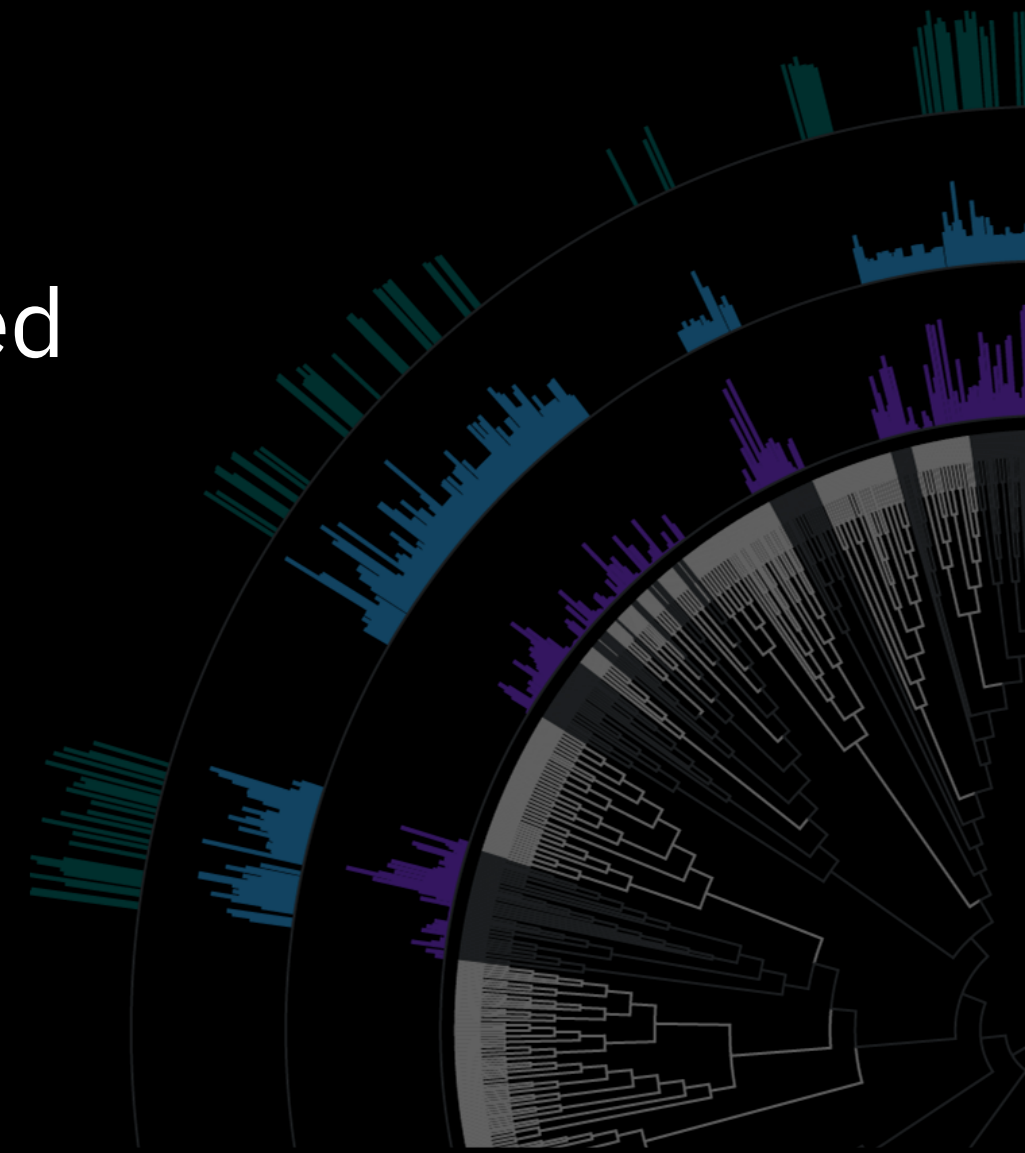


# IBM Research

## Foundation Models for Accelerated Discovery

John R. Smith, IBM Fellow, Accelerated Discovery  
IBM T. J. Watson Research Center, Yorktown Heights, NY USA  
[jsmith@us.ibm.com](mailto:jsmith@us.ibm.com)

October 2022



# IBM Research

**3,000**

researchers

**100s**

of disciplines

Almaden

Albany  
Yorktown

Cambridge

Dublin

Warrington

Zurich

Haifa

Delhi

Tokyo  
Shin-Kawasaki

Bangalore

Singapore

Nairobi

Rio de Janeiro  
Sao Paulo

Johannesburg

Cloud  
Computing

Artificial  
Intelligence

Quantum  
Computing

Semiconductors  
and Systems

Security and  
Cryptography

Physical  
Sciences

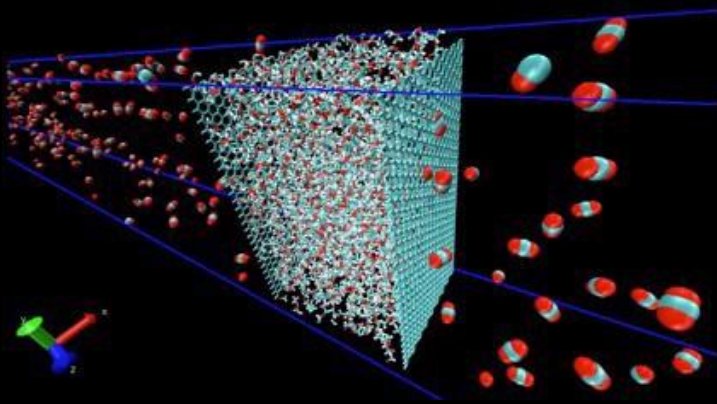
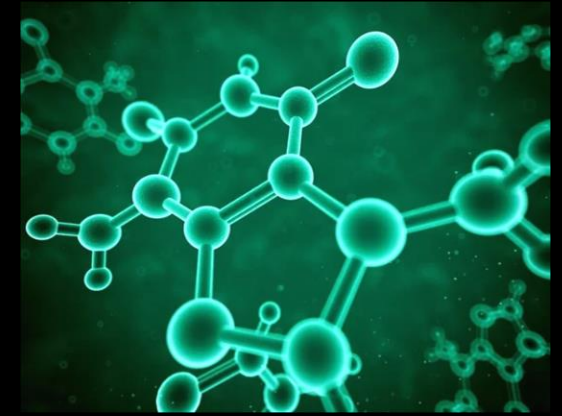
Mathematical  
Sciences

Life  
Sciences

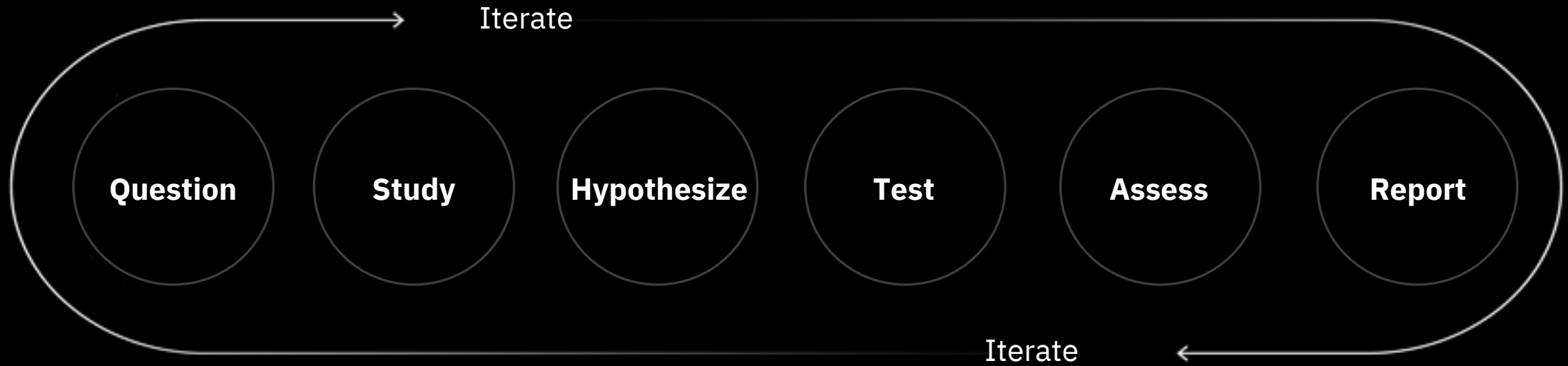
Computer  
Science

Accelerated Discovery

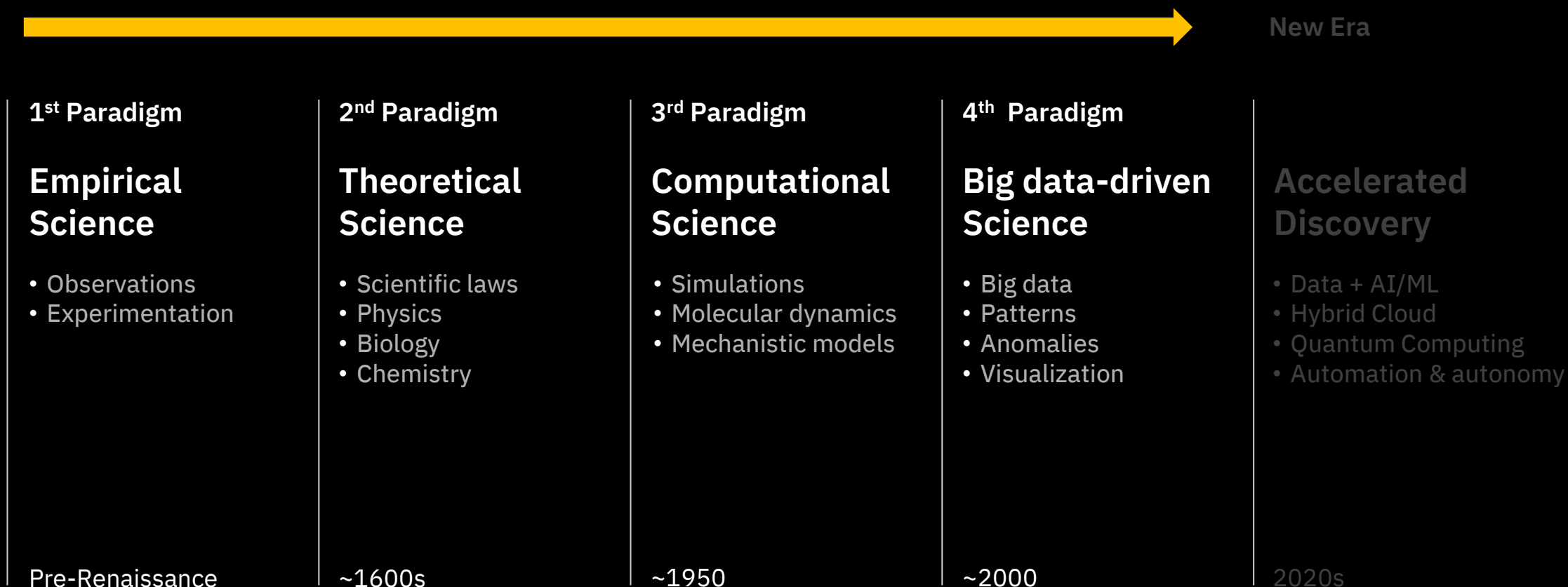
How do we discover solutions to complex problems?



