

Ai DAY

BIONTECH |  InstaDeep™

This Slide Presentation Includes Forward-Looking Statements

This presentation contains forward-looking statements within the meaning of the Private Securities Litigation Reform Act of 1995, as amended. In some cases, forward-looking statements can be identified by terminology such as “will,” “may,” “should,” “expects,” “intends,” “plans,” “aims,” “anticipates,” “believes,” “estimates,” “predicts,” “potential,” “continue,” or the negative of these terms or other comparable terminology, although not all forward-looking statements contain these words. The forward-looking statements in this presentation are neither promises nor guarantees, and you should not place undue reliance on these forward-looking statements because they involve known and unknown risks, uncertainties, and other factors, many of which are beyond BioNTech’s control and which could cause actual results to differ materially from those expressed or implied by these forward-looking statements. You should review the risks and uncertainties described under the heading “Risk Factors” in BioNTech’s Quarterly Report on Form 6-K for the period ended June 30, 2024 and in subsequent filings made by BioNTech with the SEC, which are available on the SEC’s website at <https://www.sec.gov/>. Except as required by law, BioNTech disclaims any intention or responsibility for updating or revising any forward-looking statements contained in this presentation in the event of new information, future developments or otherwise. These forward-looking statements are based on BioNTech’s current expectations and speak only as of the date hereof.

Furthermore, certain statements contained in this presentation relate to or are based on studies, publications, surveys and other data obtained from third-party sources and BioNTech’s own internal estimates and research. While BioNTech believes these third-party sources to be reliable as of the date of this presentation, it has not independently verified, and makes no representation as to the adequacy, fairness, accuracy or completeness of, any information obtained from third-party sources. In addition, any market data included in this presentation involves assumptions and limitations, and there can be no guarantee as to the accuracy or reliability of such assumptions. While BioNTech believes its own internal research is reliable, such research has not been verified by any independent source. In addition, BioNTech is the owner of various trademarks, trade names and service marks that may appear in this presentation. Certain other trademarks, trade names and service marks appearing in this presentation are the property of third parties. Solely for convenience, the trademarks and trade names in this presentation may be referred to without the ® and TM symbols, but such references should not be construed as any indicator that their respective owners will not assert, to the fullest extent under applicable law, their rights thereto.

Agenda

Introduction and Vision

14:00 Welcome & Introductory Remarks

14:05 Our Vision for AI

Part I. Scaling AI Capabilities

14:10 Computing Infrastructure

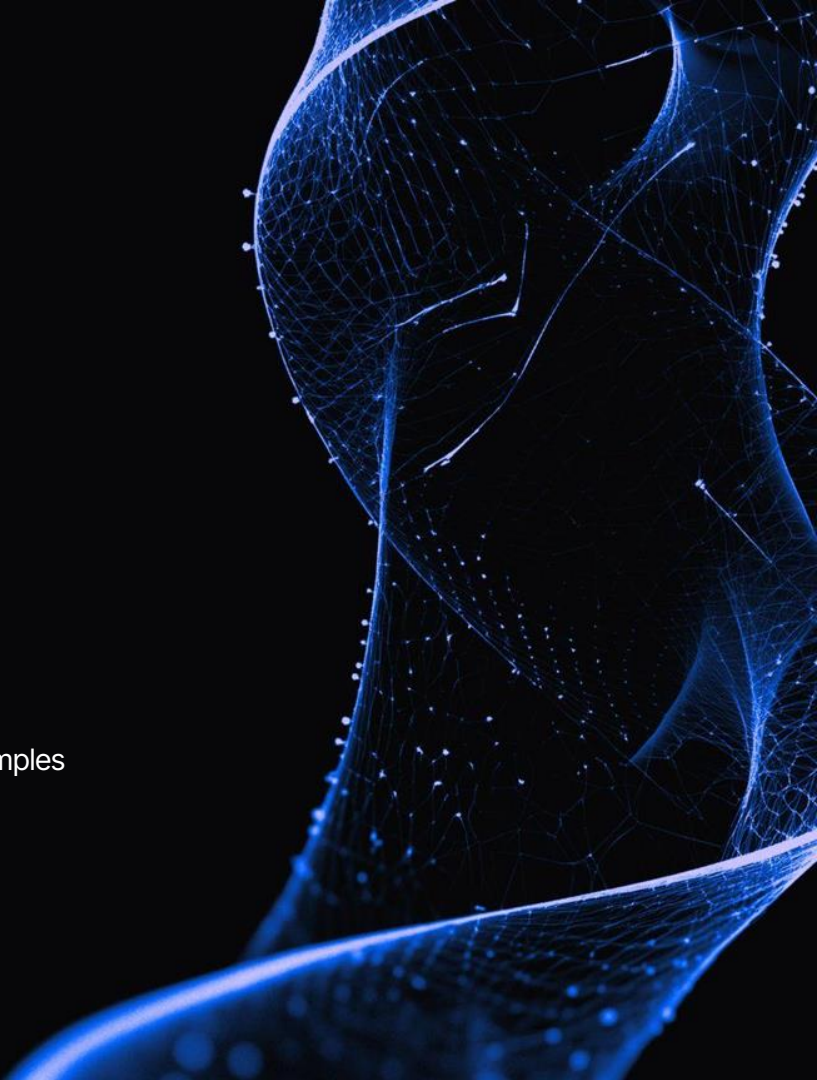
14:25 Innovation: Bayesian Flow Networks

14:45 DeepChain: One Platform, Multiple Tools

Part II. Deploying AI across the pipeline

15:00 Applying AI end-to-end to the immunotherapy pipeline: examples

15:40 Closing Remarks and Q&A



Introduction and Welcome



Ugur Sahin
Founder & CEO
BioNTech



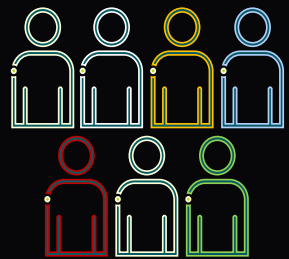
Ryan Richardson
Chief Strategy Officer
BioNTech



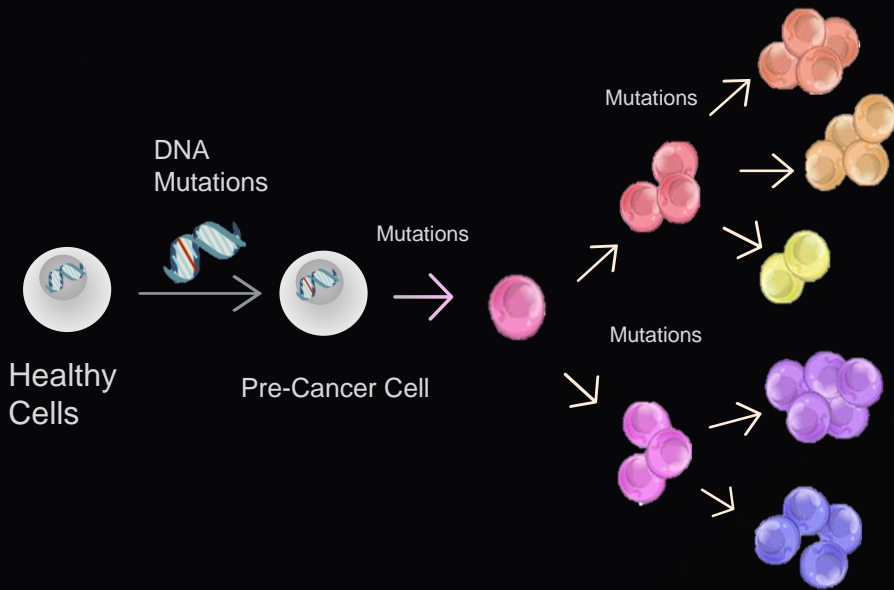
Karim Beguir
CEO
InstaDeep

Root Cause of Cancer Treatment Failure

Interindividual Variability & Intratumoral Heterogeneity



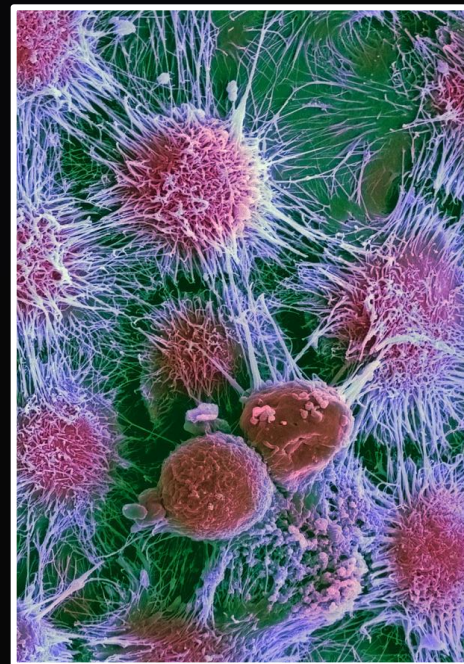
Individual Patients



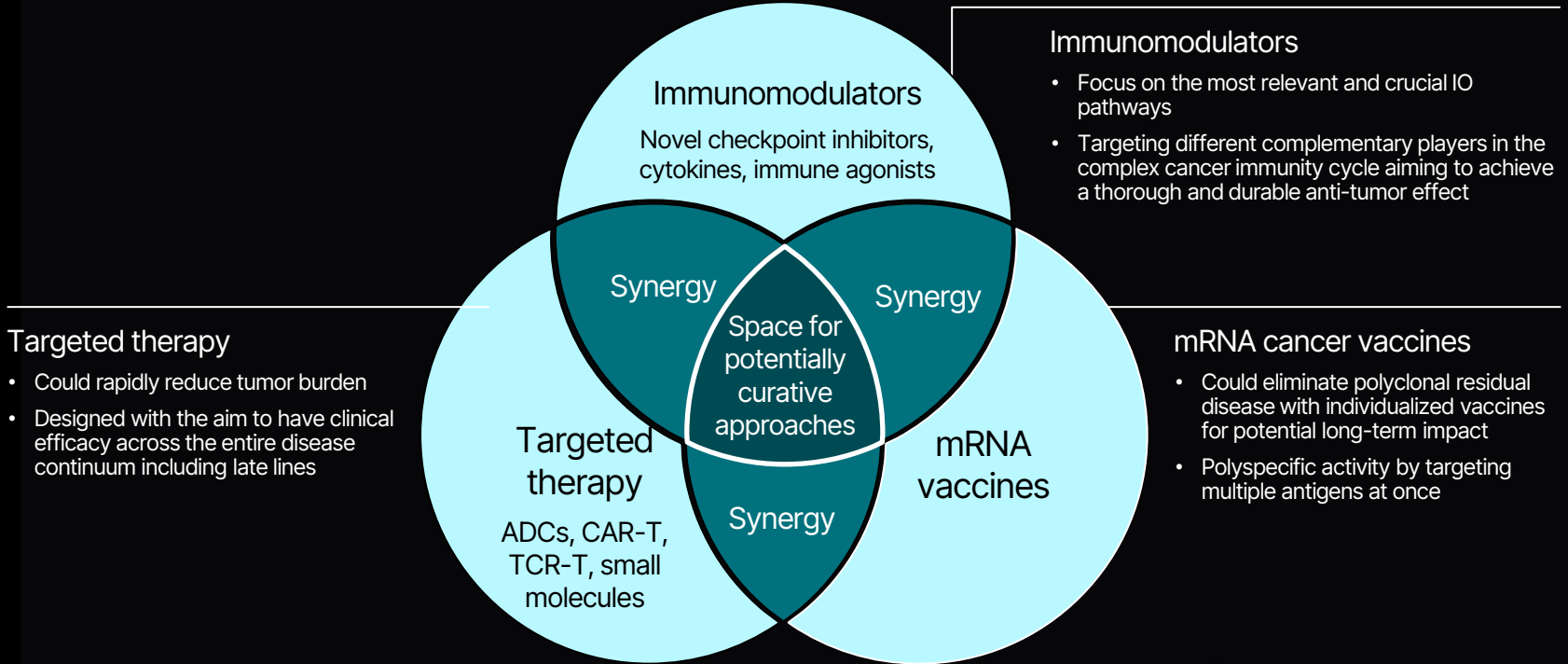
Cancer Evolution 5-20 Years – up to 10,000 Mutations

Cancer Cells

Genetically Diverse & Adaptable





Towards a Potentially Curative Approach to Cancer Based on Multiple Modalities and Differentiated Novel/Novel Therapeutic Combinations





ADC = antibody-drug conjugate; CAR = chimeric antigen receptor; TCR-T = T-cell receptor engineered T cell; IO = immune oncology.

Charting the Course for Tomorrow's Personalized Medicine

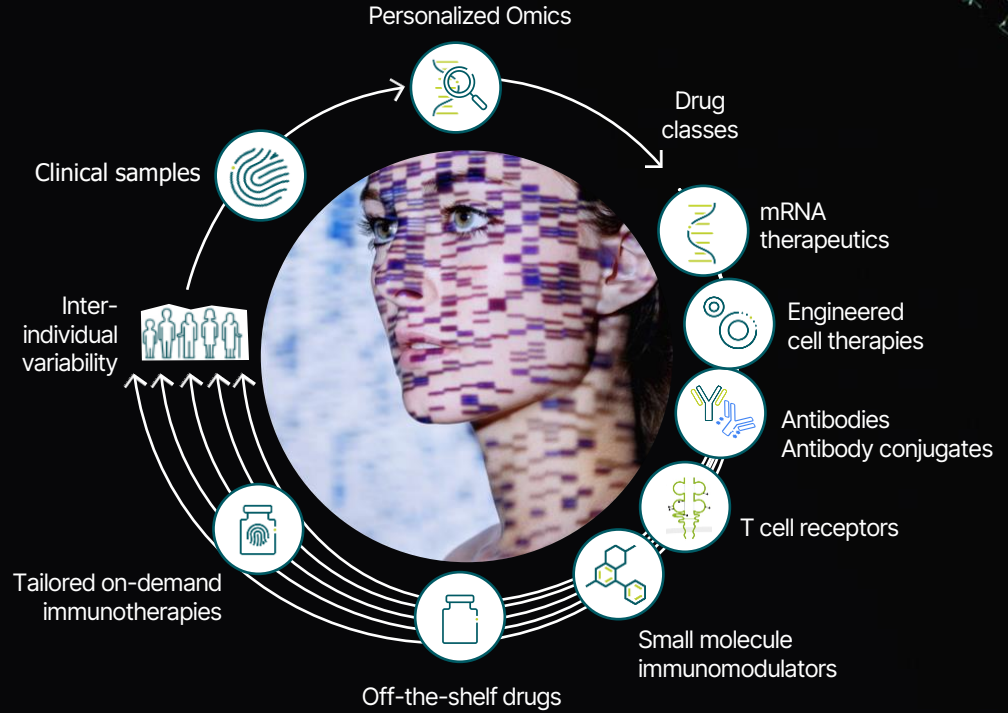
- 

Deep genomics & immunology expertise to analyse patient data
- 

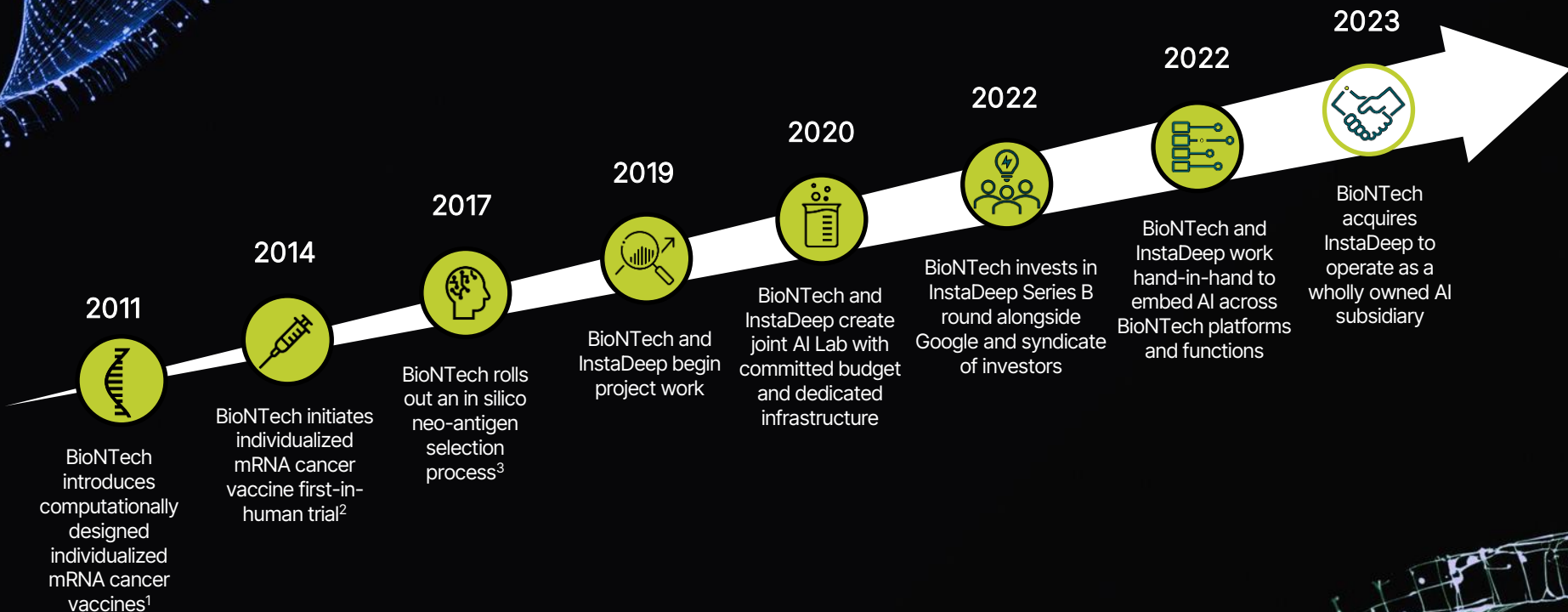
Individualized treatment platforms to address inter-individual variability
- 

