[jmichel@cytokinetics.com](mailto:jmichel@cytokinetics.com)