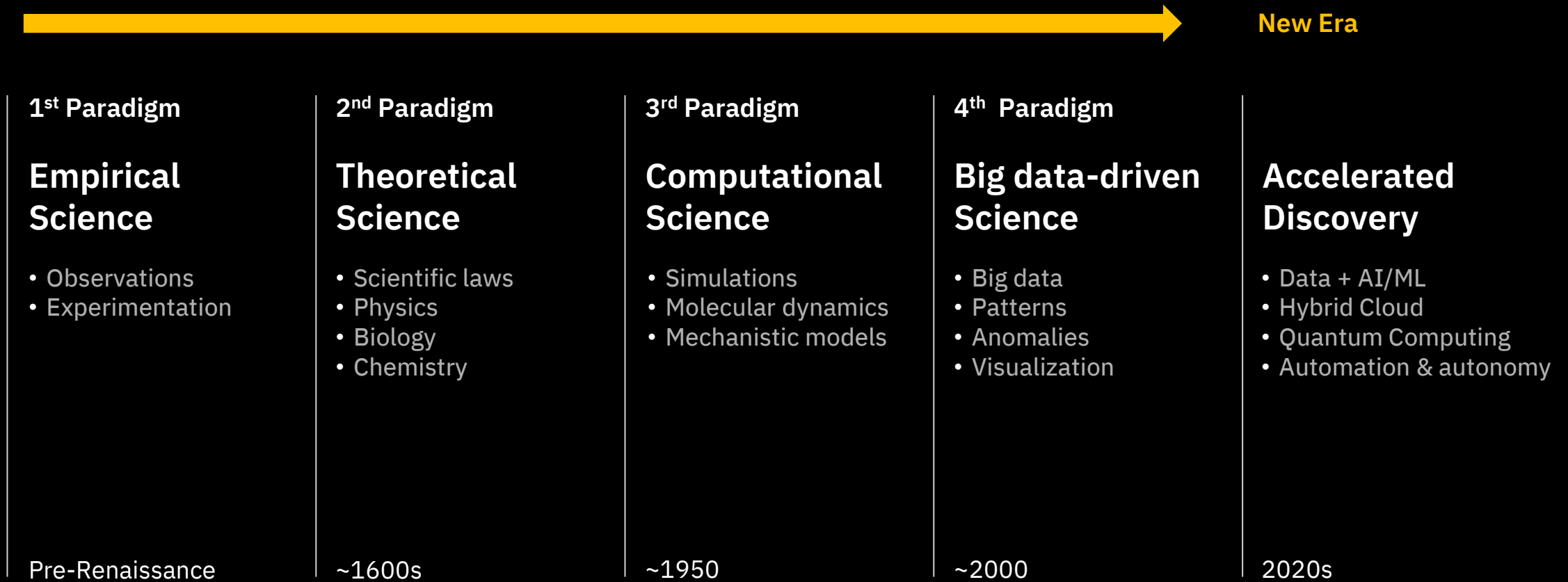
The scientific method has been our best tool for discovery...



# The scientific method has evolved over time



# We are entering a new era of discovery



# Automating and Accelerating Scientific Discovery

Extraction, integration and reasoning with knowledge at scale

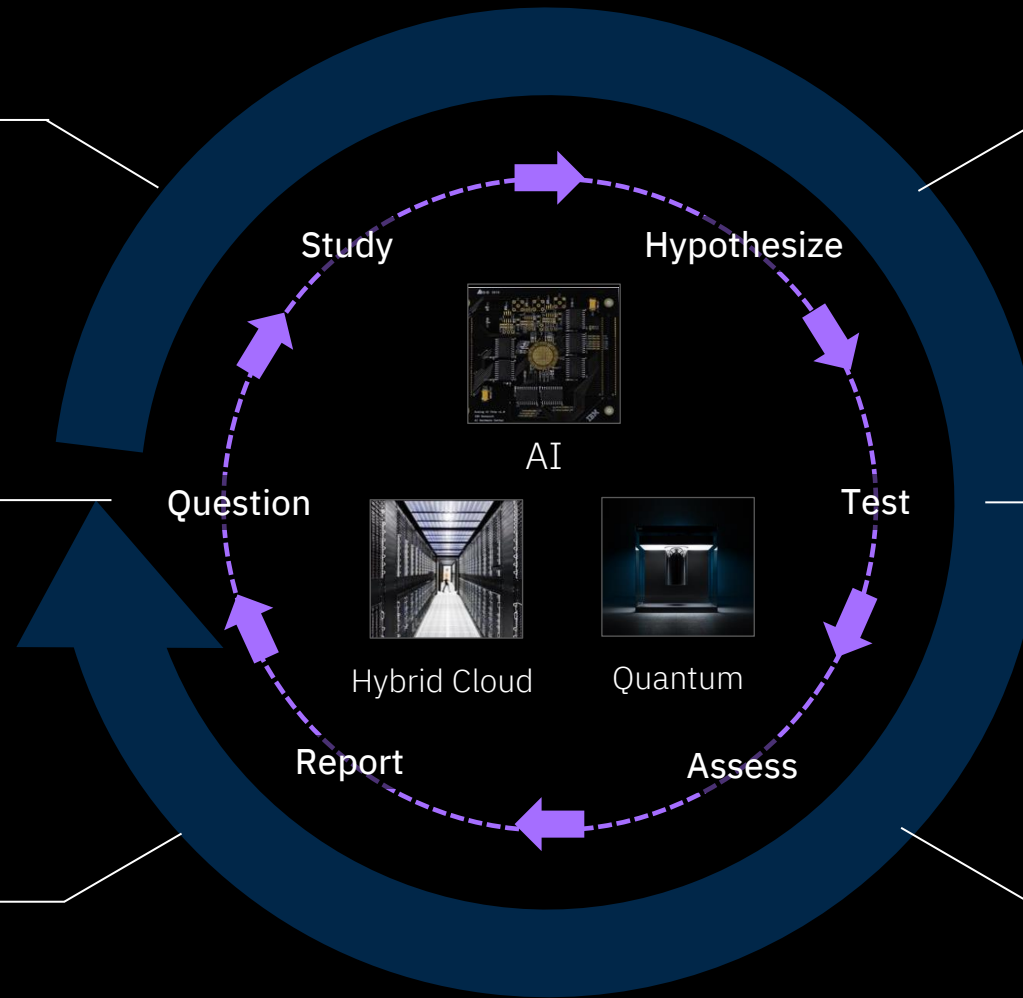
Tools help identify new questions based on needs and gaps in knowledge

Machine representation of knowledge leads to new questions

Generative models propose new hypotheses and automatically explore vast discovery spaces

Robotic labs automate experimentation and bridge digital models and physical testing

AI surrogate models integrated with simulation fill in missing information





# Examples in Materials Discovery

It takes roughly 10 years and upwards of \$10–\$100 million to discover one new material.

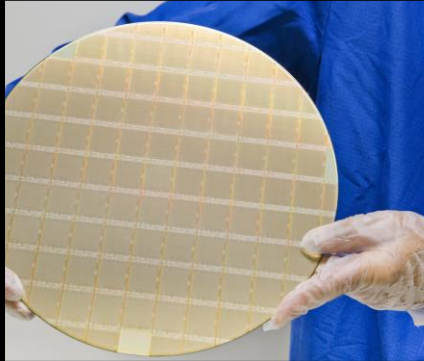


We aim to cut down both years and cost by 90%.



# Example Accelerated Discovery challenges and results for materials

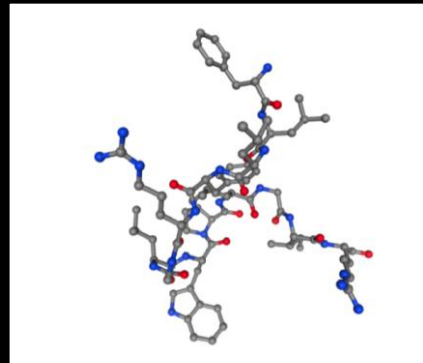
## Sustainable semiconductors



Discovery and synthesis of a new photoacid generator molecule in less than 1 year

<https://research.ibm.com/science/photoresist>

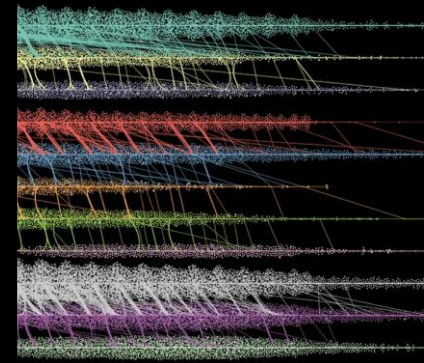
## Therapeutics



Discovery and validation of two new antimicrobial compounds in 48 days instead of 2-4 years

<https://research.ibm.com/publications/accelerated-antimicrobial-discovery-via-deep-generative-models-and-molecular-dynamics-simulations>

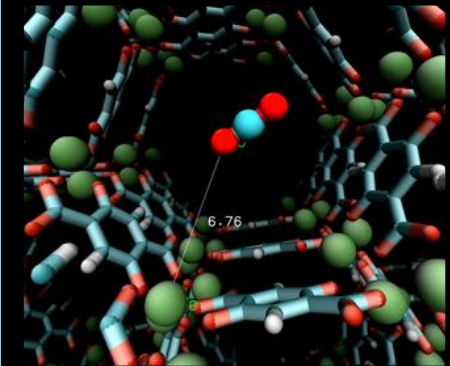
## Biomarkers



Discovered biomarkers and new disease trajectories for Type 1 diabetes

<https://research.ibm.com/blog/ai-predicting-onset-of-type-1-diabetes>

## Climate & Sustainability



Discovery of 500 molecular candidates for membranes to better separate CO<sub>2</sub> from flue gas

<https://research.ibm.com/blog/accelerating-materials-discovery>

# Discovery Technology Foundational Building Blocks

Accelerate scientific discovery through consumable “general purpose” discovery tools and platforms