AI & digitally-integrated target & drug discovery and development
- 

In-house manufacturing









BioNTech and InstaDeep – The Road to Partnership



1. Cancer Res, PMID 22237626, 2. Nature PMID 25901682, Nature PMID 28678784, 3. BioNTech's personalized cancer vaccine candidate, autogene cevumeran, is partnered with Genentech, a member of Roche Group.

Two Companies: One Mission

<p>BioNTech: >6,800¹ Employees</p> 	<p>HQ: Mainz, Germany</p> 	<p>Developing medicines to fight cancer, infectious diseases and other serious diseases.</p> 
<p>InstaDeep: >370¹ Employees</p> 	<p>HQ: London, UK</p> 	<p>Focused on productizing disruptive AI innovation</p> 

Our Goal: Building a leading AI-first, personalized immunotherapy platform

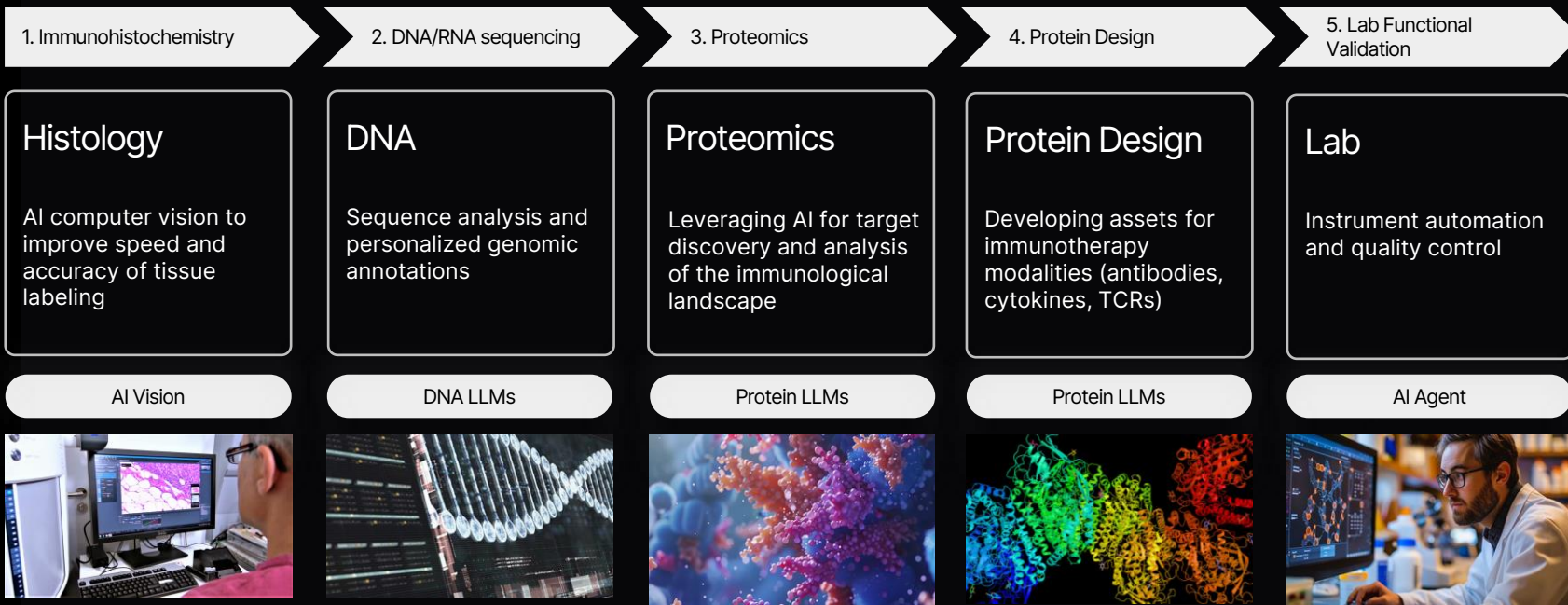
(and leveraging the breakthroughs obtained in the process)

¹. As of 30 June 2024



Our Goal: Deploying AI end-to-end in our immunotherapy pipeline

AI-first Immunotherapy Platform



Our Approach:

1. Scaling AI capabilities
2. Deploying across the pipeline

Part I.

Scaling AI Capabilities





Computing Infrastructure



InstaDeep's Supercomputing Cluster

Cluster Specifications

224 Nvidia H100 GPUs

86,000 CPU Cores

1.7 PetaBytes persistent storage

400 Gbps RoCE network

Our Supercomputing Cluster Is Nearing Exascale Levels:

InstaDeep on-premise Cluster totals to ~0.5 ExaFLOPS

Top 100 worldwide [1]

Top 20 H100 GPU clusters worldwide [2]

[1] "Top 500, The List", June 2023

[2] "State of AI Report Compute Index", August 2024

Advanced In-House Rack Design

Easy to expand with modular nodes

Consistent performance, cost, power, cooling

Optimized for large-scale AI workloads

Simplified management with consistent architecture

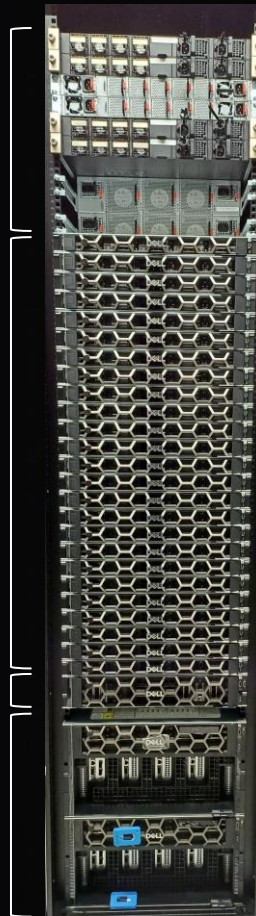
Minimize expenses with standard design

Network stack
400 Gbps

CPU Nodes
6144 CPUs

Fast Storage
122 TB

GPU Nodes
16 H100 GPUs



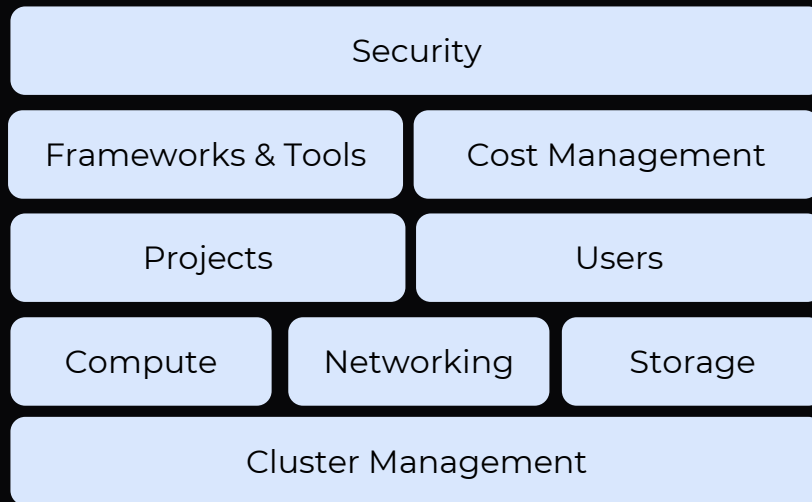
Software Stack Supporting our Cluster:



Fully tailored AI stack from hardware to experiments

Open standards

Cutting-edge tooling



Strategic Benefits from our SuperComputing Cluster

Availability when most needed

Flexibility on Sw/Hw Integration

No Vendor Locking

Repeatable design

Predictable costs

Cost efficient (**50% savings** on cloud equivalent at 60% usage)

Scaling Intelligence

Scaling Intelligence

Why

Scaling
Laws

How

Engineering
Expertise

What

Accelerating
Scientific Discovery

Scaling Intelligence

Why

Scaling
Laws

How

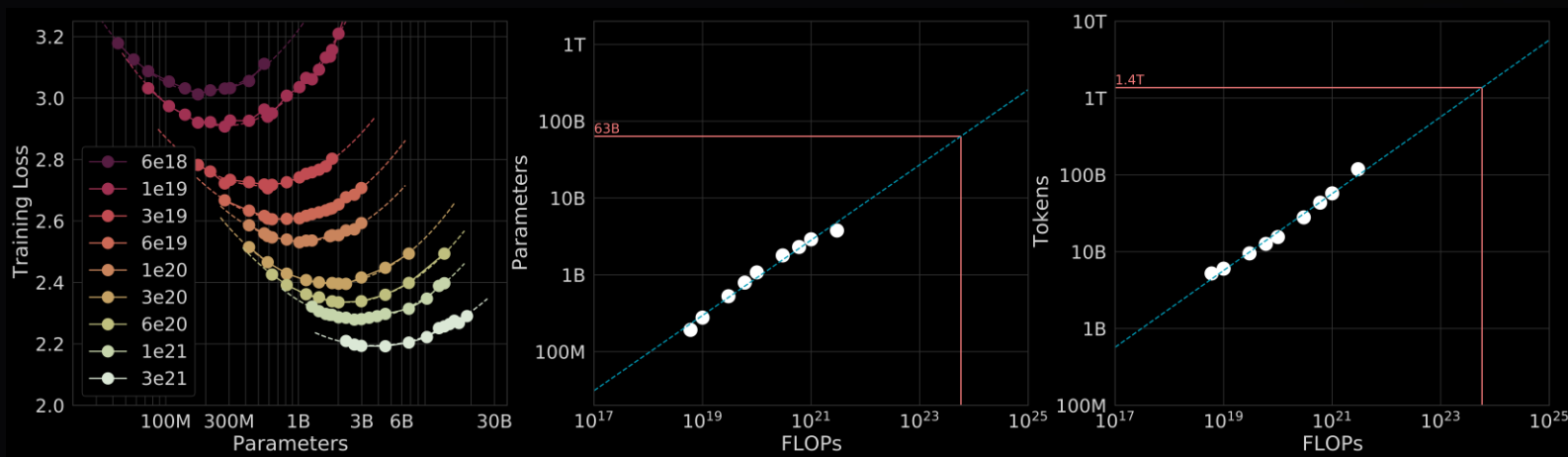
Engineering
Expertise

What

Accelerating
Scientific Discovery

Scaling Laws

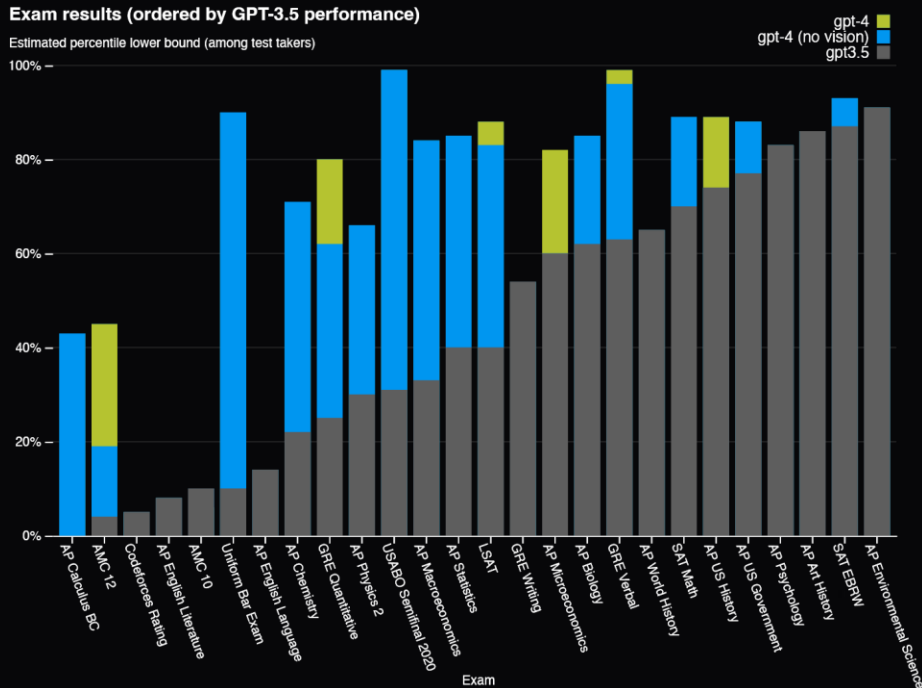
Performance of Large Language Models (LLMs) is a smooth, well-behaved and predictable function of the **number of parameters** of your model, the **amount of data** used to train it, and **computing resources**.



Source: Hoffmann, J., Borgeaud, S., Mensch, A., Buchatskaya, E., Cai, T., Rutherford, E., Casas, D.D.L., Hendricks, L.A., Welbl, J., Clark, A. and Hennigan, T., 2022. Training compute-optimal large language models.

Scaling Laws

We can expect
“more intelligence”
 by scaling existing
 algorithms.



Source: Achiam, Josh, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida et al. "Gpt-4 technical report." (2023)

Scaling Intelligence

Why

Scaling
Laws

How

Engineering
Expertise

What

Accelerating
Scientific Discovery

How to Scale Next-Generation Foundation Models

Scaling next-generation AI systems demands advanced engineering solutions, tightly with the hardware to balance training and deployment constraints.

