# Genomics Projects Will Exceed 40 Exabytes in the Next Decade

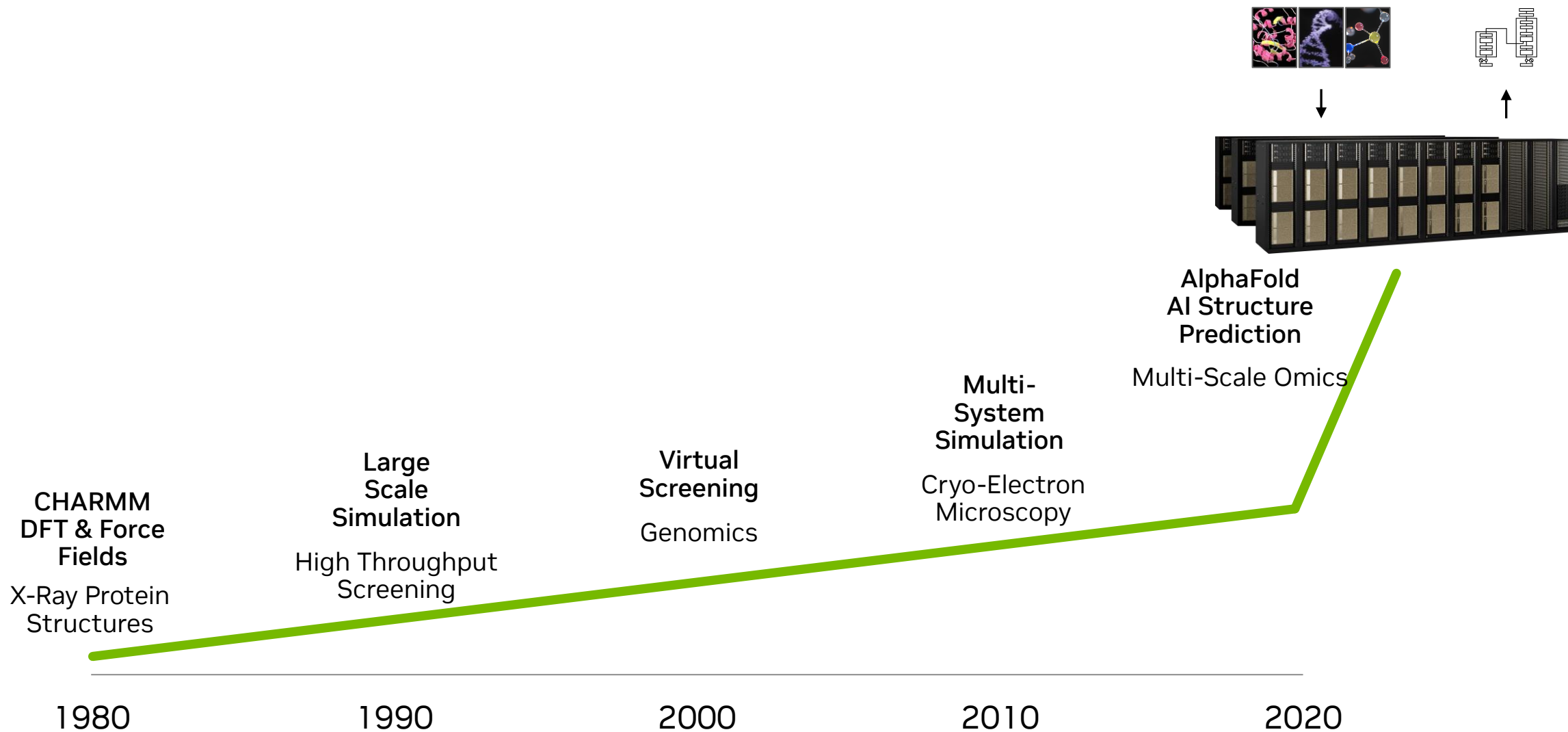## As sequencing becomes less expensive, the data deluge grows



"Our ability to sequence DNA has far outpaced our ability to decipher the information it contains, so genomic data science will be a vibrant field of research for many years to come."

National Human Genome Research Institute

nVIDIA.