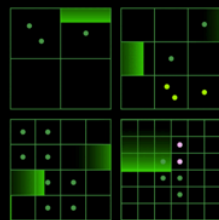
## Deep Search



<https://www.research.ibm.com/covid19/deep-search>

## Simulation (ST4SD)

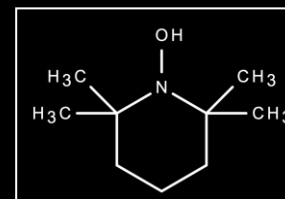
(Simulation Toolkit for Scientific Discovery)



<https://pages.github.ibm.com/st4sd/overview/>

## Generative Models (GT4SD)

(Generative Toolkit for Scientific Discovery)



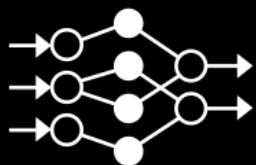
<https://github.com/GT4SD>

## Synthesis & Testing (RXN)



<https://rxn.res.ibm.com/rxn/robo-rxn>

## Foundation Models



"c1cccc1(C(=O)O)"

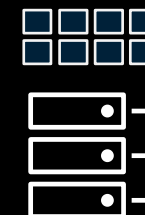
Knowledge  
Integration

AI- and  
quantum-  
enriched  
simulation

Hypothesis  
generation

Cloud-based  
AI-driven  
autonomous  
labs

## Discovery Optimization



Accelerated Discovery Workflows

# Deep Search: Structuring and Reasoning with Scientific Knowledge

## Large-scale curation and insight-extraction from unstructured multi-modal documents

Deep Search can parse large collections of scientific articles

Discovers entities and allows search via knowledge graphs

Inspection / cps-paper.pdf

Original document

FIGURE 7 The evaluation workflow to identify the petroleum system elements (PSE) in an article and filter to properties. It starts by searching for all petroleum system elements of a certain type (eg. source, reservoir or seal) and a particular report (workbooks 1 and 2). By passing graph traversals (workbook 3 & 5, 9, 11, 12) along specific edges and logical operations (workbook 4, 10, 13, 14), we are able to obtain a list of candidate formations (workbook 15), age (workbook 16) and rocks (workbook 17), ranked by their accumulated weight. Execution of this query takes less than 18 ms on average.

TABLE 1 Top-k accuracies validation of KG query results. Numbers represent the fraction in which any of the k highest ranked answers matches the expected answer

PSE	Property	Top-1	Top-2	Top-3	Top-5
Reservoir	Age	0.82	0.96	0.98	1.00
	Formation	0.93	0.98	1.00	1.00
	Rock	0.62	0.80	0.87	0.94
Seal	Age	0.73	0.91	0.94	0.97
	Formation	0.82	0.94	0.97	0.98
	Rock	0.82	0.92	0.95	0.97
Source	Age	0.75	0.92	0.96	0.97
	Formation	0.89	0.96	0.97	0.98
	Rock	0.83	0.92	0.95	0.96

CONCLUSIONS

With the introduction of the CPS platform, we demonstrate substantial benefit for domain experts and data scientists in parsing deep exploration of published knowledge in a fully integrated, yet modular cloud solution. CPS seamlessly connects to the CCS, complementing it with a highly scalable, automated pipeline to build consistent domain knowledge models and an intuitive, powerful approach to explorational queries and graph-scale analytics. This is accomplished through three fundamental design considerations: (1) We do not require manual data curation or annotation; (2) We built a scalable, efficient approach to parse the literature, metadata and multimedia; all embedded in

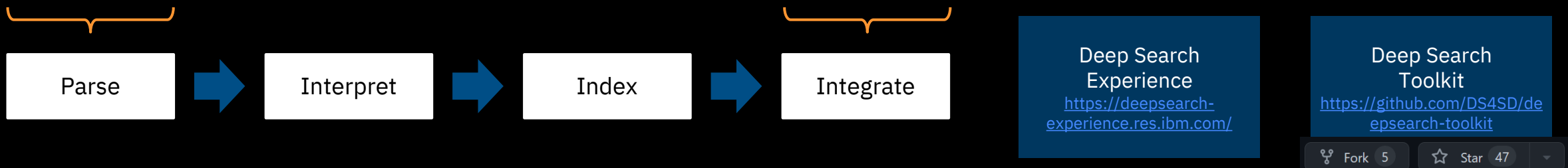
IBM Corpus Processing Service

KG-prod Search Results Workflows Downloads

Filter

Nodes 94'219

- abbreviation 867
- abstracts 91
- affiliations 0
- atomic-elements 36
- authors 0
- claims 0
- classifications 0
- documents 91
- experiment 383
- functional-groups 105
- images 0
- material 5'311
- material-attribute 35
- material-class 43
- material-property 144
- material-property-to-value-... 24'621
- material-to-abbreviation 454
- material-to-material-class 5'366
- material-to-material-property 13'382
- materials-database 12'295
- paragraphs 13'576
- publishers 0
- smiles 5'570



# Simulation Toolkit for Scientific Discovery (ST4SD)

Filling knowledge gaps with AI-enriched modeling and simulation

## ST4SD runtime

A runtime for simulation workflows that includes memoization and surrogate support as well as pluggable HPC backends

## AI surrogates

Includes and enables AI surrogate and hybrid workflows including configurable strategies for performance monitoring and risk management

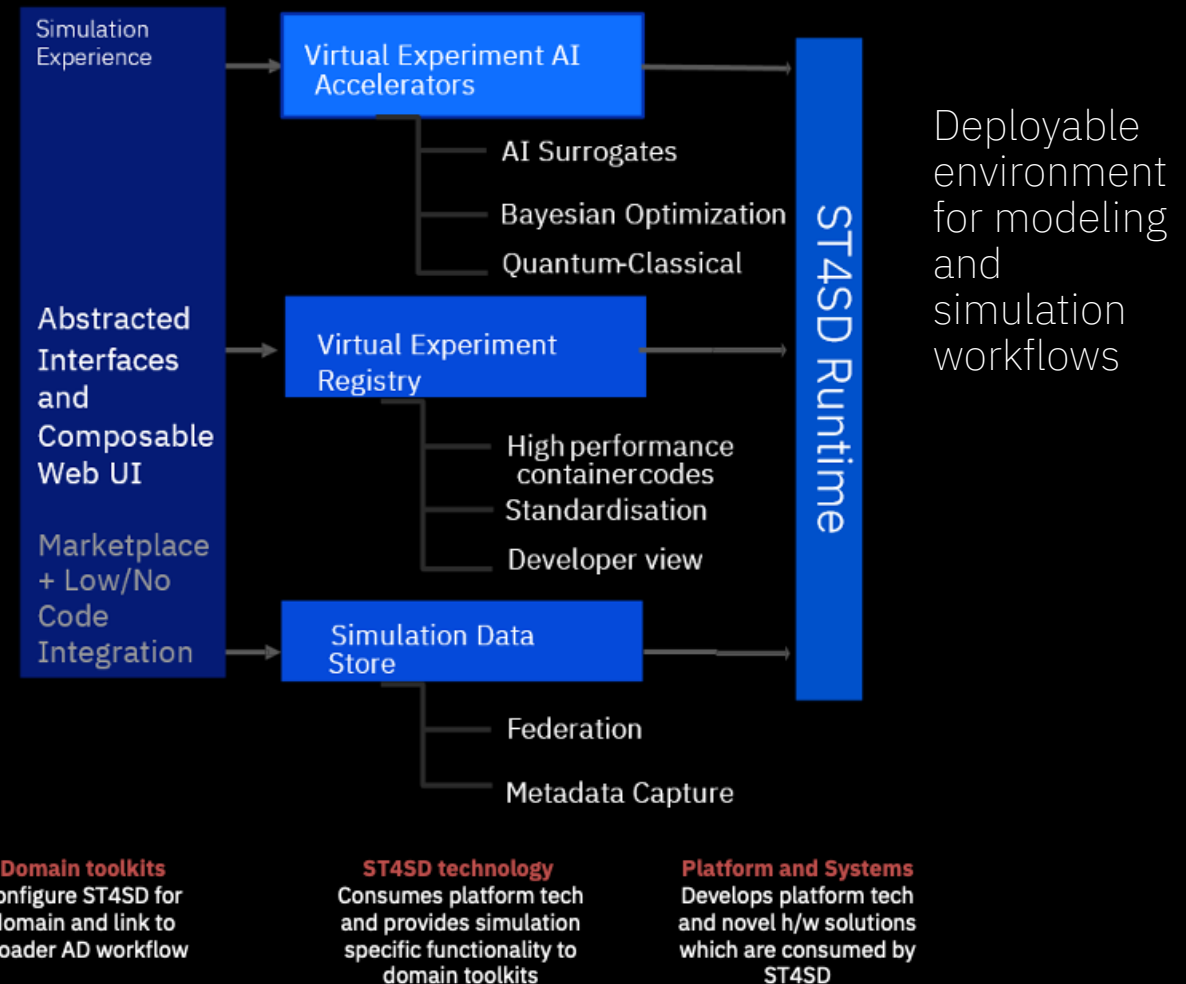
## Virtual experiment registry

Includes pre-built workflows and building blocks for composing and hosting virtual experiments

## ST4SD experience

User interactive tool with pre-packaged experiments

## ST4SD Toolkit



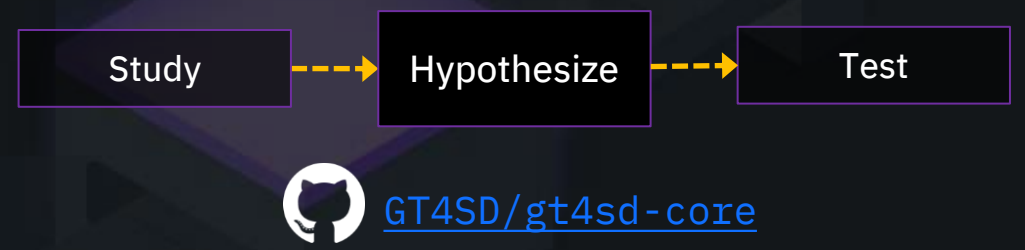
# Generative Toolkit for Scientific Discovery (GT4SD)

Open-source library to accelerate hypothesis generation in scientific discovery



Repository statistics and badges:

- pypi package 0.50.0
- Running tests: style, pytests and entry-points passing
- License MIT
- code style black
- contributions welcome
- website live
- downloads 18k
- downloads/month 2k
- launch binder
- DOI 10.5281/zenodo.7073764
- Award 2022 IEEE Open Software Services Award