Memory

Model Sharding
Rematerialization
Quantization / Precision

Network

Compute/Comm. overlap
I/O and Data processing
Hardware and Topology

Compute

XLA optimization
Kernel Fusion and Caching
Data Parallelism

Scaling Intelligence

Why

Scaling
Laws

How

Engineering
Expertise

What

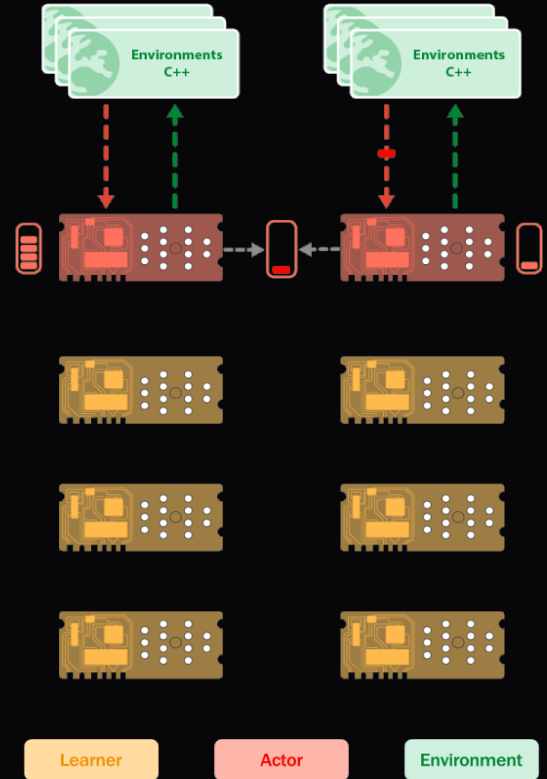
Accelerating
Scientific Discovery

#1 Accelerating Reinforcement Learning

Reinforcement Learning is the science of **learning from trials and errors**.
 A simulation engine **turns computation into data**.

Scaling Reinforcement Learning

- Multiple threads keep the **hardware accelerators** active.
- Learner cores process experience, synchronizing updates using **JAX primitives**.
- The architecture can be **replicated across a large number of nodes** to form a supercomputing cluster.
- Leverage **the high-speed inter chip interconnects** between the nodes of the hardware accelerators.



The Sebulba Architecture on 8x hardware accelerators¹

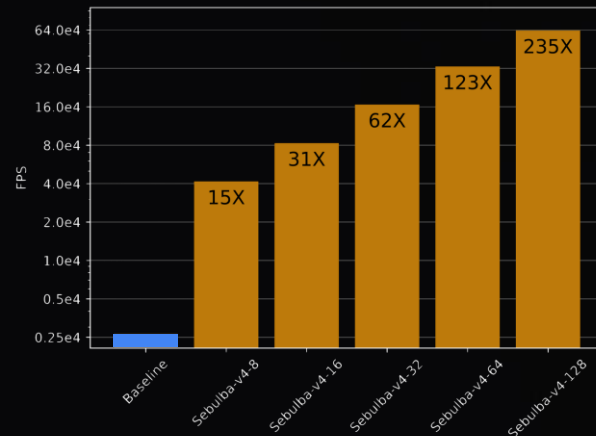
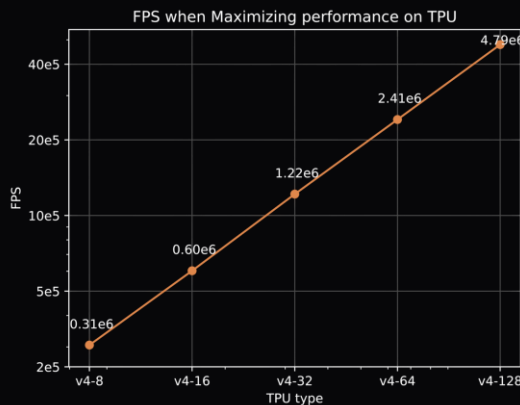
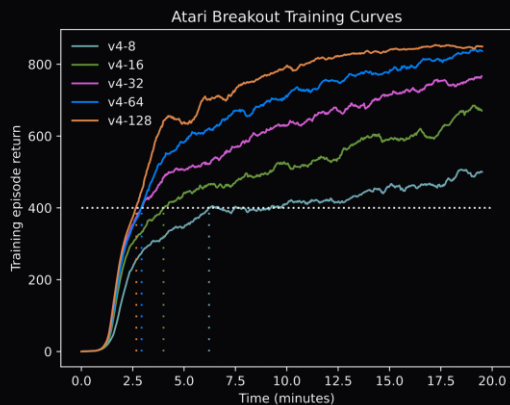
[1] "InstaDeep's scalable reinforcement learning on Cloud TPU", October 19, 2023, Google Cloud blog post
 [2] Berner, C., Brockman, G., Chan, B., Cheung, V., Debiak, P., Dennison, C., Farhi, D., Fischer, Q., Hashme, S., Hesse, C. and Józefowicz, R., 2019. Dota 2 with large scale deep reinforcement learning.
 [3] Hessel, M., Kroiss, M., Clark, A., Kemaev, I., Quan, J., Keck, T., Viola, F. and van Hasselt, H., 2021. Podracer architectures for scalable reinforcement learning.

#1 Accelerating Reinforcement Learning

Better: 50% improvement in performance as we scale the hardware and simulated data.

Cheaper: 13x cost reduction due to the more efficient use of the hardware.

Faster: 240x faster to train an RL agent up to convergence.



[1] "InstaDeep's scalable reinforcement learning on Cloud TPU", October 19, 2023, Google Cloud blog post

[2] Berner, C., Brockman, G., Chan, B., Cheung, V., Debiak, P., Dennison, C., Farhi, D., Fischer, Q., Hashme, S., Hesse, C. and Józefowicz, R., 2019. Dota 2 with large scale deep reinforcement learning.

[3] Hessel, M., Kroiss, M., Clark, A., Kemaev, I., Quan, J., Keck, T., Viola, F. and van Hasselt, H., 2021. Podracer architectures for scalable reinforcement learning.

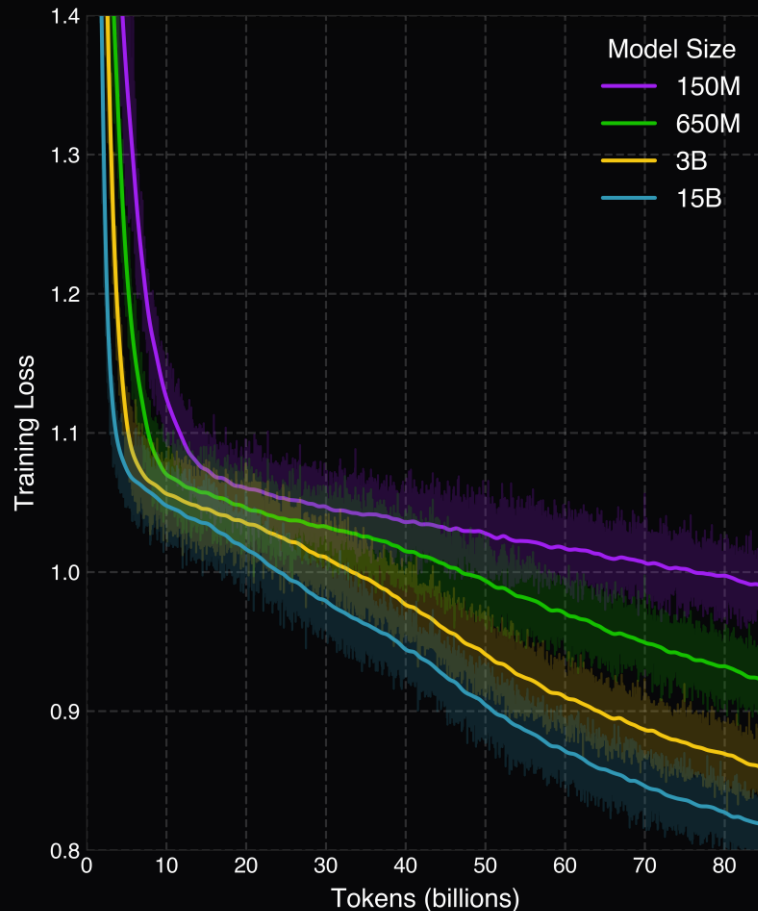
#2 Generative AI for Biology

Scaling the Next-Generation of our Generative AI Models

- in-house JAX-based software library
- best-in-class engineering for scaling LLMs
e.g. Hybrid Parallelism, mixed-precision, rematerialization, etc.

Results

- Train multi-billion parameters models (+15B)
- [Scaling Laws in Action](#)
- Hardware efficiency on par with the latest Meta Llama 3.1¹¹
i.e. Model Flop Utilization of ~50% for our 15B model



^[1] Dubey, A., Jauhri, A., Pandey, A., Kadian, A., Al-Dahle, A., Letman, A., Mathur, A., Schelten, A., Yang, A., Fan, A. and Goyal, A., 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.

Source: InstaDeep

Summary

InstaDeep's Supercomputing Cluster and advanced software stack could facilitate new scientific breakthroughs, services and products that were previously out of reach.

Scaling Laws

Engineering Expertise

Accelerating Scientific Discoveries

AI Innovation

Bayesian Flow Networks

Generative AI



Sora: generative AI
model for videos

OpenAI 2024

A stylish woman walks down a Tokyo street filled with warm glowing neon and animated city signage...

Joint Probability Distributions

For a generative model of face images, creating a new face means picking a **sample** from the **joint probability distribution** of all the pixels.

Q: Why is this so hard?

A: Because all the pixels are **interrelated**

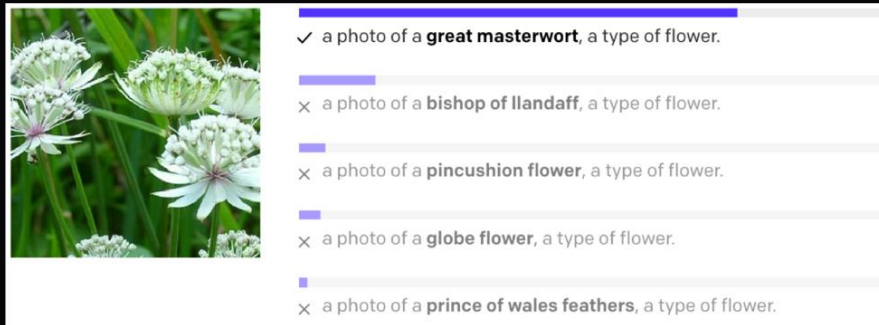


Generating Diverse High-Fidelity Images with VQ-VAE-2 (Razavi et al. 2019)

Steerable Generation

We want to **control** what we generate.

With a **multimodal model** (e.g. images and text) we can do this with **conditional sampling**: **fix** one modality, **generate** another.



Given an image, generate a caption (CLIP OpenAI 2021)

an armchair in the shape of an avocado. . . .



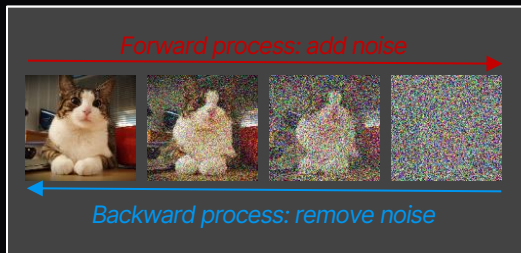
Given an prompt, generate an image (DALL·E OpenAI 2021)

One Model, Many Tasks

Learning a **joint distribution** over **many variables** then choosing which to **fix** and which to **generate** gives us **one model for many tasks**

But which Model?

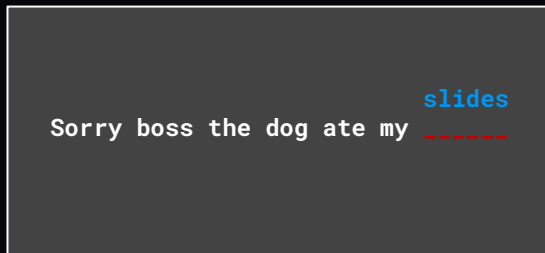
Diffusion



Pros: Continuous data (especially images), inpainting, fast gradient-based sampling

Cons: Discrete data

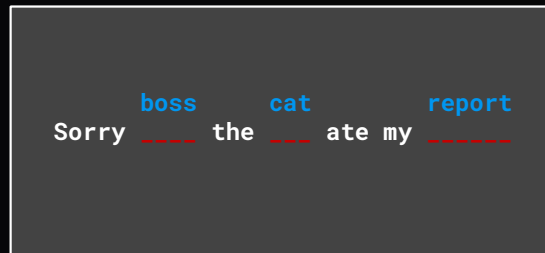
Autoregression (GPT)



Pros: Sequence data (especially text)

Cons: Unordered data, inpainting, slow sampling

Masked prediction (BERT)



Pros: Discrete data, inpainting, representation learning

Cons: Continuous data, slow sampling

Bayesian Flow Networks



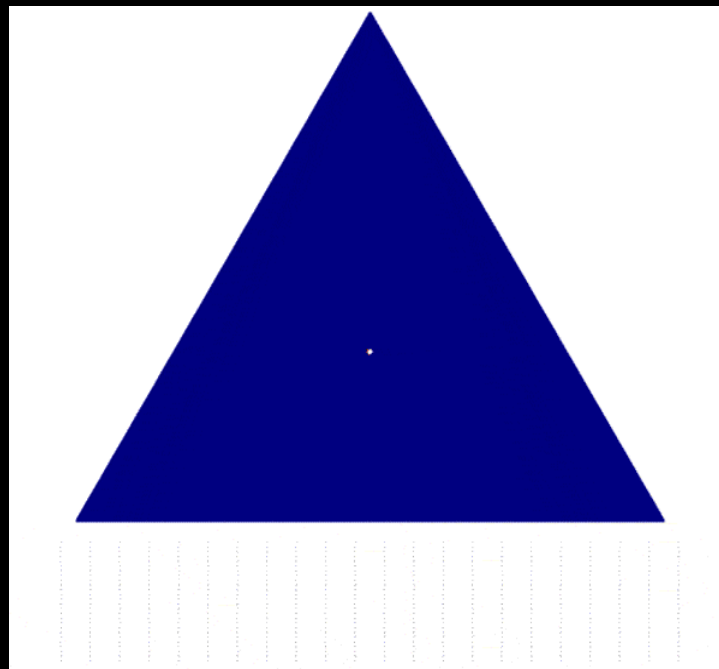
BFNs are a new class of generative model that uses **Bayesian inference** to update **beliefs** about data.



Unlike diffusion models, they generate **discrete data** in a **continuous way**, allowing for **gradient-based sampling**



This makes BFNs well-suited for **controllable generation** across **diverse data modalities**.



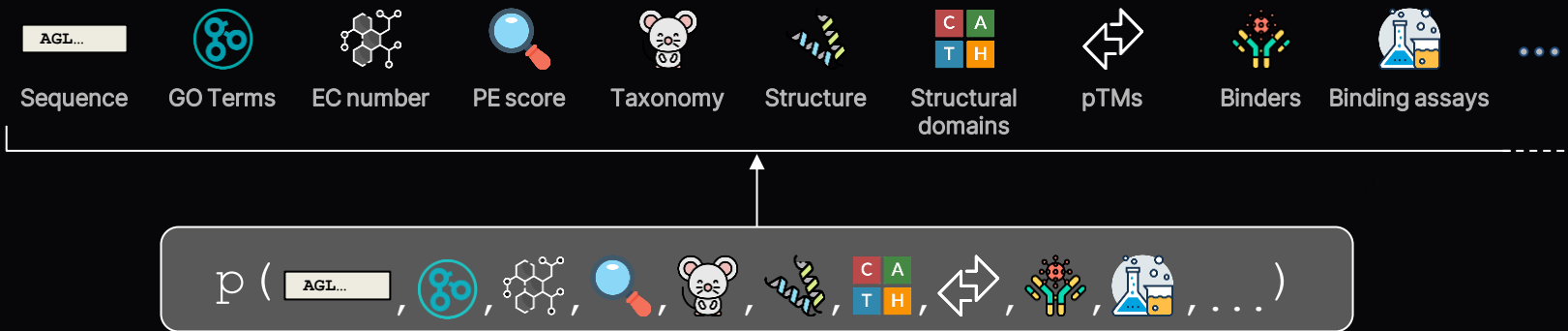
Graves, Srivastava, Atkinson, Gomez 2023

Generative Modelling

A **unifying framework** to learn useful functions from data

(1) Learn to model the joint distribution of *all* your data, and (2) conditionally sample for tasks of interest.