# Biomedical Research and Drug Discovery Is at an Inflection Point

## Computer Aided Drug Discovery is Expanding Exponentially



**AlphaFold AI Structure Prediction**

Multi-Scale Omics

**Multi-System Simulation**

Cryo-Electron Microscopy

**Virtual Screening**

Genomics

**Large Scale Simulation**

High Throughput Screening

**CHARMM DFT & Force Fields**

X-Ray Protein Structures

1980          1990          2000          2010          2020

# AI is Transforming the Drug Discovery Process

## Deep learning is an essential tool for modern R&D



**Target ID**

Gene Expression Prediction, scRNA analysis

Accelerated cryo-em & protein structure prediction

Knowledge synthesis from scientific literature

**In-Silico Discovery & Design**

Active-learning virtual screening

AI powered molecular property prediction & generation

Drug-target interaction prediction

**Clinical Development**

NLP for Clinical trial matching

Adverse event monitoring

Histopathology/Radiology/OMICS biomarker ID

*Source: https://www.nature.com/articles/s42256-022-00463-x*

# Lab in a Loop: AI to Transform Drug Discovery and Development



Aviv Regev, Head of Genentech
Research and Early Development

# Generative AI Is Used to Design Biologics

Antibody Foundation Models | Post-processing Analysis

## Antibody LLM based on ESM-1nv in BioNeMo



## Downstream processing by RAPIDS, *e.g.*, UMAP



20sec/structure — Faster protein structure prediction

100x — Faster post-training analysis

<1month — From onboarding to first pretrained protein LLM

Christopher Langmead
Director, Digital Biologics Discovery

AMGEN

NVIDIA DGX Cloud
*AI-training-as-a-service solution*

NVIDIA Base Command Platform
*for workflow management*

NVIDIA AI Enterprise
*RAPIDS for data post-processing*

NVIDIA BioNeMo
*For training and inferencing*

# NVIDIA Clara for Healthcare and Life Sciences

## World's Largest Data Industry | 36% CAGR by 2025

GENOMICS

MEDICAL DEVICES

DRUG DISCOVERY

MEDICAL IMAGING

DIGITAL HEALTH

## NVIDIA CLARA

**PARABRICKS**
Genomics

**ISAAC**
Robotics

**HOLOSCAN**
Instruments

**BIONEMO**
Biomolecules

**MONAI**
Imaging

**NEMO**
Natural Language

# BioNeMo Framework Supports Optimized Biomolecular Models

## Proteins | Small Molecules | Genomics



**ESM-1 | ESM-2**
Protein LLMs



**MegaMolBART**
Generative Chemistry Model



**ProtT5**
Protein Sequence Generation



**DiffDock | EquiDock**
Docking Prediction



**NEW: OpenFold**
3D Protein Structure Prediction



**NEW: DNABERT**
DNA Sequence Model



**NEW: MolMIM**
Molecular Generation



**BETA: Geneformer**
Single Cell Expression Model

NVIDIA

# Optimizing OpenFold Training for Drug Discovery

6x performance improvement in MLPerf HPC v3.0 Benchmark over baseline



MLPerf HPC - OpenFold Training

- Training time to reach 0.9 lDDT-Cα
  - AlphaFold2: 7 days
  - 1056 H100s: 12.4 hrs
  - 2080 H100s: 10 hrs

- MLPerf HPC v3.0 benchmark results
  - OpenFold partial training task finished in 7.51 min, **6x** faster than baseline

# Build Generative AI Virtual Screening Workflows with NVIDIA NIM

## Use composable NVIDIA NIMS to build workflows for CADD applications

# NVIDIA Parabricks for Alignment & Variant Calling

## Speed, Scale, Accuracy

**Alignment**
- BWA-MEM
- Minimap2
- STAR

FastQ →

**Gold Standard Processing and Quality Control**
- Sort BAM
- Mark duplicates
- BQSR
- BAM Metrics
- Collect multiple metrics
- BAM2FQ

BAM/CRAM →

**High-Accuracy Variant Calling**
- DeepVariant
- Mutect2
- GenotypeGVCF
- HaplotypeCaller
- STAR-Fusion
- IndexGVCF

BAM/CRAM →

VCF/gVCF

**Universal Analysis**
Industry-standard tools for all major sequencers, ported to GPU

**Up to 100x Acceleration**
Up to 100x faster for WGS compared to CPU-only

**Up to 50% Lower Cost**
Up to 50% lower compute cost for WGS compared to CPU-only

**Higher Accuracy with AI**
The power of deep learning for customized high accuracy analysis

# A Universal Analysis Solution

## Short-Read



## Long-Read

# Germline Analysis from 18 hours to 10 minutes

## 108x Acceleration using H100s Dynamic Programming Core



**Germline workflow runtime per whole genome**
(HG002, 30x Illumina)

# NATIONAL BIOBANK OF THAILAND ACCELERATES GENOMIC ANALYSIS BY 30X

## Challenge

The National Biobank of Thailand (NBT) is the leading HPC facility and computational science R&D center in the ASEAN region.