## 1. Train generative models

```
gt4sd-trainer --training_pipeline_name paccmann-vae-trainer --epochs 25
```

## 2. Create inference pipelines

```
gt4sd-saving --training_pipeline_name paccmann-vae-trainer --model_path
```

## 3. Run inference pipelines

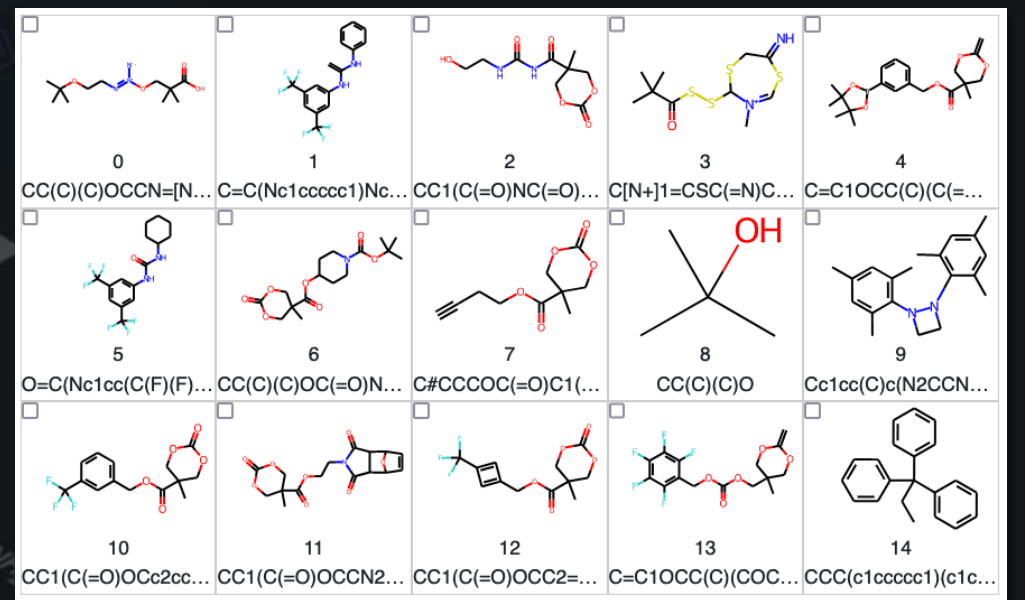
```
gt4sd-inference --algorithm_name PaccMannGP --algorithm_application PaccMannGP
```

## 4. Share your models with the community

```
gt4sd-upload --training_pipeline_name paccmann-vae-trainer --model_path
```

Applications include hypothesis generation for inverse design and discovery of materials

## Example molecules generated using GT4SD



# IBM RoboRXN – intelligent lab automation for the cloud

## AI + Hybrid Cloud + Robotic Labs for automated synthesis prediction and execution

35,000  
Users via cloud

9+ Million  
Reaction  
predictions

7  
Industrial  
partners

Synthesizing new molecule

Started: Nov 30 2020, 6:49am PT

Live from IBM RoboRXN

Action 2 Overview

Adding  $C_2H_3F_3O_3S$

In this action, the molecule methyl trifluoromethane sulfonate is added to Reactor 2.

Methyl trifluoromethane sulfonate  $C_2H_3F_3O_3S$

Methyl trifluoromethane sulfonate is a brown liquid. Insoluble in water. This material is a very reactive methylating agent, also known as methyl triflate.

**NOW**

10 ml of reagent containing methyl trifluoromethane sulfonate is being moved from Vial 61 and added to Reactor 2.

Position of the robot arm  
Moving to Vial 61

00:06:00 || LIVE

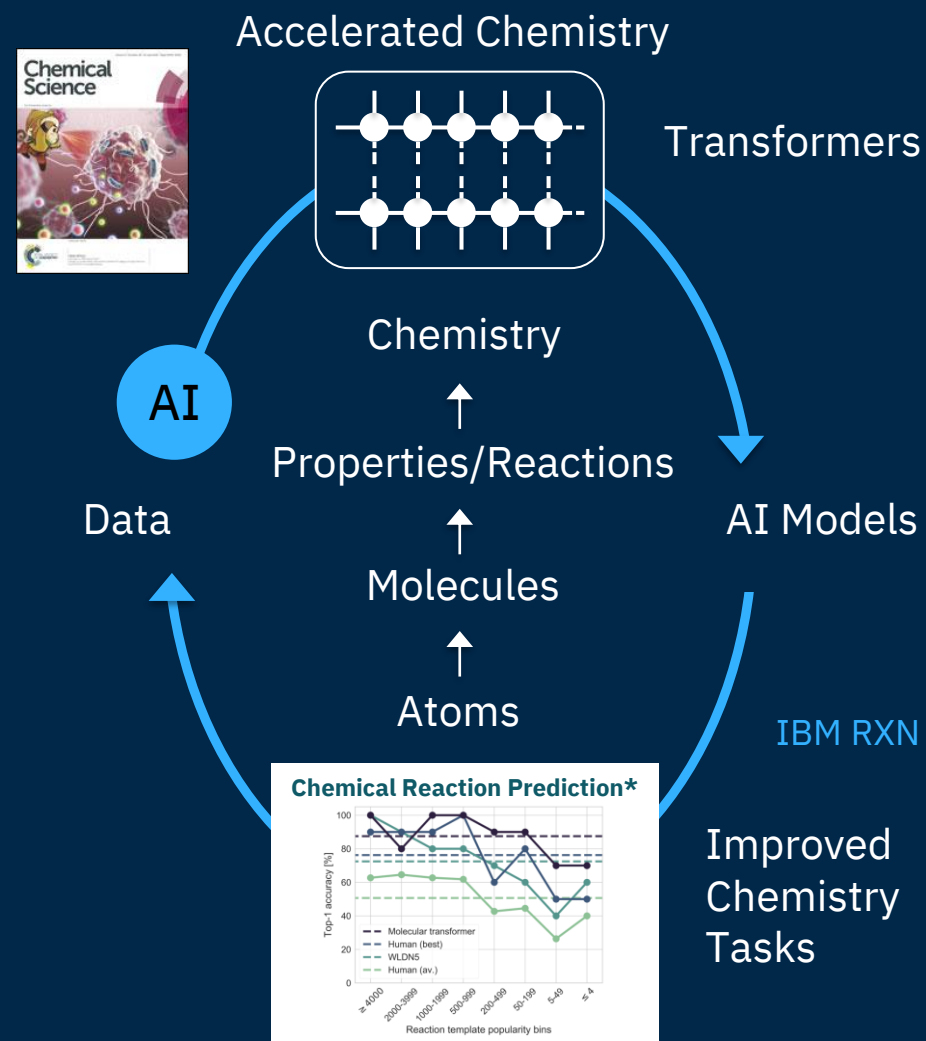
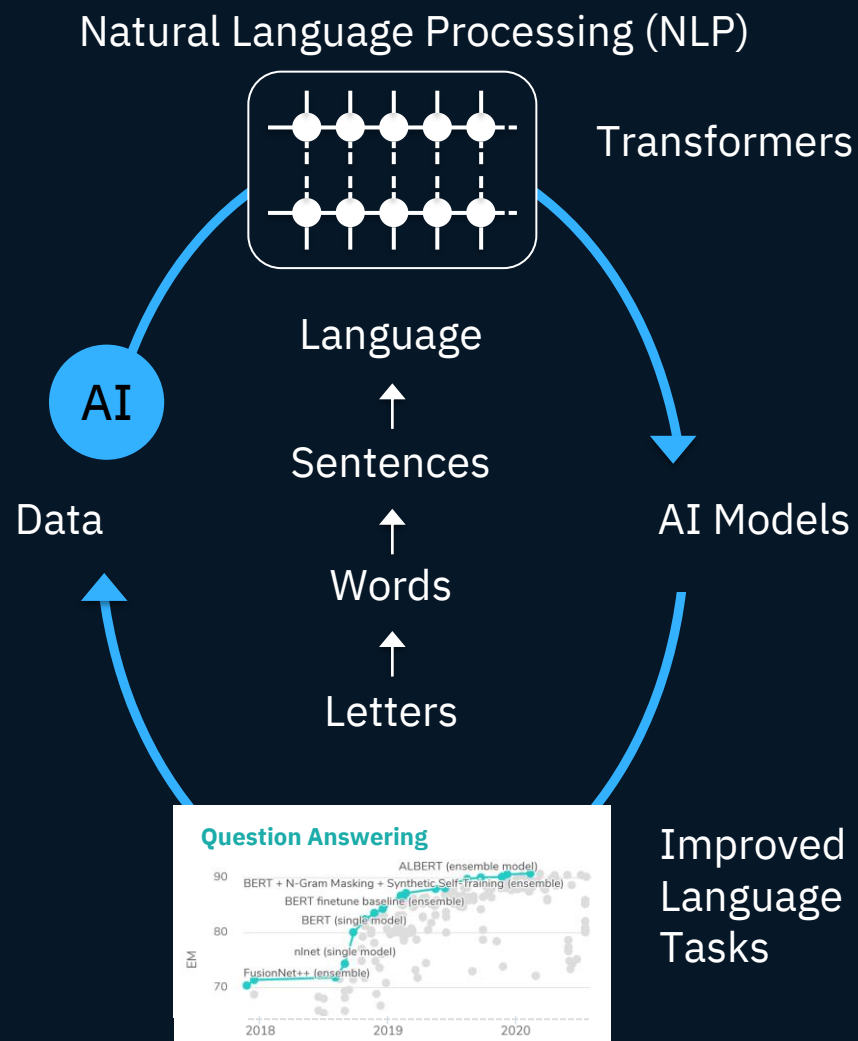
IBM Research



# Foundation Models

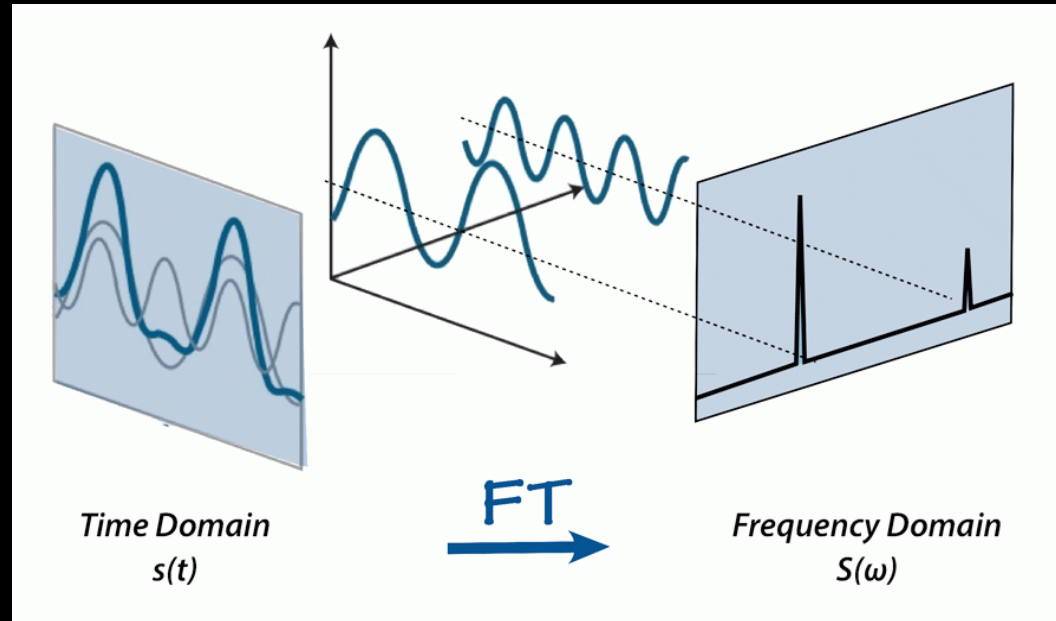
# The same AI breakthroughs happening in language are changing scientific discovery