Proteomics: Joint Modelling

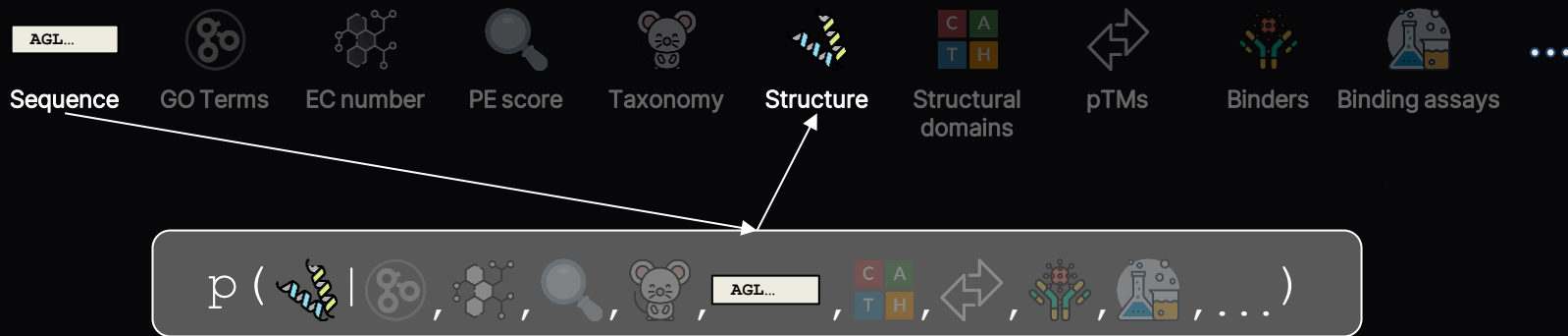


Generative Modelling

A **unifying framework** to learn useful functions from data

(1) Learn to model the joint distribution of *all your data*, and (2) conditionally sample for tasks of interest.

Proteomics: Protein Folding

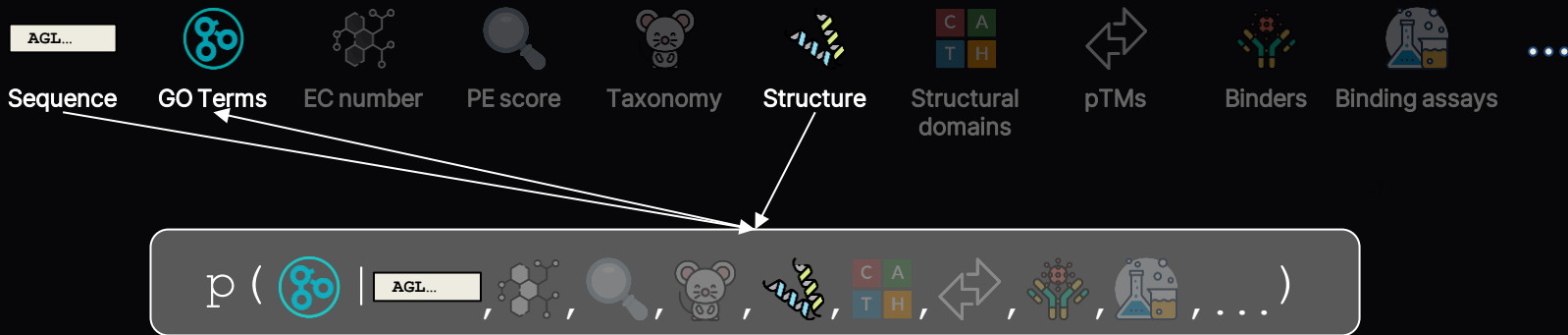


Generative Modelling

A **unifying framework** to learn useful functions from data

(1) Learn to model the joint distribution of *all* your data, and (2) conditionally sample for tasks of interest.

Proteomics: Function Prediction

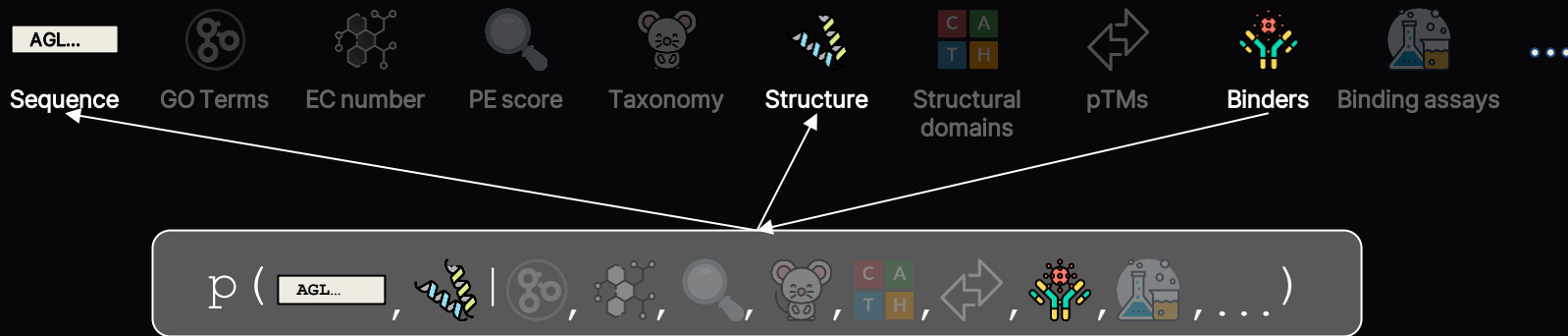


Generative Modelling

A **unifying framework** to learn useful functions from data

(1) Learn to model the joint distribution of *all* your data, and (2) conditionally sample for tasks of interest.

Proteomics: Antibody Design



Generative Modelling

A **unifying framework** to learn useful functions from data

(1) Learn to model the joint distribution of *all* your data, and (2) conditionally sample for tasks of interest.

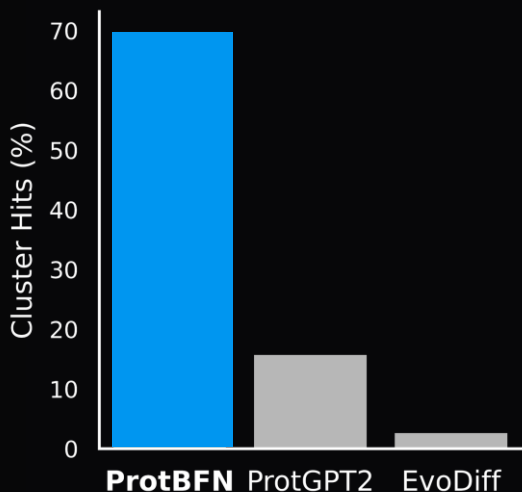
Proteomics: Sequence Generation

			De Novo	Conditional
AGL...				
Sequence	Autoregression (GPT)	R MPPR____...	Yes	Limited
	Masked prediction (BERT)	RRS MPP__IV...	No	Yes
	"Discrete" diffusion	P R I MMPRSSPV...	Limited applicability to discrete data	

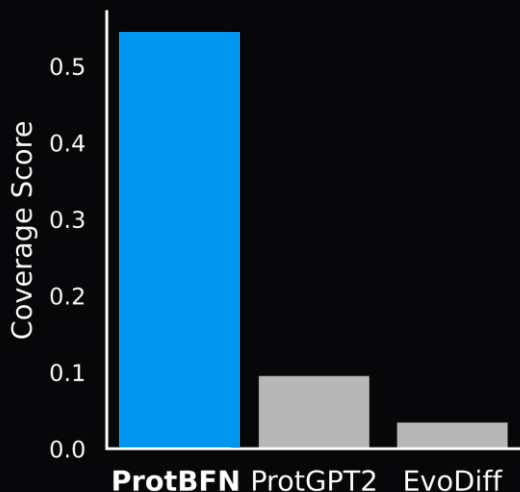
Natural, Diverse & Novel Protein Sequences

ProtBFN learns statistical and biochemical properties of natural proteins with high-fidelity.

More natural...¹



...more diverse...¹



...and highly novel.²



1. 10,000 generated sequences from each model are matched to clusterings from UniRef50. A hit is determined as a match with >50% sequence identity. Coverage score is the ratio of the number of unique clusters hit to the expected number if sequences were drawn i.i.d. from the models training distribution. ProtGPT2 (huggingface.co/nferruz/ProtGPT2) and EvoDiff (github.com/microsoft/evodiff) sequences are sampled using publicly available code and model weights provided by the authors.

2. Identity of ProtBFN generated sequences to the best matching protein sequence found in the models training data. Any identity < 100% is a novel sequence that the model has not seen before.

Globular Structural Motifs With Novel Sequences

Predicted structures of generated sequences show natural, globally coherent and functionally diverse folds.

Structure largely determines function in nature.



Single and multi-domain proteins.

Globally coherent generations with inter-domain interactions.

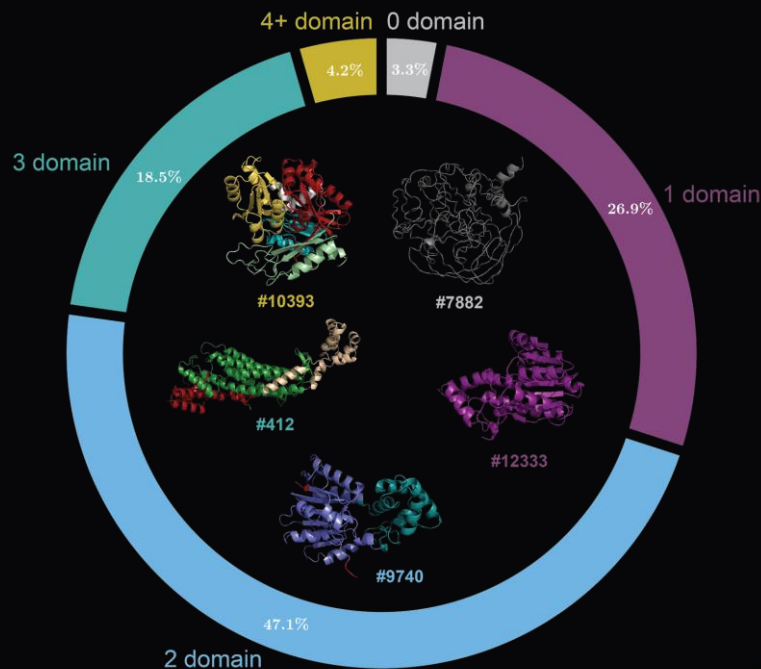
Spans diversity of known structures and tree-of-life.

Alpha Helical, Beta Sheet, Alpha-Beta and Irregular domains.

Small and large domains.

Transmembrane Proteins (porins and transporters) and Enzymes.

Domains specific to Archaea, Bacteria, Eukarya (Plants, Humans).



BFN for Protein Sequences



Outperforms or matches **task-specific autoregressive, diffusion** and **BERT** models.



Improved **naturalness, diversity** and **novelty**.



Uses **zero-shot conditioning** of model.



Patent application filed.



2024-8-23

Protein Sequence Modelling with Bayesian Flow Networks

Timothy Atkinson^{1,1}, Thomas D. Barrett^{2,1},

Scott Cameron¹, Bora Guloglu¹, Matt Greenig¹, Louis Robinson¹, Alex Graves¹, Liviu Copoiu¹ and Alexandre Laterre¹

¹Equal contributions, ¹InstaDeep

Exploring the vast and largely uncharted territory of amino acid sequences is crucial for understanding complex protein functions and the engineering of novel therapeutic proteins. Whilst generative machine learning has advanced protein sequence modelling, no existing approach is proficient for both unconditional and conditional generation. In this work, we propose that Bayesian Flow Networks (BFNs), a recently introduced framework for generative modelling, can address these challenges. We present ProtBFN, a 650M parameter model trained on protein sequences curated from UniProtKB, which generates natural-like, diverse, structurally coherent, and novel protein sequences, significantly outperforming leading autoregressive and discrete diffusion models. Further, we fine-tune ProtBFN on heavy chains from the Observed Antibody Space (OAS) to obtain an antibody-specific model, AbBFN, which we use to evaluate zero-shot conditional generation capabilities. AbBFN is found to be competitive with, or better than, antibody-specific BERT-style models, when applied to predicting individual framework or complimentary determining regions (CDR).

Released only a few days ago!¹

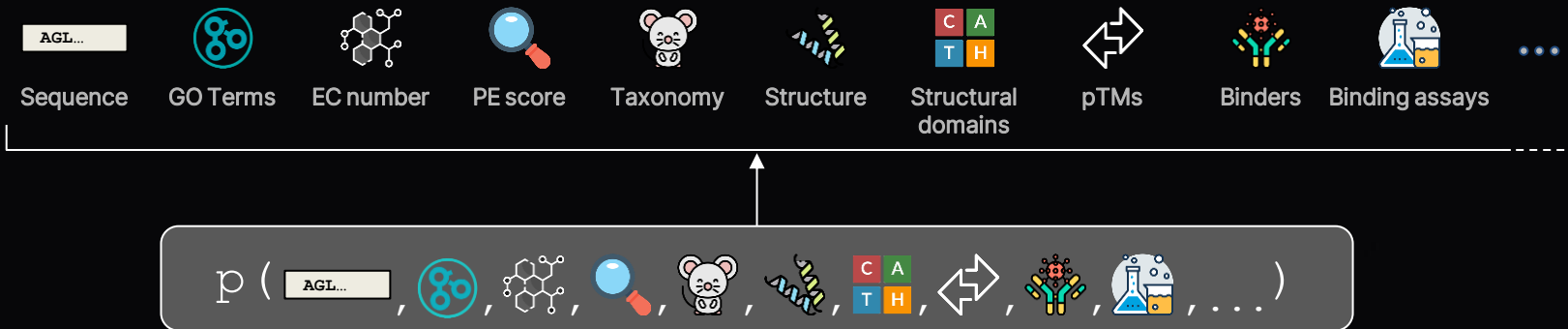
1. Available at <https://www.biorxiv.org/content/10.1101/2024.09.24.614734v1>

Going Beyond Sequence-Only Models

Our goal is to **model everything**: building foundational models of the joint distribution of heterogeneous scientific data.

Performance across multiple data types and sources.

Flexibility in the hands of scientists with task-specific inference.



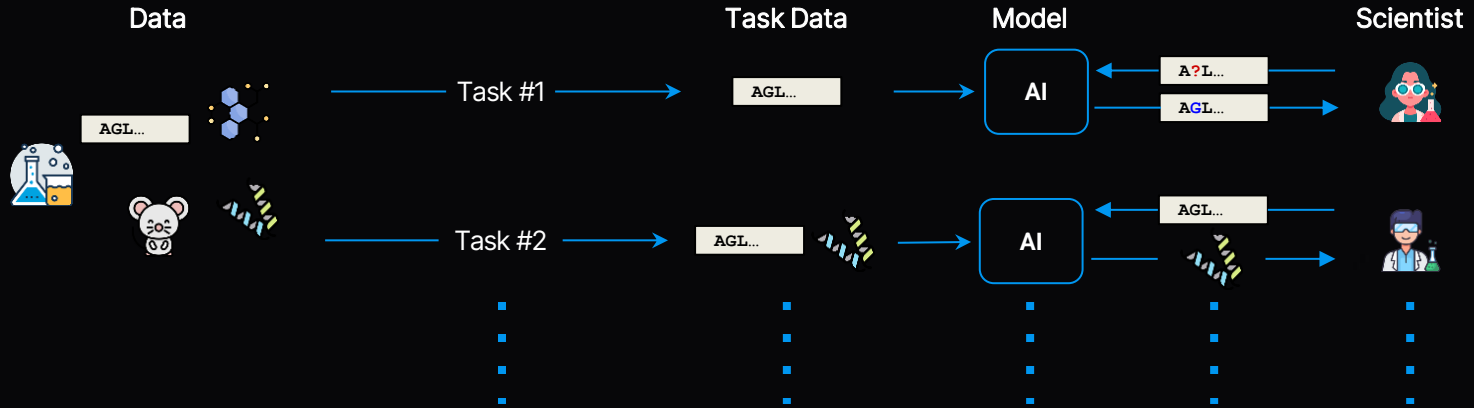
Going Beyond Sequence-Only Models

Our goal is to **model everything**: building foundational models of the joint distribution of heterogeneous scientific data.