Tasked with analyzing massive genome sequencing data from over **50,0000** individuals.

Their goal was to perform **whole genome sequencing (WGS)** to help accurately identify causative mutations and rare variants.

## Solution

NBT leveraged **NVIDIA Parabricks** for genomic analysis on NVIDIA DGX A100, processing 5 PB of data in parallel with speed and accuracy.

The solution accelerated genomic analysis from 30 hours per individual to 1-2 hours.

NBT was able to reduce the whole genome sequencing (WGS) by 4 months, leading to faster genomic discoveries. NBT continues to use DGX A100 for their AI related projects.

### NVIDIA DGX A100
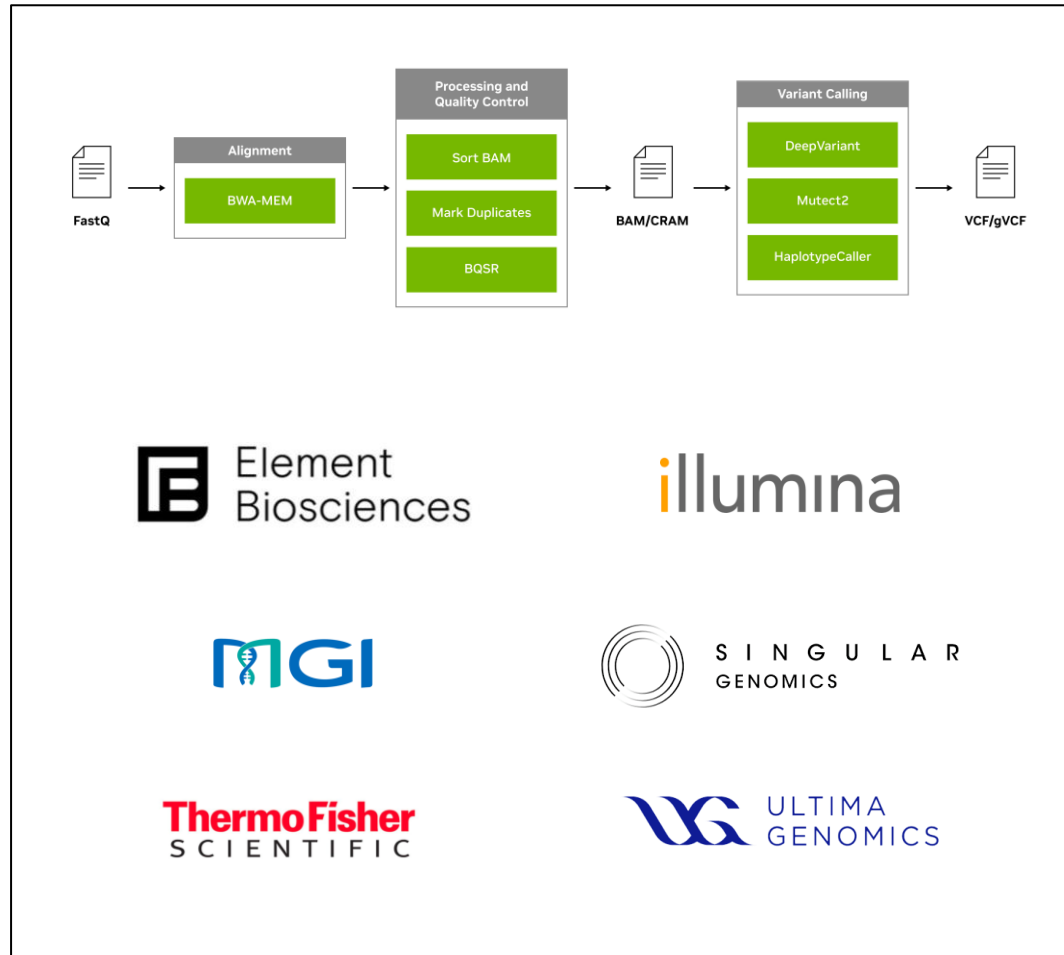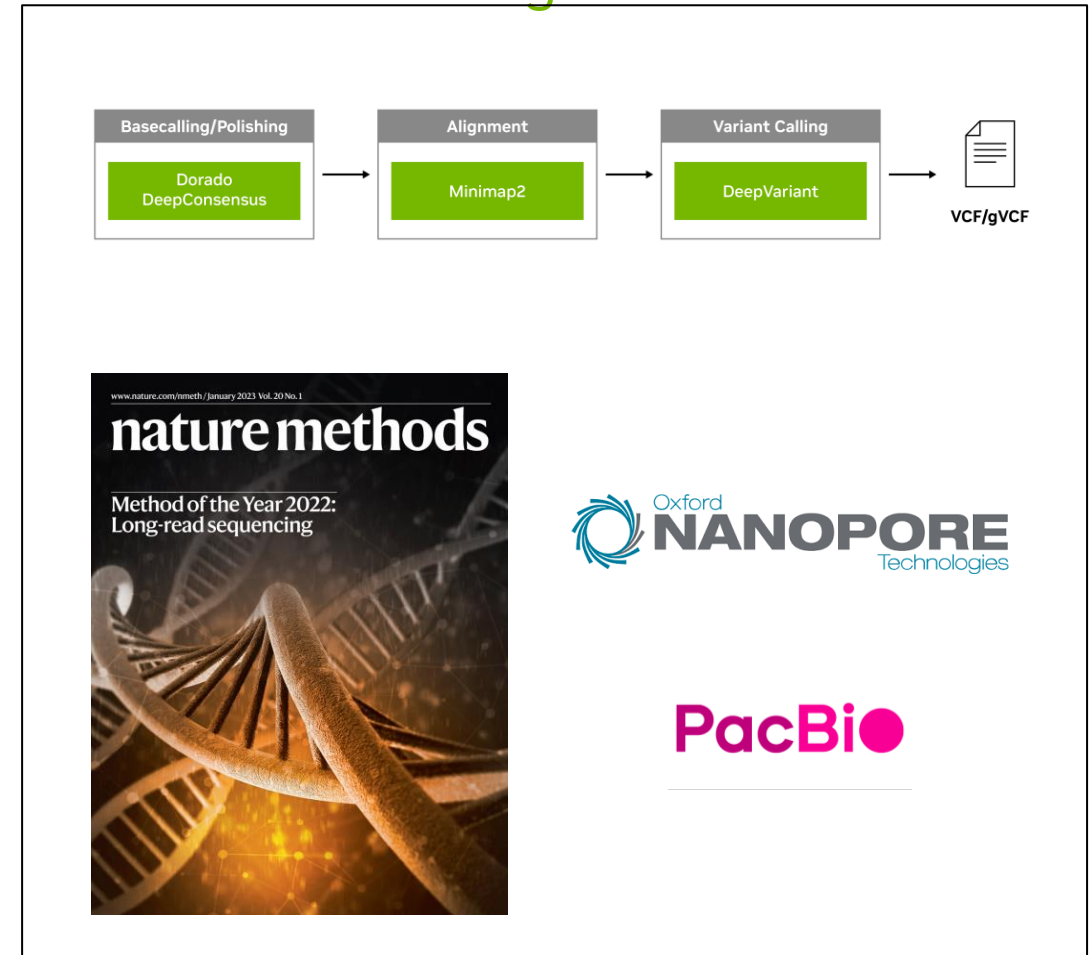*Unprecedented compute performance in the world's first 5 petaFLOPS AI system*