Foundation models are powering new capabilities



# Well chosen representations can simplify complex problems

*Example from the world of signal processing*



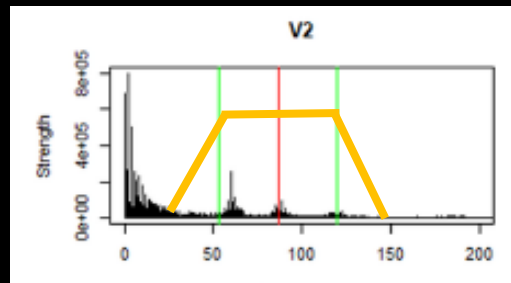
*Data represented in the frequency domain can reduce complexity for filtering and denoising*

Time-domain



FT  
→

Frequency-domain

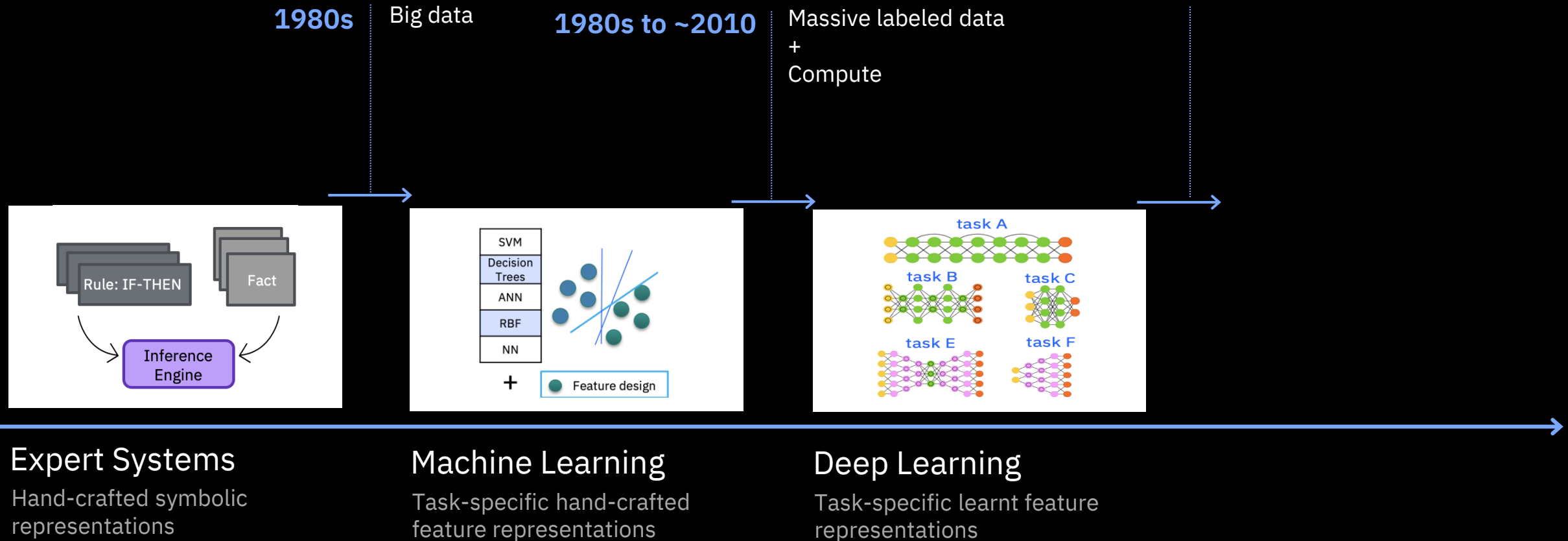


$FT^{-1}$   
→

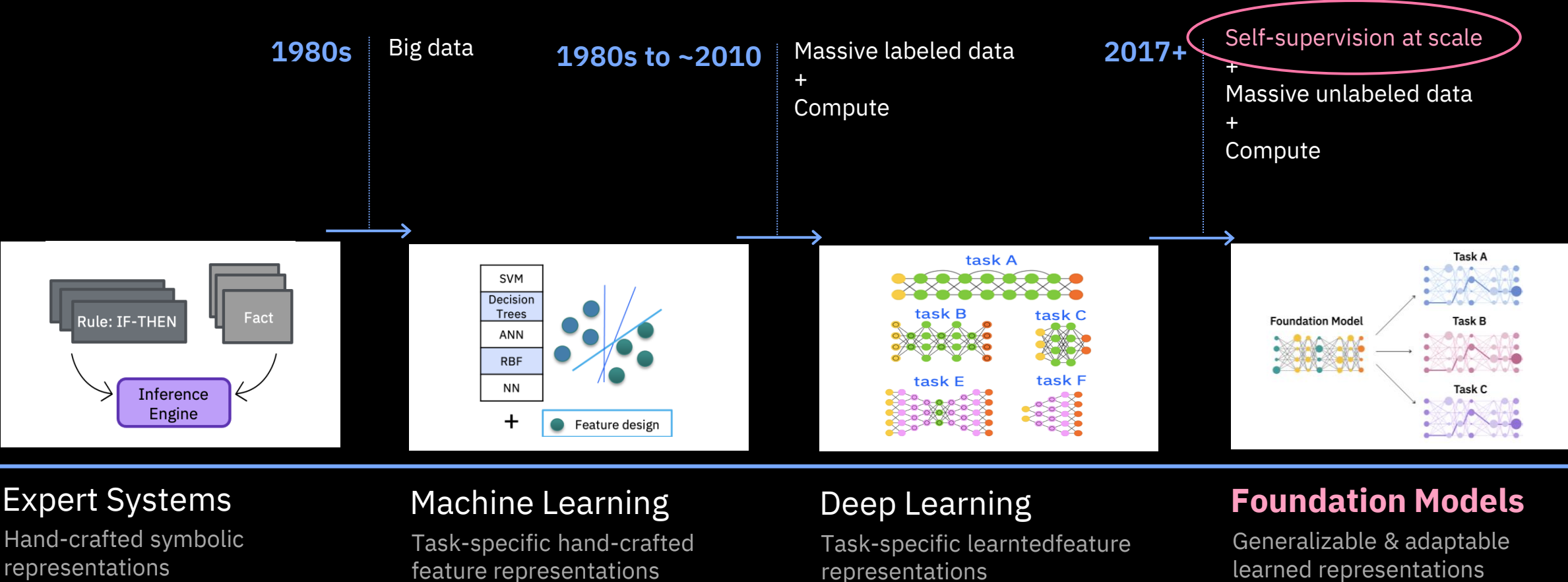
Time-domain



# Story of AI has also been a story of data representations

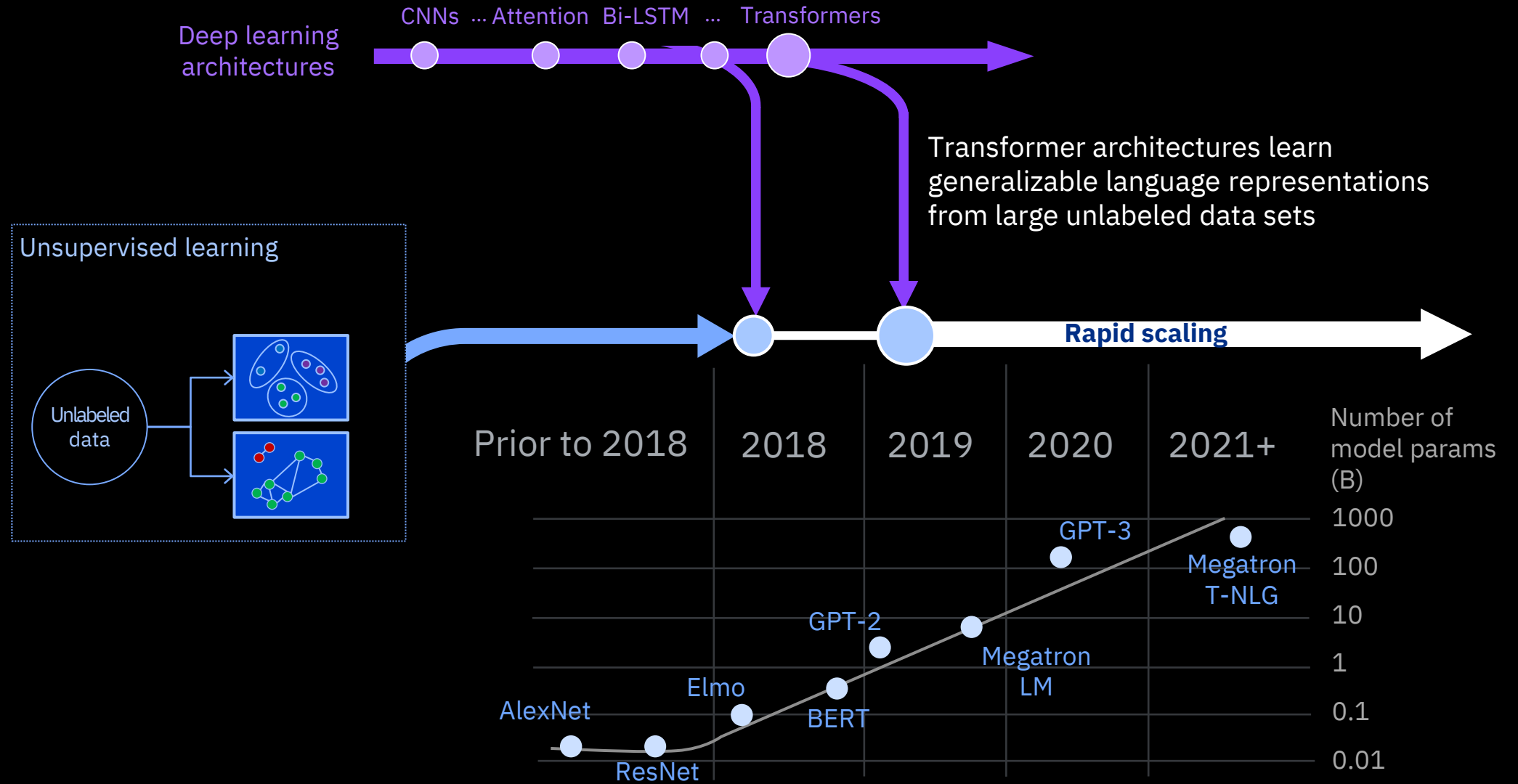


# Foundations models are an emerging class of representation that is changing the AI landscape

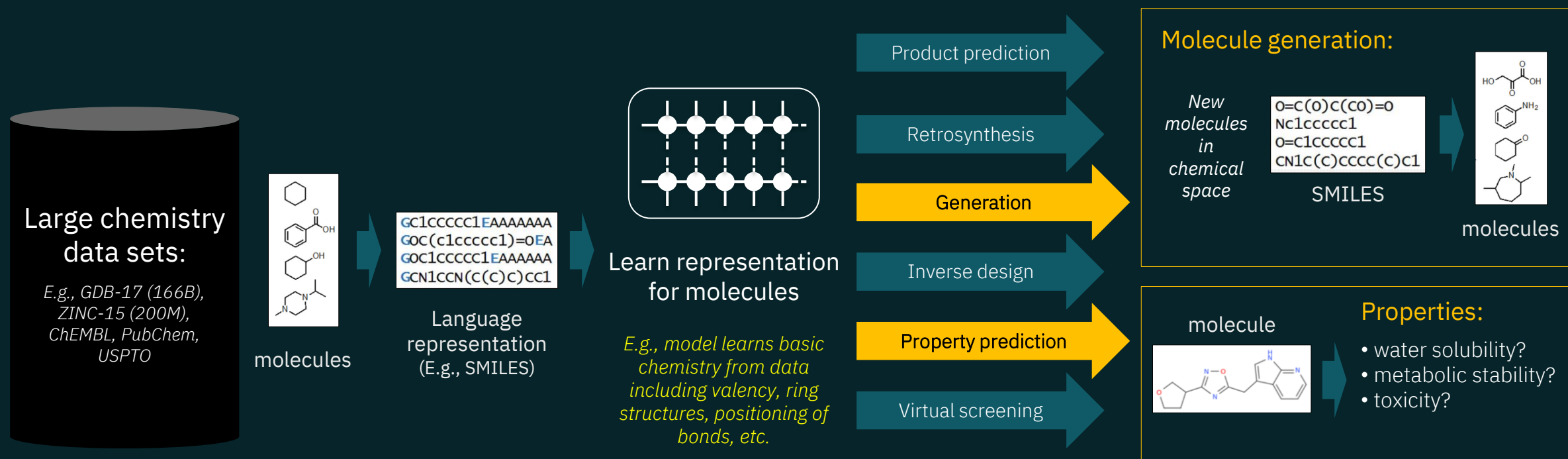


# Success of Foundation Models for Language from self-supervision *at scale*

Enabled by a novel learning architectures + data + compute



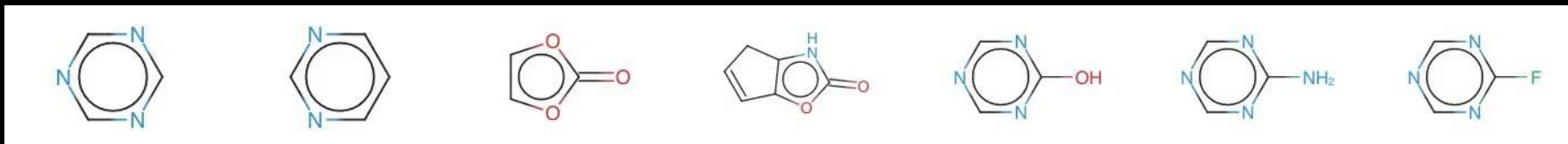
Foundation models for molecules enable multiple downstream tasks for molecule design, property prediction, reaction prediction, and more.



# Generative models are a powerful tool for molecular inverse design and discovery

Learn from data to generate hypothetical and novel candidates for targeted properties

IBM Research, “Molecular Inverse-Design Platform for Material Industries”, *KDD*, Aug 2020



IBM Research, “Accelerated antimicrobial discovery via deep generative models and molecular dynamics simulations,” *Nature Biomedical Engineering* 2021

IBM Research, “Target-Specific and Selective Drug Design for COVID-19 Using Deep Generative Models,” May 2020