Performance across multiple data types and sources.

Flexibility in the hands of scientists with task-specific inference.

Conventional ML



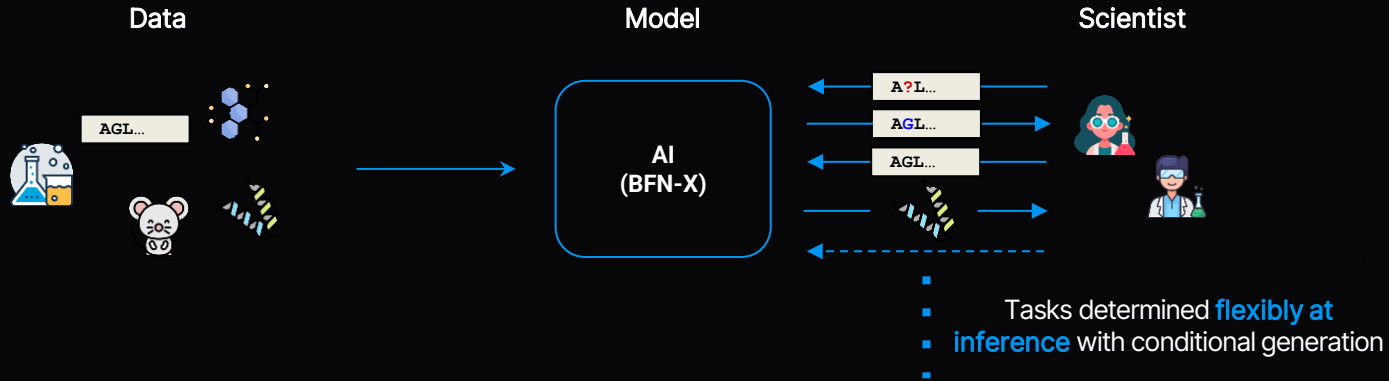
Going Beyond Sequence-Only Models

Our goal is to **model everything**: building foundational models of the joint distribution of heterogeneous scientific data.

Performance across multiple data types and sources.

Flexibility in the hands of scientists with task-specific inference.

Our vision



Introducing AbBFN-X



First look at our **multimodal model** for antibodies.

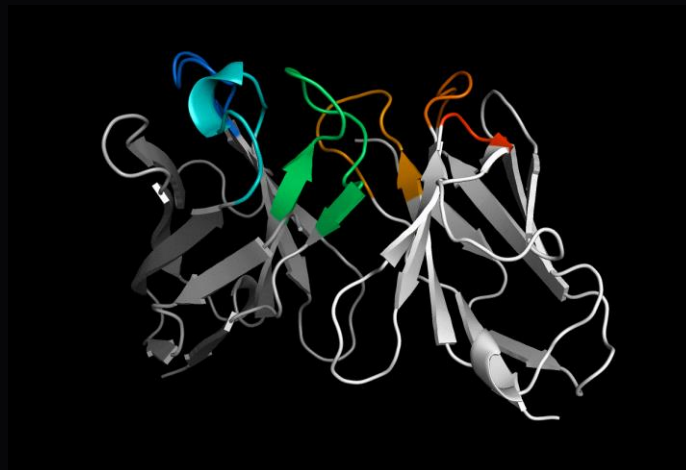
36 different attributes jointly modelled: **sequence, genetic, biophysical**



Empowers scientists with **tunable generation**, being **highly flexible** across many tasks.

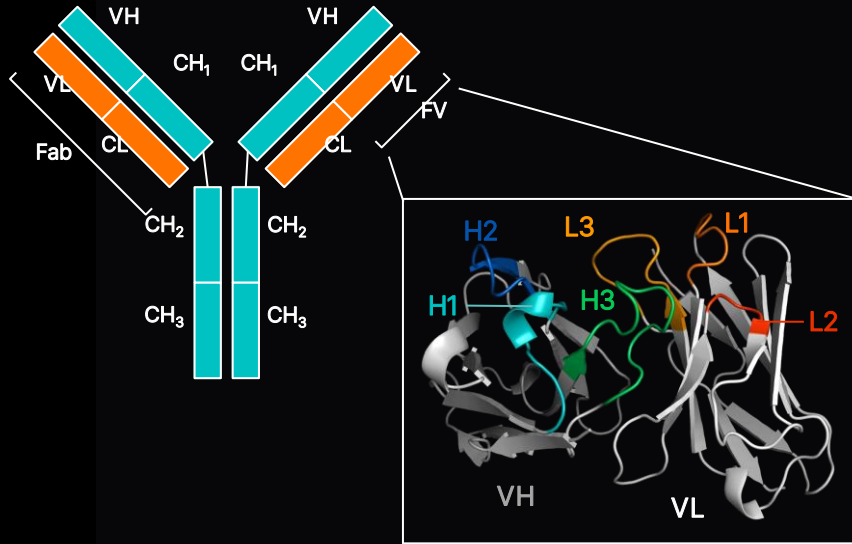


Today's use cases go **beyond standard AI-enhanced antibody design** workflows.



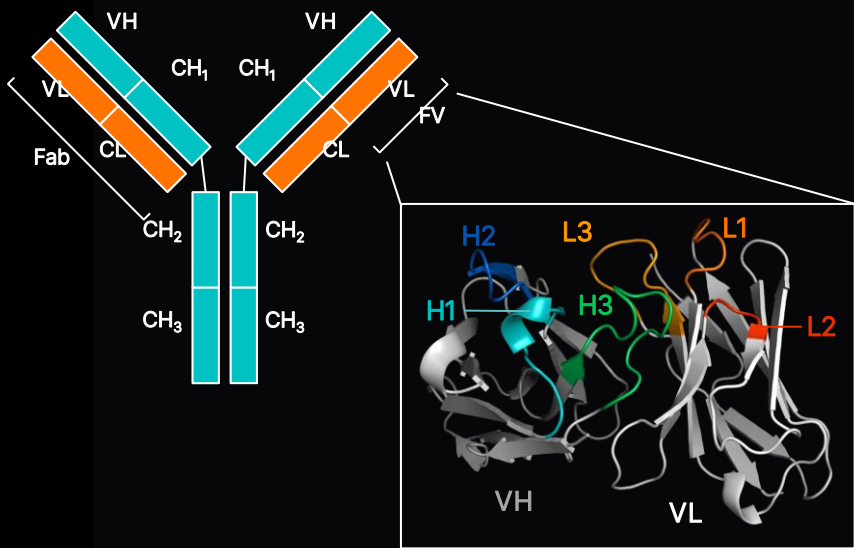
Source: Generated Image

AbBFN-X



CDR-H1 **CDR-H2** **CDR-H3**
VH: EVQLLESGGGLVQPGGSLRLSCAAS **GFTFSSYA**MSWVRQAPGKGLEWWSA **ISWNSGSI**YADSVKGRFTISRDN SKNTLYLQMNSLR AEDTAVYYC **ARGWSQVDTAMDLDY**GQGTLVTVSS
V gene *D gene* *J gene*
CDR-L1 **CDR-L2** **CDR-L3**
VL: DIQMTQSPSSVSASVGD RVTITCRAS **QSVSSN**LAWYQQKPGKAPKLLIY **GAS**SLQSGVPSRFSGSGSGTDFTLTISSLPEDFATYYC **QQYNNWLT**FGQGTRLEIK
V gene *J gene*

AbBFN-X



Amino acid sequence

AGL...	FWR-H1	AGL...	FWR-L1
AGL...	CDR-H1	AGL...	CDR-L1
AGL...	FWR-H2	AGL...	FWR-L2
AGL...	CDR-H2	AGL...	CDR-L2
AGL...	FWR-H3	AGL...	FWR-L3
AGL...	CDR-H3	AGL...	CDR-L3
AGL...	FWR-H4	AGL...	FWR-L4

Genetic Attributes

	HV gene		HV seq. identity
	HD gene		HD seq. identity
	HJ gene		HJ seq. identity
	LV gene		LV seq. identity
	LD gene		LD seq. identity
	LC locus		Species

Biophysical Attributes

	Hydrophobicity
	Positive Patches
	Negative Patches
	Charge Imbalance

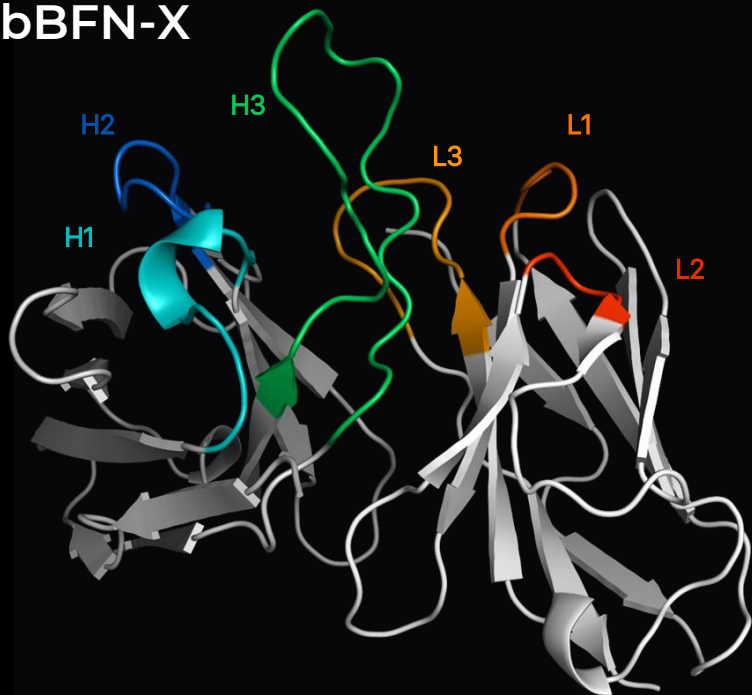
Length Attributes

	CDR-H1 length		CDR-L1 length
	CDR-H2 length		CDR-L2 length
	CDR-H3 length		CDR-L3 length
	VH length		VL length

VH: EVQLLESGGGLVQPGGSLRLSCAAS **QFTFSSYA**MSWVRQAPGKGLEWVSA **ISWNSSGI**YADSVKGRFTISRDNKNTLYLQMNSLRAEDTAVYYC **ARGWSQVDTAMDLDY**GQGTLVTVSS
V gene *D gene* *J gene*

VL: DIQMTQSPSSVSASVGDRTITCRAS **QSVSSN**LAWYQQKPGKAPKLLIY **GAS**SLQSGVPSRFSGSGSDFTLTISSLQPEDFATYYC **QQYNNWLT**FGQGTREIK
V gene *J gene*

AbBFN-X



Amino acid sequence

AGL...	FWR-H1	AGL...	FWR-L1
AGL...	CDR-H1	AGL...	CDR-L1
AGL...	FWR-H2	AGL...	FWR-L2
AGL...	CDR-H2	AGL...	CDR-L2
AGL...	FWR-H3	AGL...	FWR-L3
AGL...	CDR-H3	AGL...	CDR-L3
AGL...	FWR-H4	AGL...	FWR-L4

Genetic Attributes

	HV gene		HV seq. identity
	HD gene		HD seq. identity
	HJ gene		HJ seq. identity
	LV gene		LV seq. identity
	LD gene		LJ seq. identity
	LC locus		Species

Biophysical Attributes

	Hydrophobicity
	Positive Patches
	Negative Patches
	Charge Imbalance

Length Attributes

	CDR-H1 length		CDR-L1 length
	CDR-H2 length		CDR-L2 length
	CDR-H3 length		CDR-L3 length
	VH length		VL length

CDR-H1

CDR-H2

CDR-H3

VH: EVQLLESGGGLVQPGGSLRLSCAASGFTFSSYAMSWVRQAPGKGLEWWSAISWNSGSIYADSVKGRFTISRDNKNTLYLQMNSLRAEDTAVYYCAKDLLGSFPYDASGYDYFDYWGQGTLVTVSS

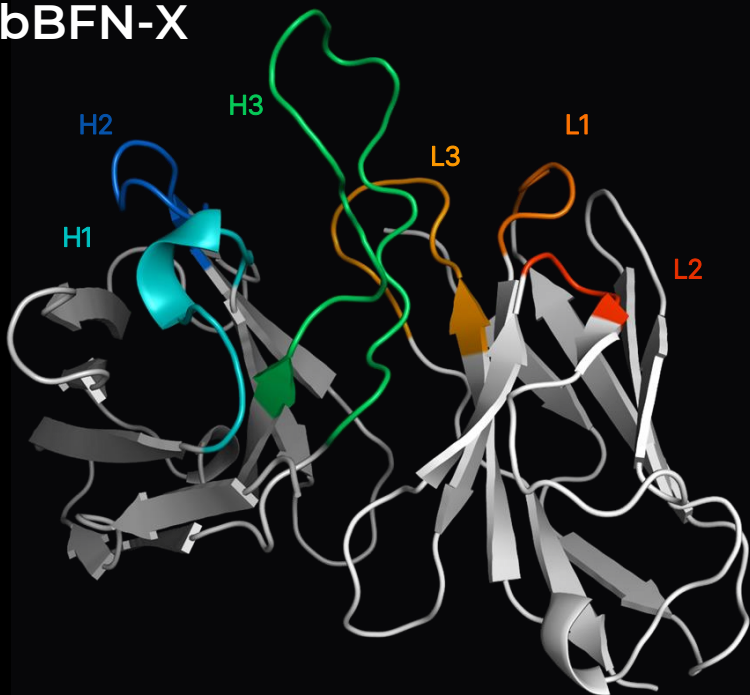
CDR-L1

CDR-L2

CDR-L3

VL: DIQMTQSPSSVSASVGDRTITCRASQSVSSNLAWYQQKPGKAPKLLIYGASSLQSGVPSRFSGSGSDFTLTISSLQPEDFATYYCQQANSFPPTFGQGTRLEIK

AbBFN-X



Amino acid sequence

AGL...	FWR-H1	AGL...	FWR-L1
AGL...	CDR-H1	AGL...	CDR-L1
AGL...	FWR-H2	AGL...	FWR-L2
AGL...	CDR-H2	AGL...	CDR-L2
AGL...	FWR-H3	AGL...	FWR-L3
AGL...	CDR-H3	AGL...	CDR-L3
AGL...	FWR-H4	AGL...	FWR-L4

Genetic Attributes

	HV gene		HV seq. identity
	HD gene		HD seq. identity
	HJ gene		HJ seq. identity
	LV gene		LV seq. identity
	LD gene		LJ seq. identity
	LC locus		Species

Biophysical Attributes

	Hydrophobicity
	Positive Patches
	Negative Patches
	Charge Imbalance

Length Attributes

	CDR-H1 length		CDR-L1 length
	CDR-H2 length		CDR-L2 length
	CDR-H3 length		CDR-L3 length
	VH length		VL length

CDR-H1

CDR-H2

CDR-H3

VH: EVQLLESGGGLVQPGGSLRLSCAASGFTFSSYAMSWVRQAPGKGLEWWSAISWNSGSIYADSVKGRFTISRDNKNTLYLQMNSLRAEDTAVYYCAKDLLGSFPYDASGYDYFDYWGQGTLVTVSS