### NVIDIA Parabricks
*Computational genomic analysis framework supporting DNA & RNA*

**93%** Reduction in WGS data processing time per individual

**30x** Faster genomic analysis vs CPU

NVIDIA.

# NEW: Single Cell & Spatial Omics Workflow

Build generative foundation models | Segment at high accuracy | Extract morphology embeddings

# BioNeMo Microservices Activates Partner Ecosystem

Hosting Partners building models with BioNeMo and contributing as NVIDIA NIMS



**Recursion**

Phenom-Beta model for cell morphology

- Phenom-Beta: First vision transformer model targeting cellular data
  - RxRx3 dataset: 17,063 CRISPR-KO genes, 2.2M HUVEC cell images, 1674 compounds, 8 dilutions

- BioHive-1 Supercomputer with NVIDIA DGX SuperPod reference architecture
  - 500 NVIDIA H100 TensorCore GPUs

# NVIDIA NIM – A New Layer on the NVIDIA Clara Stack

## IMAGING

VISTA-3D

VISTA-2D

MAISI

## DRUG DISCOVERY

Phenom-Beta

ESM-2

ESM-Fold

AlphaFold-2

OpenFold

ProtGPT2

NeuralPlexer

DiffDock

MolMIM

MoFlow

## GENOMICS

Deep Variant

FQ2BAM

## DIGITAL HUMAN

Audio2Face

Riva ASR

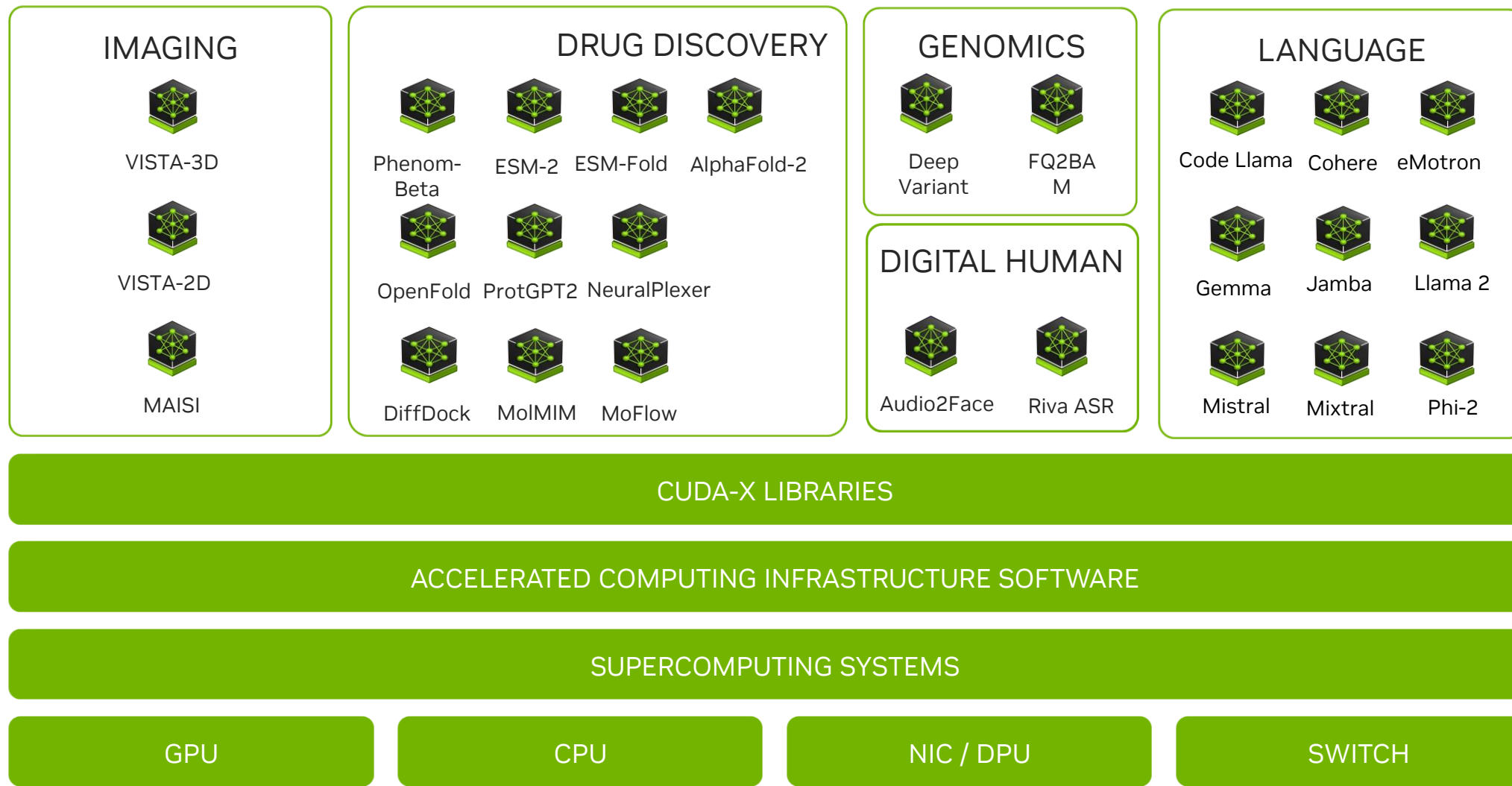## LANGUAGE

Code Llama

Cohere

eMotron

Gemma

Jamba

Llama 2

Mistral

Mixtral

Phi-2

---

**CUDA-X LIBRARIES**

**ACCELERATED COMPUTING INFRASTRUCTURE SOFTWARE**

**SUPERCOMPUTING SYSTEMS**

| GPU | CPU | NIC / DPU | SWITCH |

NVIDIA

# Summary

## Start accelerating your biomedical research with BioNeMo & Parabricks



- Technology is reshaping biomedical research

- BioNeMo provides a suite to tools for DNA, proteins, small molecules, and single-cell / spatial omics analysis

- Parabricks accelerates genome sequencing analysis to <10 min / WGS

- Join world-class leaders like Genentech, Amgen, National BioBank of Thailand in the AI accelerated biodiscovery journey!