Explore novel drug candidates for COVID-19

Biological target: SARS-CoV-2 Main Protease

Subset: GEN

Sort by: Target Affinity (AFF)

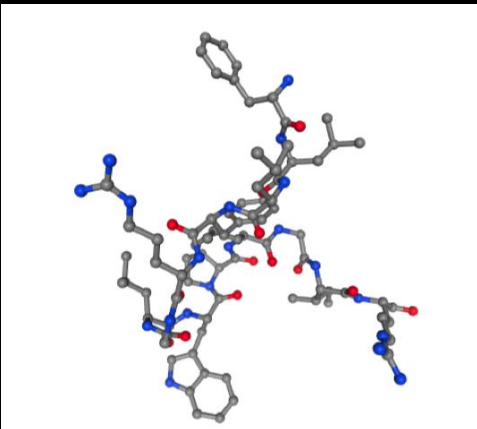
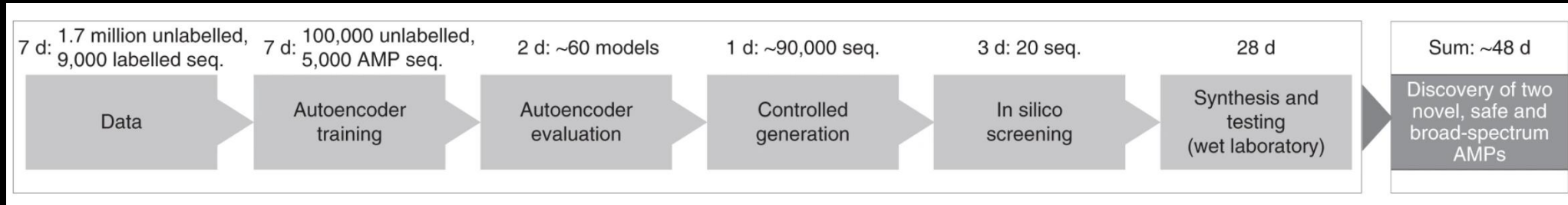
- Target Affinity (AFF)
- Drug-likeness (QED)
- Docking Energy (Dock)
  - Above
  - Below
- Synthetic Accessibility (SA)

GEN	AFF	QED	Dock	SA	LogP	NOV	SEL	TOX	MolW
1 GEN980	9.21	0.59	-7.3	2.15	3.17	0.13	1.36	1	307.35
2 GEN153	9.15	0.72	-8	2.73	3.14	0.02	2.09	0	386.41
3 GEN949	9.05	0.77	-7	2.61	3.42	0.13	1.21	0	350.2
4 GEN522	9.04	0.45	-7.5	2.61	2.22	0.1	1.34	1	433.49



# Example: Generative Models for New Antibiotic Discovery

Accelerating end-to-end anti-microbial discovery with generative models and molecular dynamics. Results demonstrate a 48 day time to discover and validate two new compounds, versus 2-4 years.



- Deep-neural-net-based generative models screen for antimicrobial function, broad-spectrum efficacy, presence of secondary structure and toxicity.
- Simulation confirms mode of action
- Synthesis and test reveal two new potent compounds

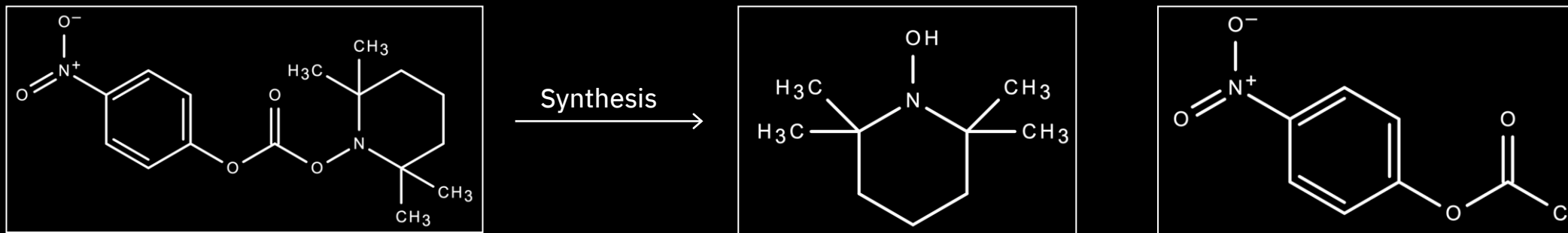
Two compounds displayed high potency against diverse Gram-positive and Gram-negative pathogens and a low propensity to induce drug resistance in *Escherichia coli*.

IBM Research & Oxford, "Accelerated antimicrobial discovery via deep generative models and molecular dynamics simulations," [Nature Biomedical Engineering](#) 2021