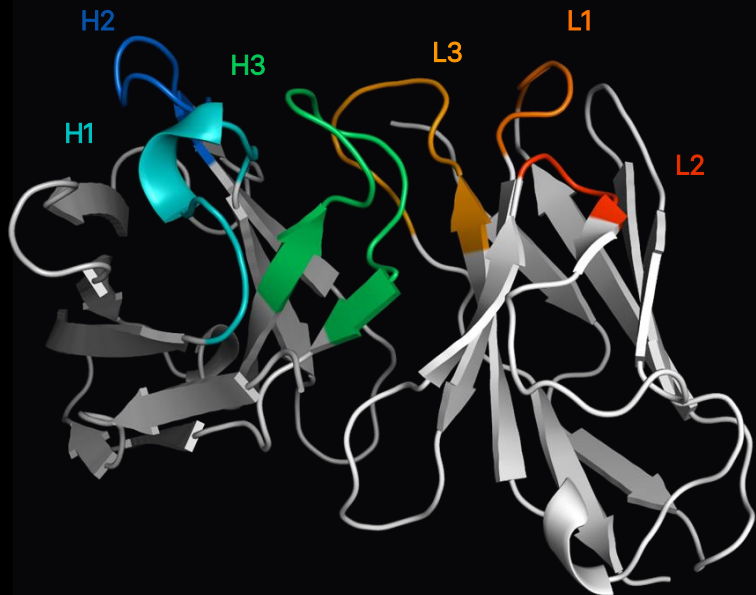
CDR-L1

CDR-L2

CDR-L3

VL: DIQMTQSPSSVSASVGDRTITCRASQSVSSNLAWYQQKPGKAPKLLIYGASSLQSGVPSRFSGSGSDFTLTISSLQPEDFATYYCQQANSFPPTFGQGTRLEIK

AbBFN-X



Amino acid sequence

AGL...	FWR-H1	AGL...	FWR-L1
AGL...	CDR-H1	AGL...	CDR-L1
AGL...	FWR-H2	AGL...	FWR-L2
AGL...	CDR-H2	AGL...	CDR-L2
AGL...	FWR-H3	AGL...	FWR-L3
AGL...	CDR-H3	AGL...	CDR-L3
AGL...	FWR-H4	AGL...	FWR-L4

Genetic Attributes

	HV gene		HV seq. identity
	HD gene		HD seq. identity
	HJ gene		HJ seq. identity
	LV gene		LV seq. identity
	LD gene		LJ seq. identity
	LC locus		Species

Biophysical Attributes

	Hydrophobicity
	Positive Patches
	Negative Patches
	Charge Imbalance

Length Attributes

	CDR-H1 length		CDR-L1 length
	CDR-H2 length		CDR-L2 length
	CDR-H3 length		CDR-L3 length
	VH length		VL length

CDR-H1

VH: EVQLLESGGGLVQPGGSLRLSCAAS

GFTFSSYA

MSWVRQAPGKGLEWWSA

CDR-H2

ISWNSGSI

YADSVKGRFTISRDN SKNTLYLQMNSLRAEDTAVYYC

CDR-H3

AKDRGGNWALIDY

WGQGTLVTVSS

CDR-L1

VL: DIQMTQSPSSVSASVGDRTITCRAS

QSVSSN

LAWYQQKPGKAPKLLIY

CDR-L2

GAS

SLQSGVPSRFSGSGSGTDFTLTISSLPEDFATYYC

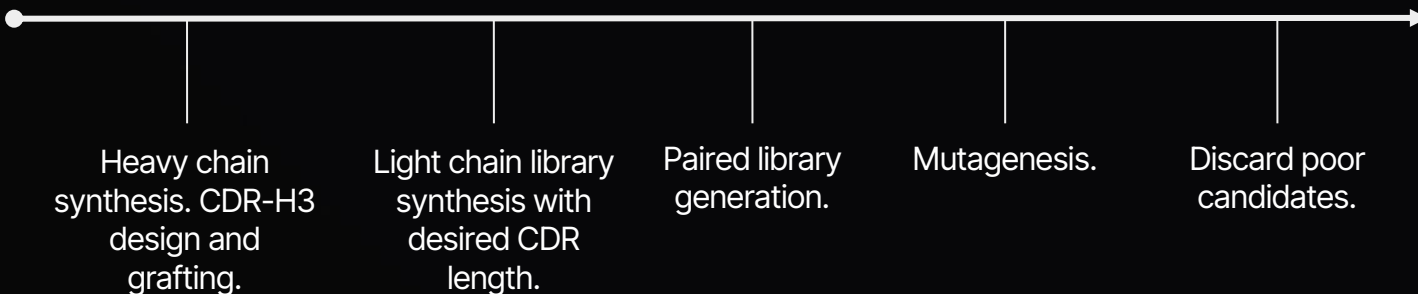
CDR-L3

QQANSFPPT

FGQGTRLEIK

Example Task 1: Anti-HIV Antibody Library Design

Generating a library of rare antibodies against HIV:



Example Task 1: Anti-HIV Antibody Library Design







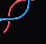





Generating a library of rare antibodies against HIV:

- (1) identify target attributes.
- (2) conditionally **sample for rare, desired antibodies.**


Amino acid sequence

AGL...	FWR-H1	AGL...	FWR-L1
AGL...	CDR-H1	AGL...	CDR-L1
AGL...	FWR-H2	AGL...	FWR-L2
AGL...	CDR-H2	AGL...	CDR-L2
AGL...	FWR-H3	AGL...	FWR-L3
AGL...	CDR-H3	AGL...	CDR-L3
AGL...	FWR-H4	AGL...	FWR-L4

Genetic Attributes

 HV gene	 HV seq. identity
 HD gene	 HD seq. identity
 HJ gene	 HJ seq. identity
 LV gene	 LV seq. identity
 LD gene	 LJ seq. identity
 LC locus	 Species

Length Attributes

 CDR-H1 length	 CDR-L1 length
 CDR-H2 length	 CDR-L2 length
 CDR-H3 length	 CDR-L3 length
 VH length	 VL length

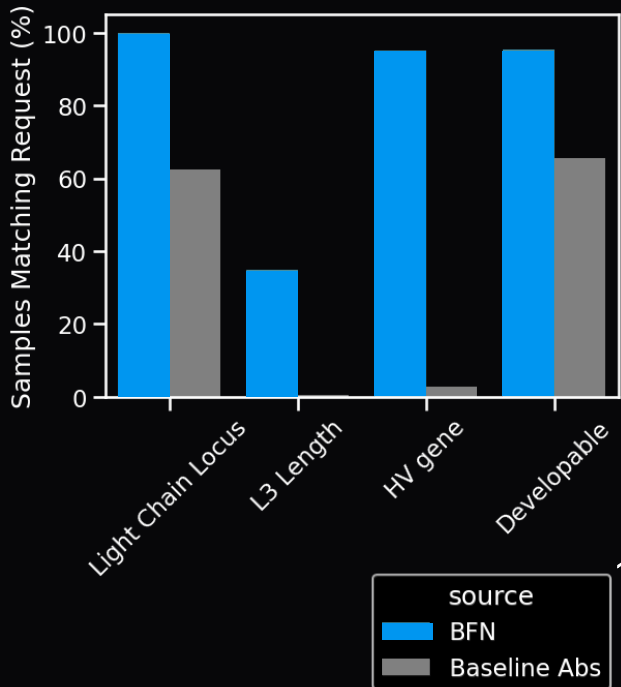
Biophysical Attributes

 Hydrophobicity	 Negative Patches
 Positive Patches	 Charge Imbalance

Example Task 1: Anti-HIV Antibody Library Design

Generating a library of rare antibodies against HIV:

- (1) identify target attributes.
- (2) conditionally **sample for rare, desired antibodies.**



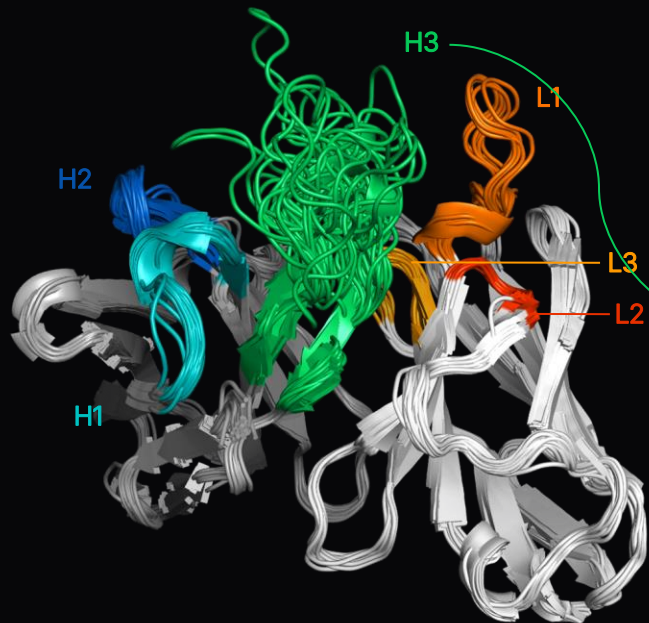
AbBFN-X antibodies are **5600x more likely²** to have **all desired characteristics.**

1. "Baseline Abs" refers to data sets of natural antibodies (Olsen *et al*, 2021, Prot. Sci.), "BFN" refers to samples generated by AbBFN-X.
2. compared to the rate of finding antibodies with the correct characteristics in data sets of natural antibodies (Olsen *et al*, 2021, Prot. Sci.)

Example Task 1: Anti-HIV Antibody Library Design

Generating a library of rare antibodies against HIV:

- (1) identify target attributes.
- (2) conditionally **sample for rare, desired antibodies.**



```

ARDEIYFLEWLISY
AKVRLGELPYEAFDI
ARGVRVQ SYNWFDP
ASGEYFFDTSSYPN
ARSSFVYPKSGYDFYFDY
ARDIAVDPESTAYFDY
AKGFSYGDGWADY
VRLRVGVLPGAFDI
ARDGGHYSH
ASGSGDSRYAQPLWFTTAFDI
ATSLNYGVIISD
ASGKMAVAYYFDY
AREGMDASMYYFDY
ARDMGYHDGALVFDN
:
    
```

Unique CDRs¹

100%

Unique CDR-H3s

100%

Unique CDR-L3s

99%

Unique non H/L3 CDRs

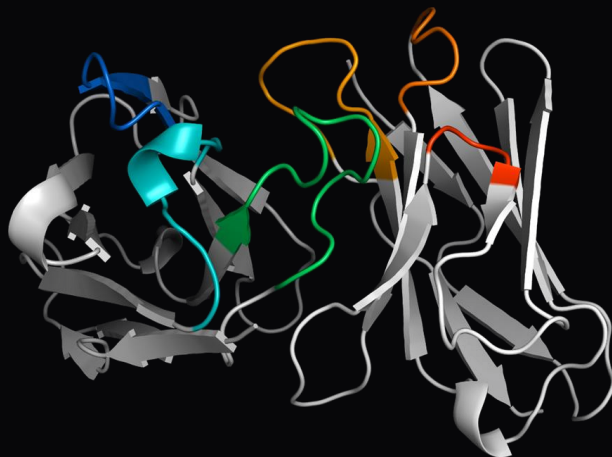
52%

1. 128 samples were generated, uniqueness assessed by considering all relevant regions at once, excluding framework regions.

Example Task 2: Heavy-Light Pairing

Generating a set of **developable** light chains that will **pair** with a heavy chain:

- (1) **Condition** on desired properties and heavy sequence.
- (2) **Sample** for **stable, diverse** sequences.



Example Task 2: Heavy-Light Pairing

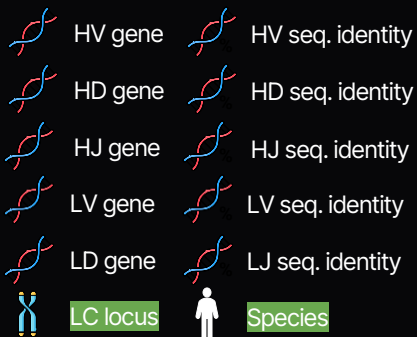
Generating a set of **developable** light chains that will **pair** with a heavy chain:

- (1) **Condition** on desired properties and heavy sequence.
- (2) **Sample** for **stable, diverse** sequences.

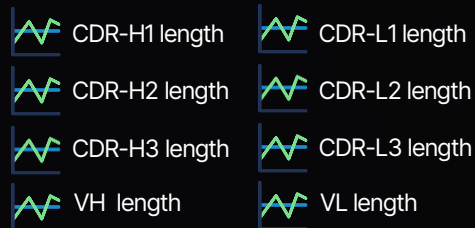
Amino acid sequence



Genetic Attributes



Length Attributes



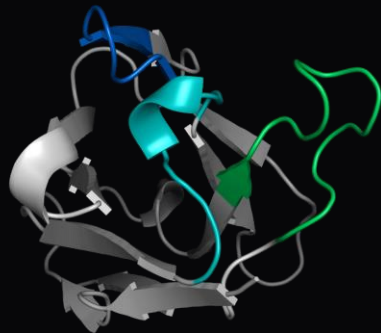
Biophysical Attributes



Example Task 2: Heavy-Light Pairing

Generating a set of **developable** light chains that will **pair** with a heavy chain:

- (1) **Condition** on desired properties and heavy sequence.
- (2) **Sample** for **stable, diverse** sequences.



Conditioning on the required heavy chain

Example Task 2: Heavy-Light Pairing

Generating a set of **developable** light chains that will **pair** with a heavy chain:

- (1) **Condition** on desired properties and heavy sequence.
- (2) **Sample** for **stable, diverse** sequences.

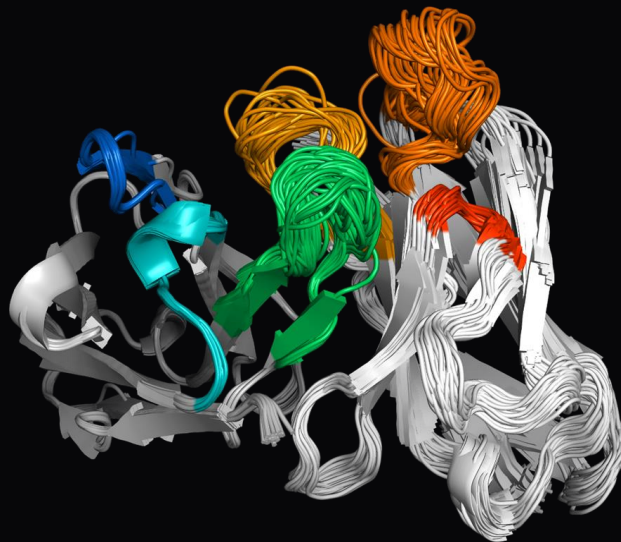


Conditioning on the required heavy chain
results in **diverse light chains** in length and
structure

Example Task 2: Heavy-Light Pairing

Generating a set of **developable** light chains that will **pair** with a heavy chain:

- (1) **Condition** on desired properties and heavy sequence.
- (2) **Sample** for **stable, diverse** sequences.



Conditioning on the required heavy chain

results in **diverse light chains** in length and structure

while respecting the **requested sequence** and generating **stable pairs**.

Generative AI for Proteomics

InstaDeep is developing **next-generation GenAI models** across the whole stack, from **fundamental ML research**, to modelling scientific data and **enabling new capabilities** for scientists.

Bayesian Flow Networks

Unified modelling of multi-modal data



All modalities as first class citizens



Joint learning of **heterogeneous data**



Task-specific **conditional generation**

Protein Sequence Modelling

Published demonstrations



Leading performance across tasks



Diverse & novel **de novo generation**



Zero-shot inpainting at inference

BFN-X

In-development foundation models



Sequence, genetic and **biophysical**



Learns **rational antibody principles**

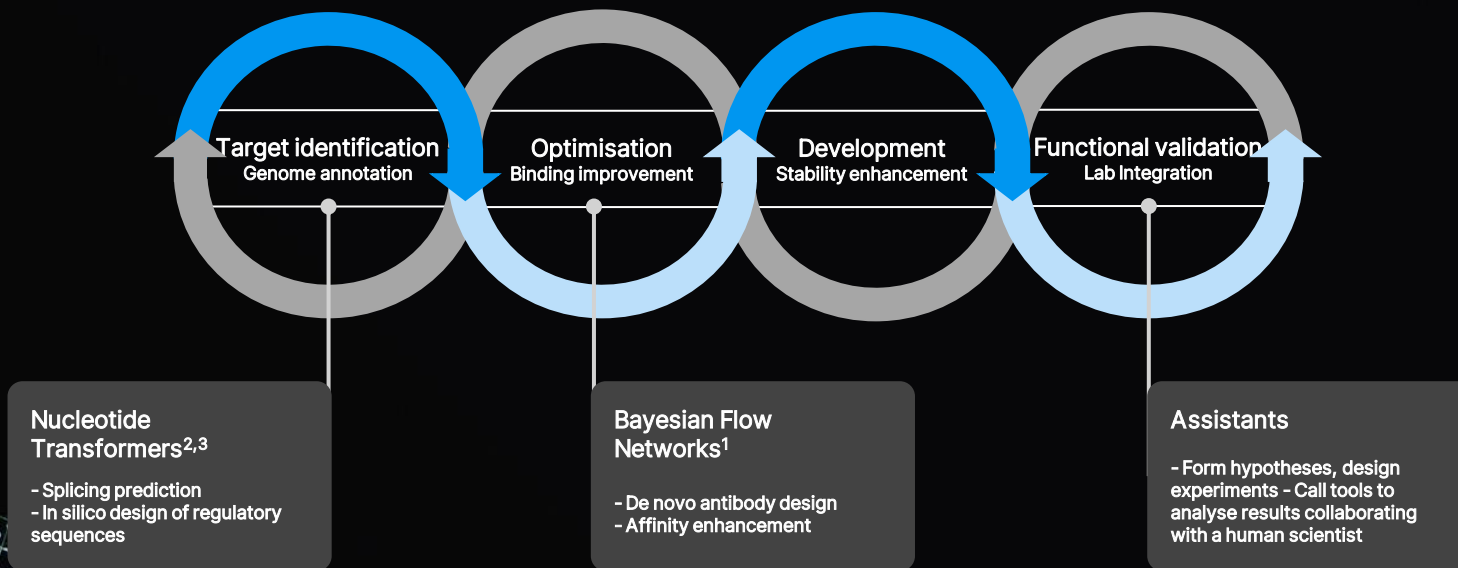


Enables **diverse suite of tasks**

DeepChain

One platform, Multiple tools

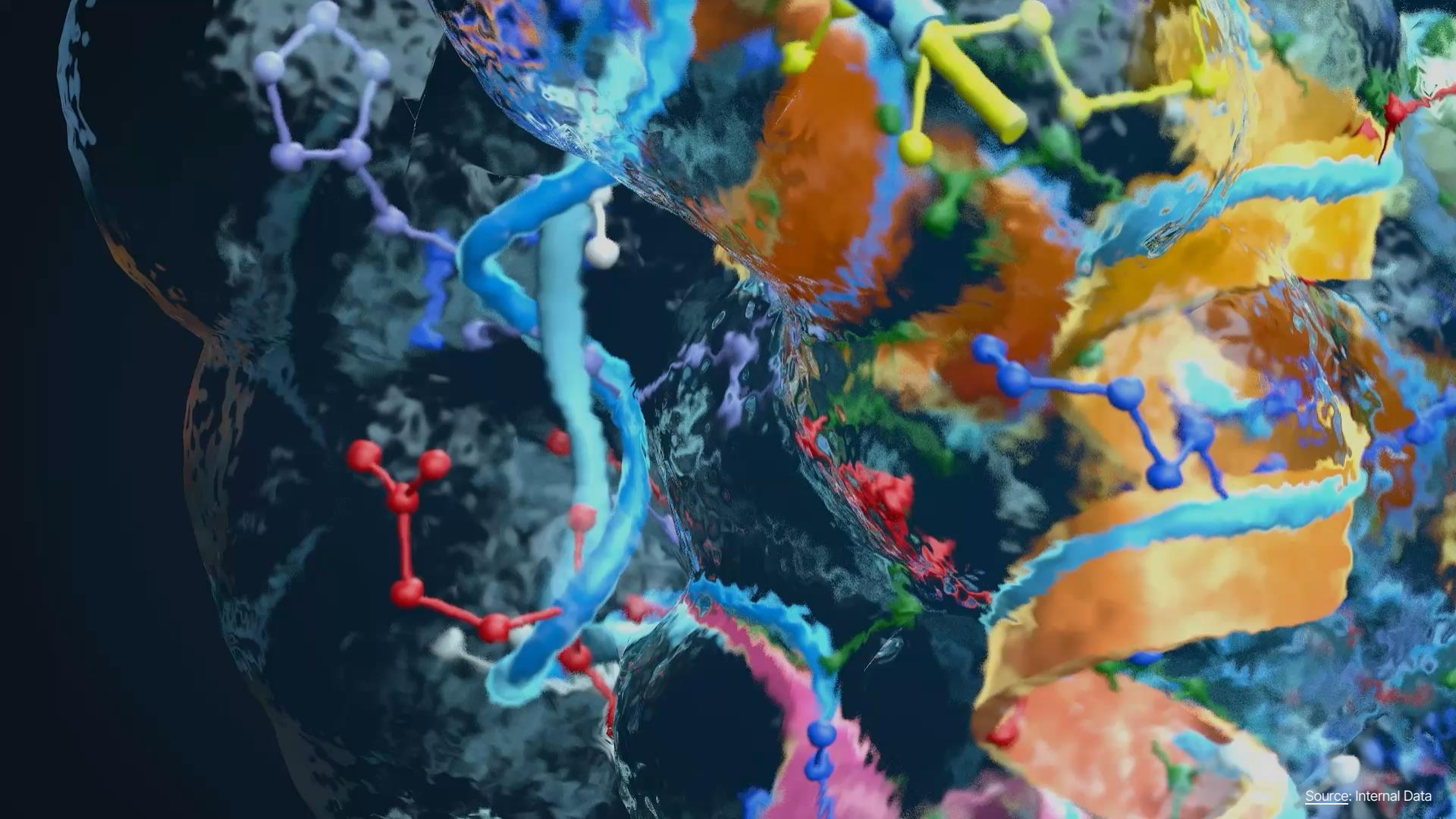
By combining State-of-the-art Science with Engineering, Our AI Tools aim to Accelerate the R&D pipeline



[1] Alice Sends Amino Acids to Bob: Protein Sequence Modelling with Bayesian Flow Networks, Barrett et al., Under Review (2024).

[2] The Nucleotide Transformer: Building and Evaluating Robust Foundation Models for Human Genomics, Dalla-Torre et al., Under Review (2023).

[3] SegmentNT: annotating the genome at single-nucleotide resolution with DNA foundation models, Almeida et al., Under Review (2024).



We Are Releasing our Flagship Models on DeepChain

ProtBFN & AbBFN¹

State-of-the-art generative protein models

- ✓ Generates natural-like, diverse, structurally coherent, and novel protein sequences
- ✓ Outperforms leading autoregressive and discrete diffusion models
- ✓ Enables flexible conditional generation in a zero-shot manner

Nucleotide Transformer² & SegmentNT³

Our foundation models for DNA

- ✓ Single nucleotide resolution
- ✓ Up to 50kb context length without performance drop
- ✓ Generalizes across species

[1] Alice Sends Amino Acids to Bob: Protein Sequence Modelling with Bayesian Flow Networks, Barrett et al., Under Review (2024).

[2] The Nucleotide Transformer: Building and Evaluating Robust Foundation Models for Human Genomics, Dalla-Torre et al., Under Review (2023).

[3] SegmentNT: annotating the genome at single-nucleotide resolution with DNA foundation models, Almeida et al., Under Review (2024).

We Are Releasing our Flagship Models on DeepChain

ProtBFN & AbBFN¹

State-of-the-art generative protein models

- ✓ Generates natural-like, diverse, structurally coherent, and novel protein sequences
- ✓ Outperforms leading autoregressive and discrete diffusion models
- ✓ Enables flexible conditional generation in a zero-shot manner

Nucleotide Transformer² & SegmentNT³

Our foundation models for DNA

- ✓ Single nucleotide resolution
- ✓ Up to 50kb context length without performance drop
- ✓ Generalizes across species

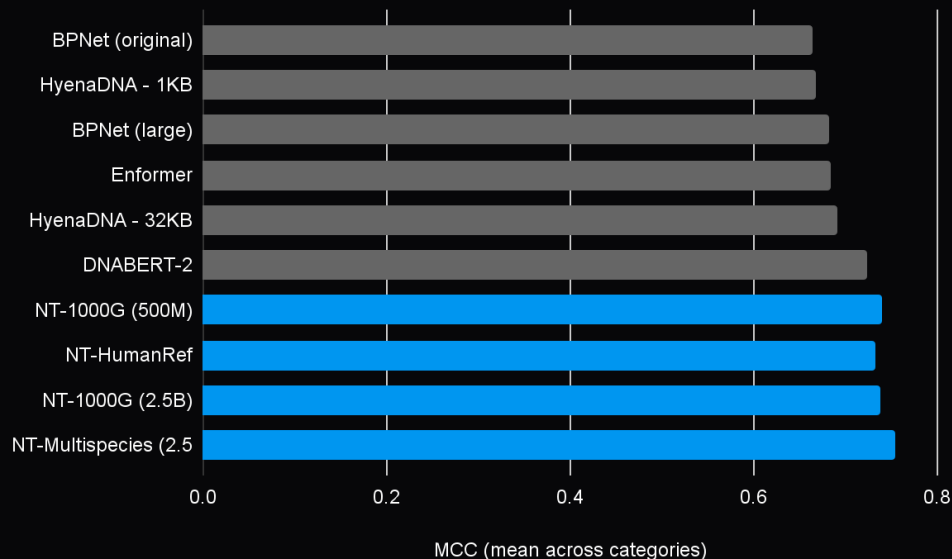
[1] Alice Sends Amino Acids to Bob: Protein Sequence Modelling with Bayesian Flow Networks, Barrett et al., Under Review (2024).

[2] The Nucleotide Transformer: Building and Evaluating Robust Foundation Models for Human Genomics, Dalla-Torre et al., Under Review (2023).

[3] SegmentNT: annotating the genome at single-nucleotide resolution with DNA foundation models, Almeida et al., Under Review (2024).

Our Foundation Models for Genomics are State-of-the-Art

Nucleotide transformer models are state-of-the-art in the space¹



[1] The Nucleotide Transformer: Building and Evaluating Robust Foundation Models for Human Genomics, Dalla-Torre et al., Nat. methods (press)

+700K
Downloads

Across model sizes¹

One of the Most
Downloaded Genomics AI
Models on Hugging Face²

[1] Cumulative Downloads for Nucleotide Transformer models, 50-2.5B parameter sizes, September 2024. Hugging Face Statistics. Models Release Date: April 2023.

[2] Count by family of models under the "Genomics" official Hugging Face tag: <https://huggingface.co/models?other=genomics&sort=downloads>, September 2024.

**We are releasing capabilities
to build and scale with our AI
models**

01

Optimized Setup

02

03



01

Optimized Setup

Get access to our hardware-accelerated workflows to run models with a few lines of code

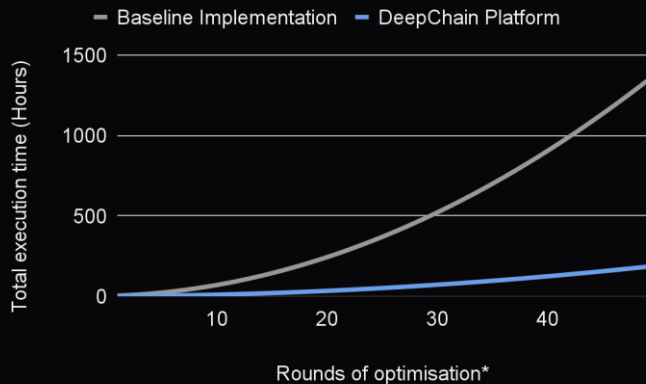
Send a request to the Inference API and receive a fast response containing the model's output

02

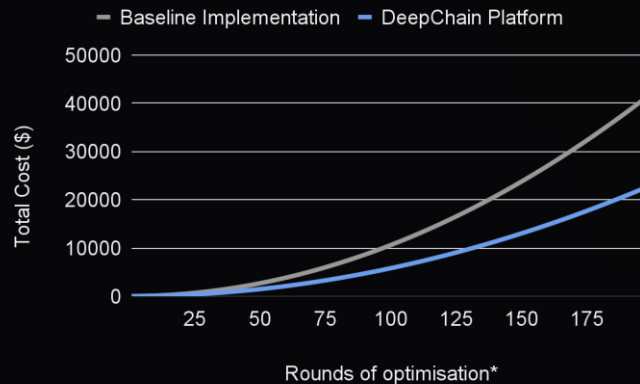
03

Running inference with DeepChain is 7X faster and 2X cheaper for *in silico* design of regulatory sequences*

Improved Speed



Reduced Cost



* reference methodology: Jores, T., Tonnie, J., Wrightsman, T. et al. Synthetic promoter designs enabled by a comprehensive analysis of plant core promoters. Nat. Plants 7, 842–855 (2021).

* Test implementation: Sequence length: 6kbp and 2.1kbp sequences, Parameters: --num_indels=8000, --prop_indels=0.5, --random_indels=True, --min_indels_size=2, --max_indel_size=5, --tissue_optimize_idx=1, --opt_metric=increase, --num_rounds=30

* Baseline implementation set-up: 1 NVIDIA V100 Tensor Core GPUs, using [published](#) Pytorch implementation available on Hugging Face

01

Optimized Setup

02

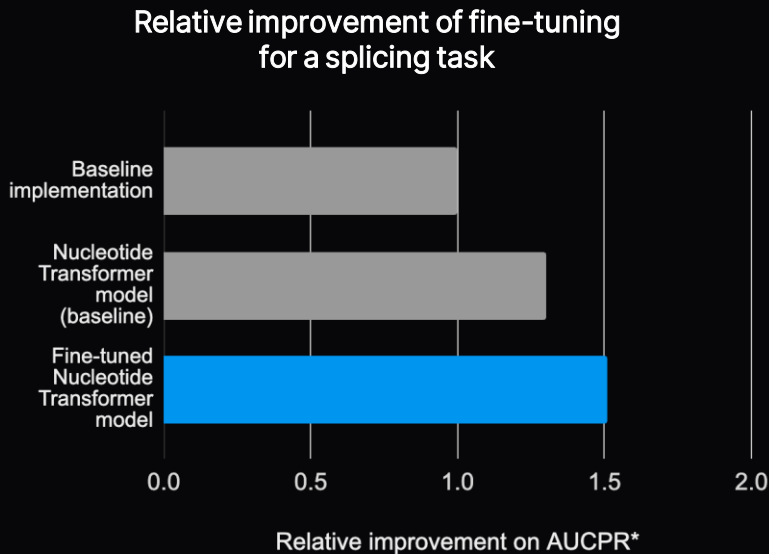
Customization

Customize models for your needs

Customize a model for a specific task via supervised fine-tuning with our proprietary parameter efficient fine-tuning methods

03

Fine-tuning a model on a specialised data set increased performance $\approx 1.5X$ for a splicing prediction use-case



* AUCPR: metric that measures the overall performance of a binary classification model by plotting precision against recall at different threshold settings, providing a more accurate assessment of performance for imbalanced classes.

* Dataset used for customer-specific calculation: Shiraishi Y, Kataoka K, Chiba K, Okada A, Kogure Y, Tanaka H, Ogawa S, Miyano S. A comprehensive characterization of cis-acting splicing-associated variants in human cancer. *Genome Res.* 2018 Aug;28(8):1111-1125. doi: 10.1101/gr.231951.117. Epub 2018 Jul 16. PMID: 30012835; PMCID: PMC6071634.