# IBM RXN for Chemistry uses AI foundation models to predict chemical reactions and chemical procedures

## Chemical Reaction predictions

1



## Chemical Procedure (recipe) predictions

2

**Add** 2,2,6,6-Tetramethylpiperidin-1-ol (2.0 g)

**Mix** 4-nitrophenyl chloroformate (1.1 g) and dichloromethane (5 ml) in a separate vial

**Add** mixture dropwise at -10°C

**Stir** for 1 hour at -10°C

**Quench** with aqueous Na<sub>2</sub>CO<sub>3</sub> (15 ml)

**Extract** with dichloromethane (20 ml)

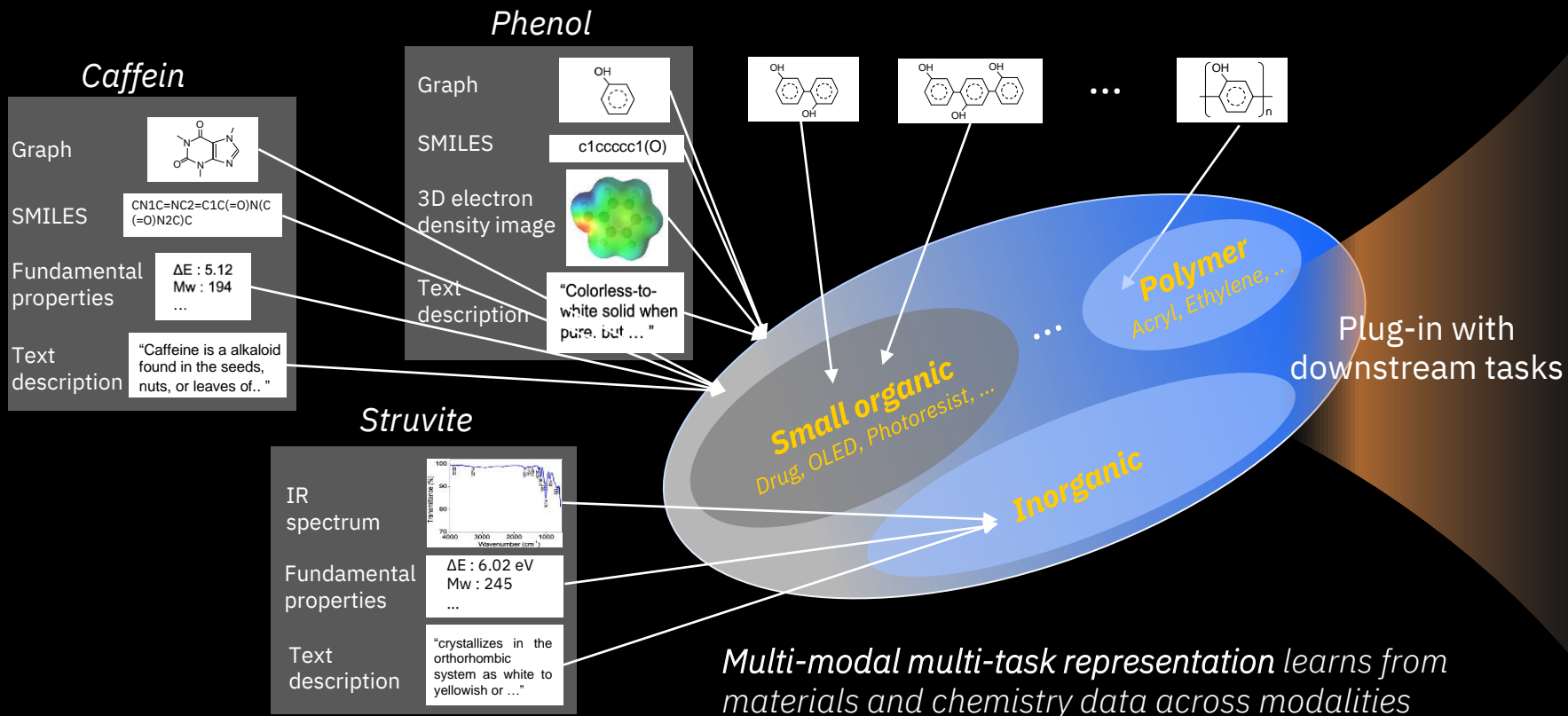
**Concentrate**

Next Steps

# (1) Learning universal representations for materials, chemistry, biology, etc

## Multi-modal foundation models that can power a broad field of downstream tasks

- Comprehensive reusable representation across domains – breaks down silos of independent modeling
- Integration of multiple modalities fuses diverse sources of knowledge and data



### Generative model

- Generate new materials' information with multiple modalities (image, text description, etc.)
- Target properties with limited data
- Interpretable generation

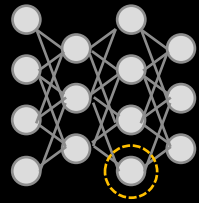
### Predictive model

- Predict properties in small/zero data domains
- Interpretable prediction

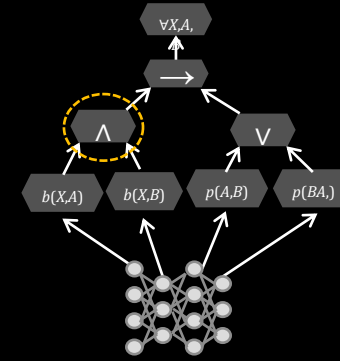
### Simulation accelerator

- Intellectual speculative skip of simulation steps by orbital information (surrogate model)

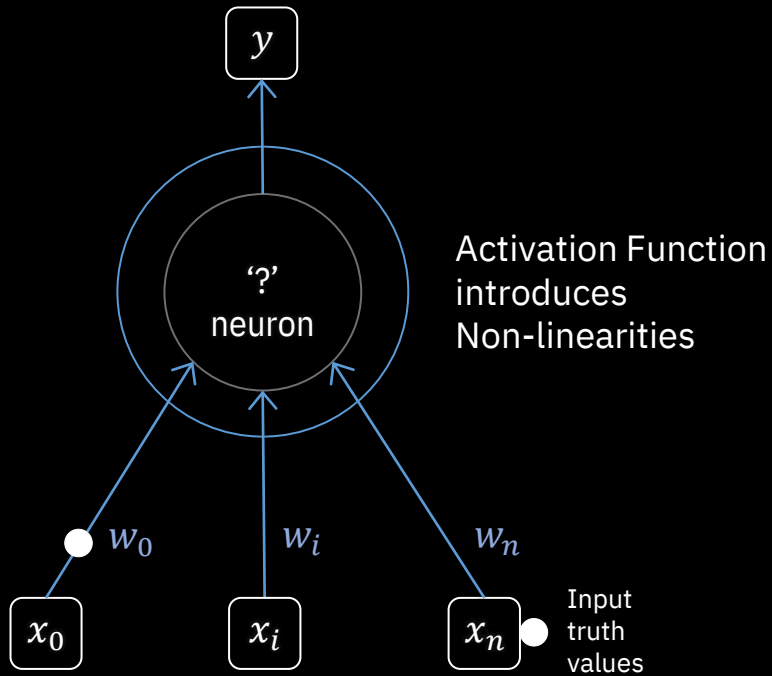
# (2) Integrating Real-valued Logic into Foundation Model Representations using **Logical Neural Networks (LNNs)** for Improved Explainability



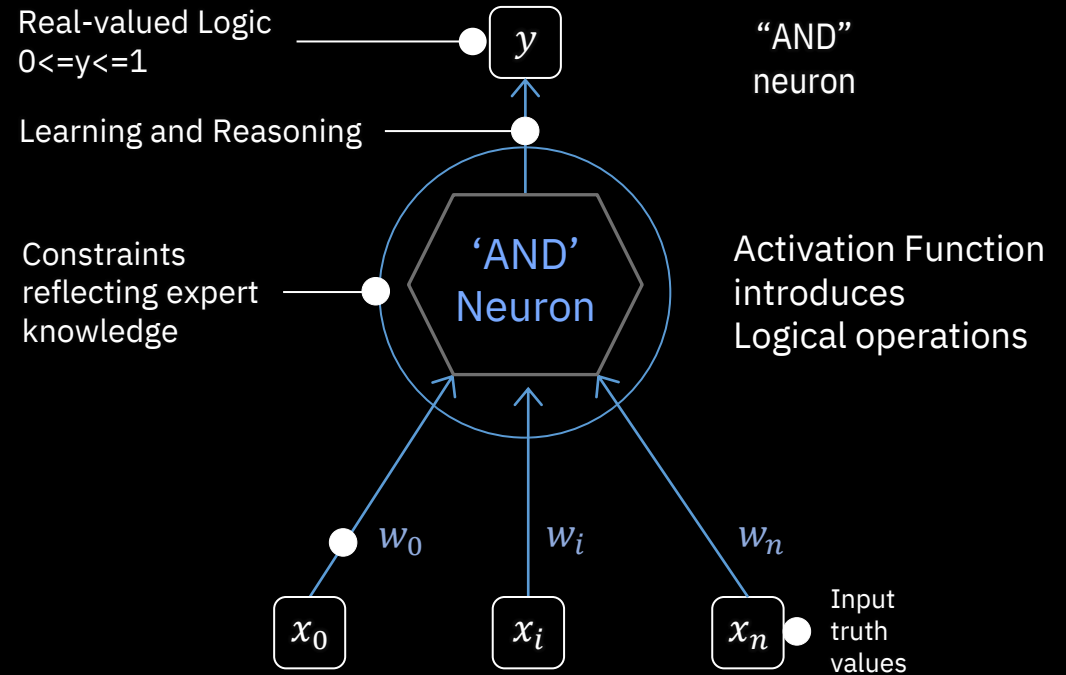
Standard Neuron



Neuron architecture in Logical Neural Networks

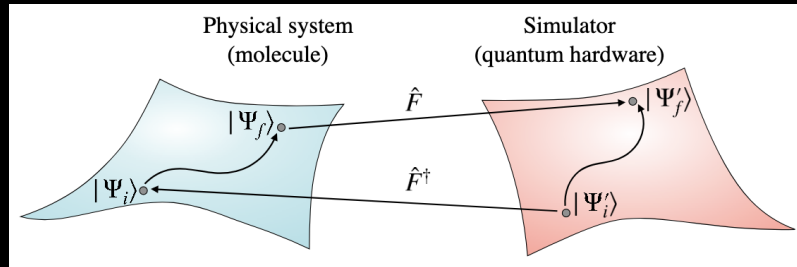


Werbos, 1974  
Rumelhart, Hinton and Williams, 1986



Riegel, et al., 2020, Logical Neural Networks

### (3) Quantum circuit-based representations can potentially compute properties in physics and chemistry more naturally, accurately and efficiently



Quantum computers are universal simulators of physics and chemistry – can provide efficient representations of dynamics by quantum circuits



Molecule

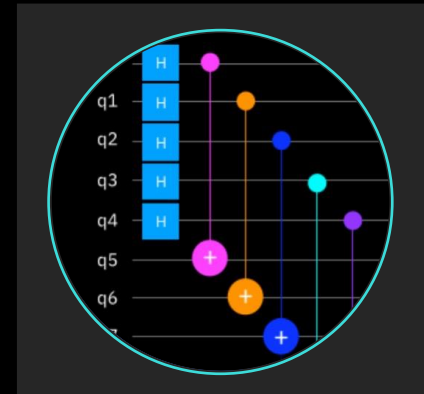


$$H = \sum_{p,q} h_{pq} c_p^\dagger c_q + \frac{1}{2} \sum_{p,q,r,s} h_{pqrs} c_p^\dagger c_q^\dagger c_q c_s$$

↓

$$H = \sum_{\mu} c_{\mu} P_{\mu}$$

Mapped to qubit operators



Mapped to quantum circuits representation

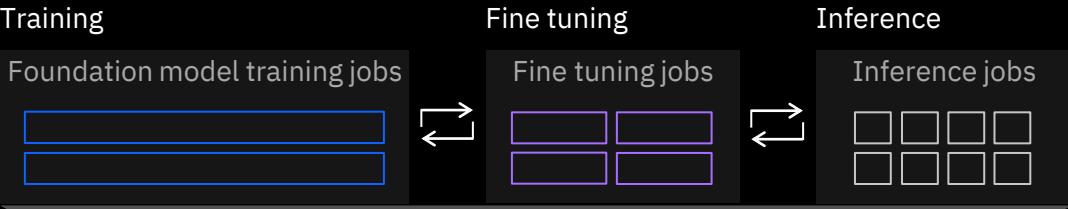


$$P = \sum$$

Properties measured on quantum computer

Potential applications include drug discovery, next-gen battery design, corrosion analysis, structural analysis, new materials design, solar conversion, catalysts and enzyme design


# (4) IBM tools and infrastructure for foundation models



**Foundation Model Pipelines and Tools**  
(Pre-built workflows, trust algorithms, evaluation metrics, testing & validation flows...)

**Scale-out Middleware**  
(Serverless ML, Workflows/Pipelines, AI Automation...)

**Hybrid Cloud Platform**  
(Distributed resource manage, placement, scheduling)



**Scale-out Infrastructure**

GPU Compute Cluster    GPUs  
AI Accelerators  
High speed networking

## Scale-out infrastructure

- Optimized container networking implementations
- AI accelerators

## Scalable middleware stack for distributed training

- Store and compute disaggregation
- Elastic and fault-tolerant distributed training
- Data caching
- Efficient resource management (job placement & scheduling)

## Tools and workflows for train, test, and consume

- Tools to support data scientist productivity
- Standardized reusable pipelines for train, test, validate, fine tune
- Tools for automatic metric collections for trust and quality control
- Scalable tools for model exploration & evaluation

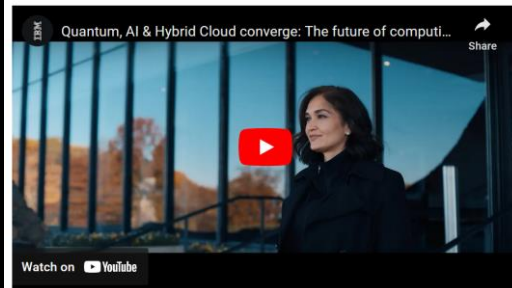
<https://research.ibm.com/topics/foundation-models>

# For more information ...

## Accelerated Discovery

### What's next in computing: The era of accelerated discovery

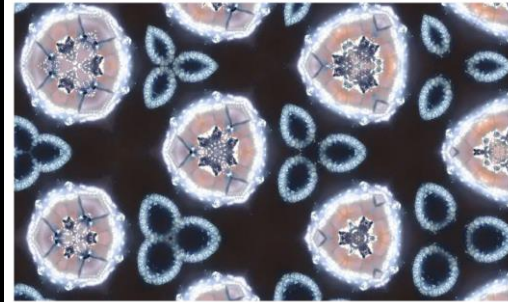
To meet the growing challenges of an ever-shifting world, the ways we have discovered new ideas in the past won't cut it moving forward. A convergence of computing revolutions taking place right now will help accelerate the rate of scientific discovery like nothing before.



<https://research.ibm.com/blog/what-is-accelerated-discovery>

### Accelerating discovery for societal and economic impact

Much of the world's ability to mitigate the effects of climate change will come down to our ability to quickly identify new materials that can be created, consumed, and recycled with minimal environmental impact.

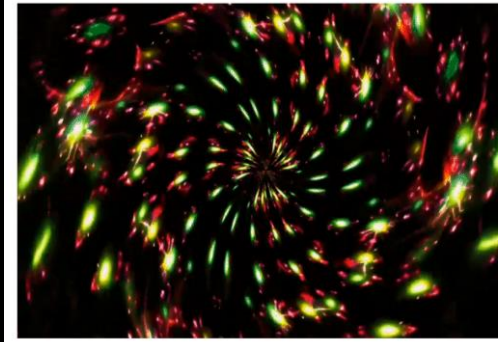


<https://research.ibm.com/blog/new-sustainable-materials>

## Foundation Models

### What are foundation models?

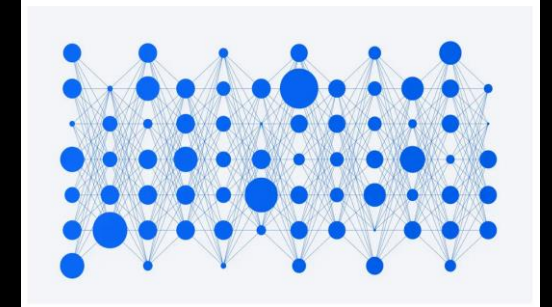
The future of AI is flexible, reusable AI models that can be applied to just about any domain or industry task.



<https://research.ibm.com/blog/what-are-foundation-models>

### IBM Research unveils two key advances for foundation models

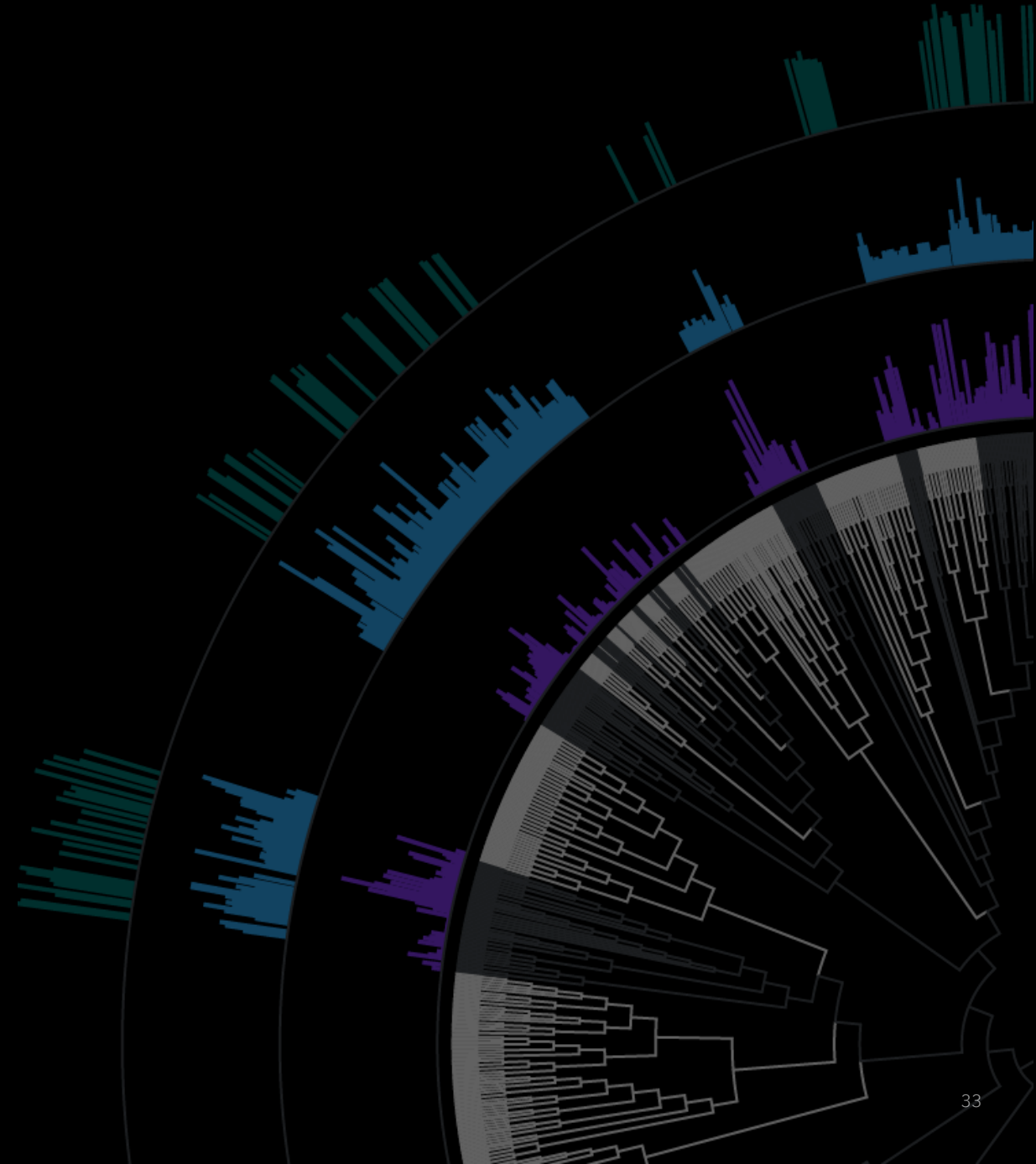
At this year's Ray Summit, researchers at IBM showed off two features, running on top of Ray, that make it easier to set up and run foundation models for AI workloads.



<https://research.ibm.com/blog/ray-summit-codeflare-foundation-models>



# Thank You!



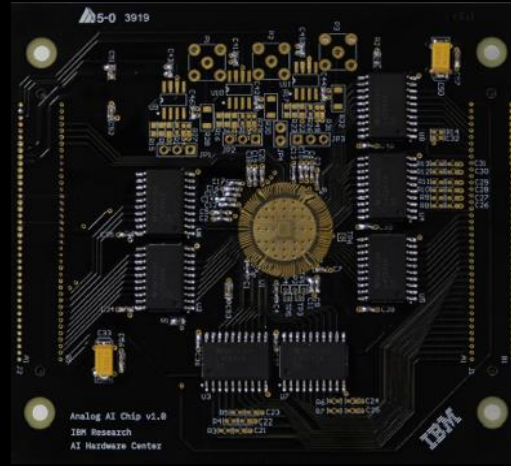
# What's Next in Computing?



Bits

**Hybrid Cloud:** cloud + high performance computing + robotics, instruments and labs

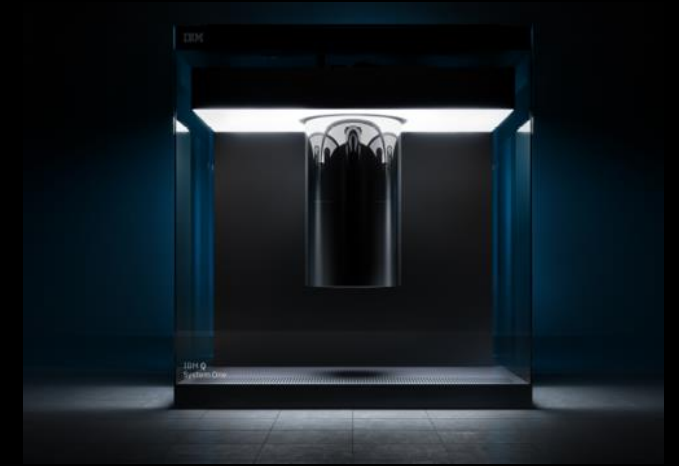
+



Neurons

**Artificial Intelligence (AI):** advanced data-driven modeling, analytics and automation

+



Qubits

**Quantum Computing:** making intractable problems tractable through a new computing modality