* Baseline implementation: Predicting Splicing from Primary Sequence with Deep Learning. Jaganathan, Kishore et al. *Cell*, Volume 176, Issue 3, 535 - 548.e24

01

Optimized Setup

02

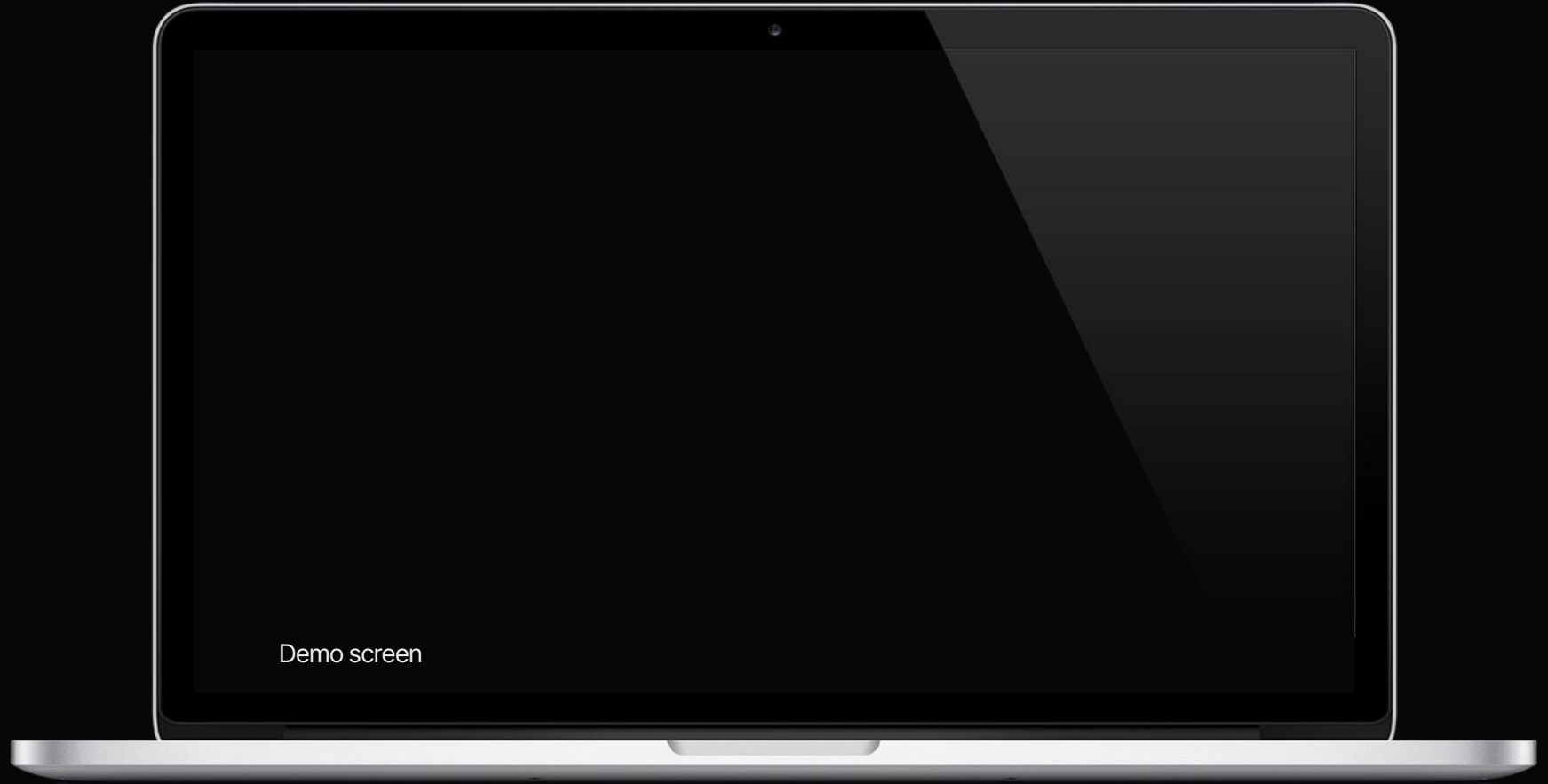
Customization

03

Assistants

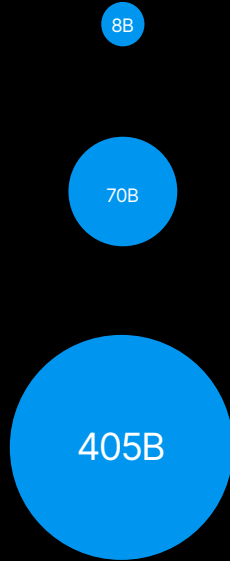
Build using natural language with
Laila

Laila can form hypotheses, design experiments,
and call tools to analyse results collaboratively with
a human scientist

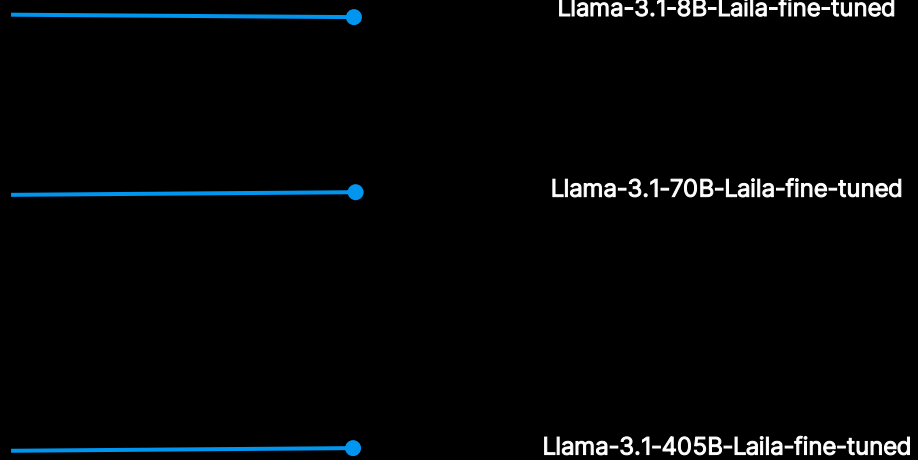


The Laila Series of AI Agents is built with Meta Llama 3.1

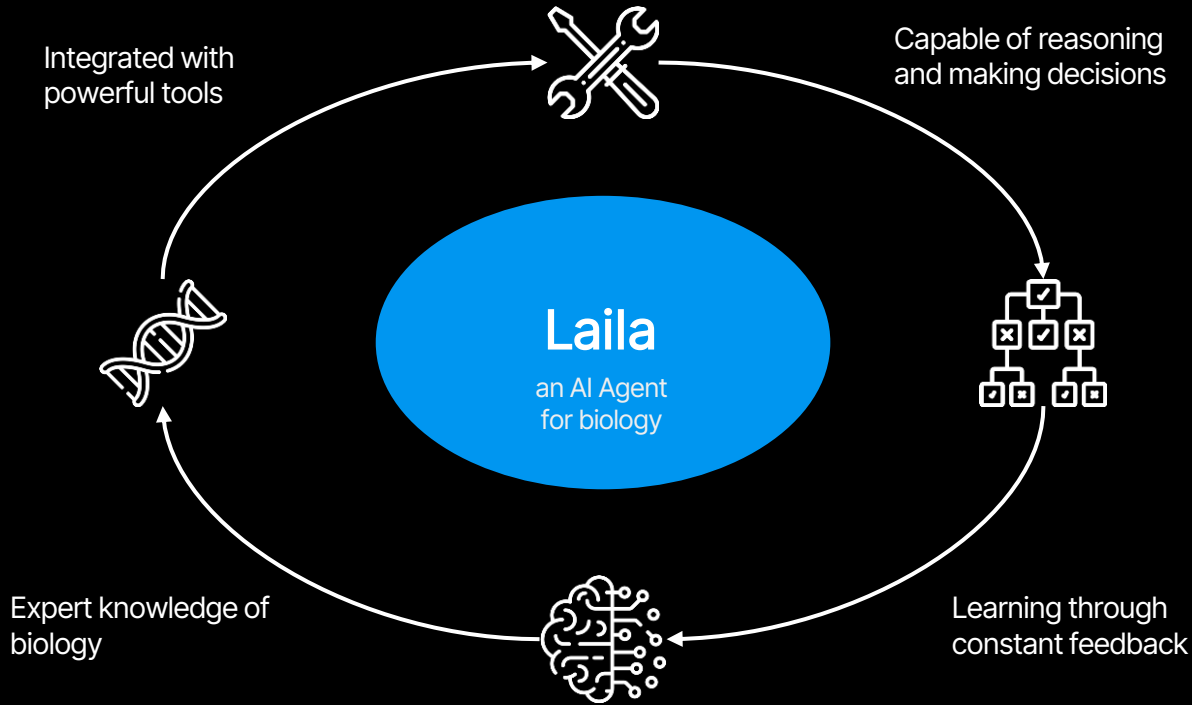
Model Sizes *in billions of parameters*



Model Versions *internally fine-tuned by InstaDeep*



Laila is more than just a chatbot...



Accelerate your R&D Pipeline with Our Flagship AI Models

Industry-leading AI solutions for the life sciences

[Get Started with DeepChain](#)

Go to deepchain.bio to gain early access

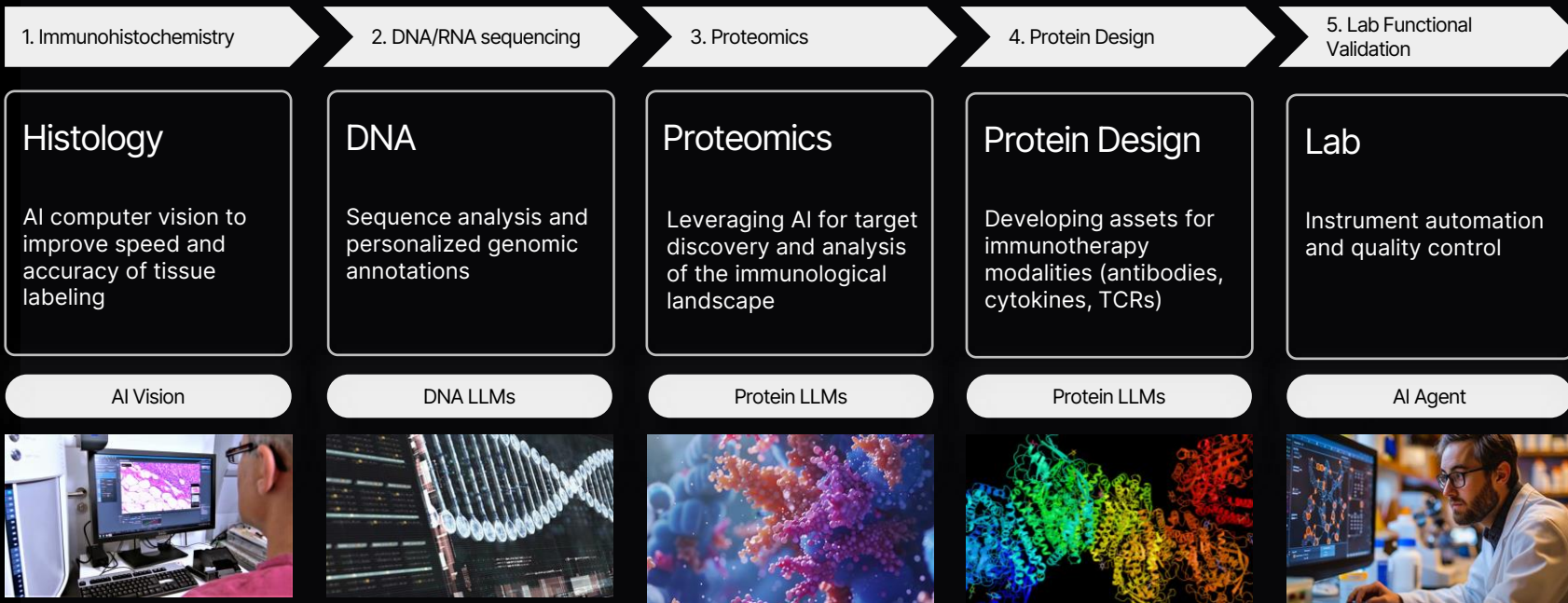
Part II.

Deploying AI Across The Pipeline

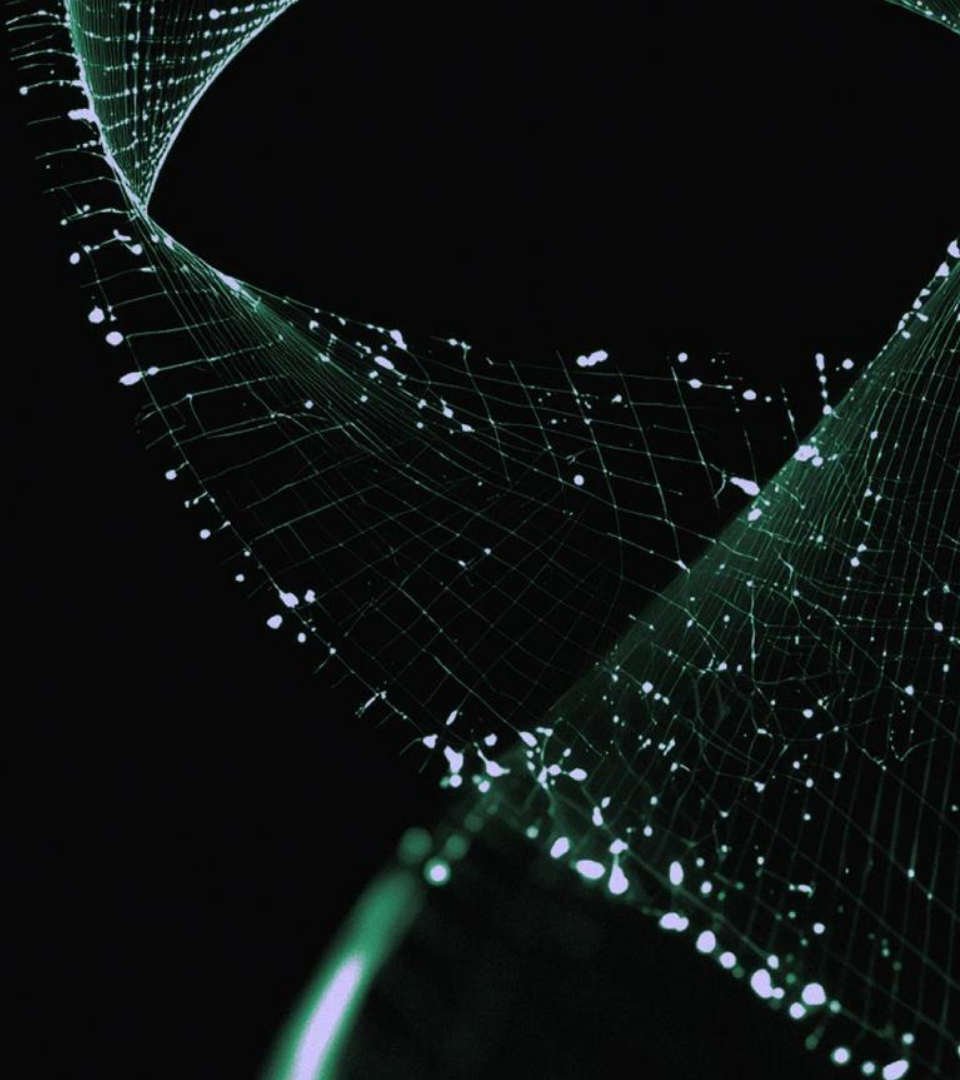


Our Goal: Deploying AI end-to-end in our immunotherapy pipeline

AI-first Immunotherapy Platform



Step 1: Histology



The Challenge

Pathologists face increasing workloads as demand grows for precise tumor and tissue analysis.



Source: Own Data

2 mm

16451.54, 6953.85 μ m
157, 79, 178

Our Approach

Develop tools using **state-of-the-art AI** to enhance pathologists' workflows by producing **faster** and **higher-quality** annotations.

AI-Assisted Tissue Annotation Tool

Aiming to enhance Pathologists' Precision And Efficiency Through Human-AI Collaboration

5X Speed-Up

By **increasing efficiency fivefold** compared to manual annotation, our AI tool allows pathologists to complete annotations in a fraction of the time, optimizing resource utilization and accelerating research and development efforts.

Superior Annotation Quality

Enabling pathologists to refine annotations at **different magnification levels** in Whole Slide Images, resulting in **higher-quality** annotations.

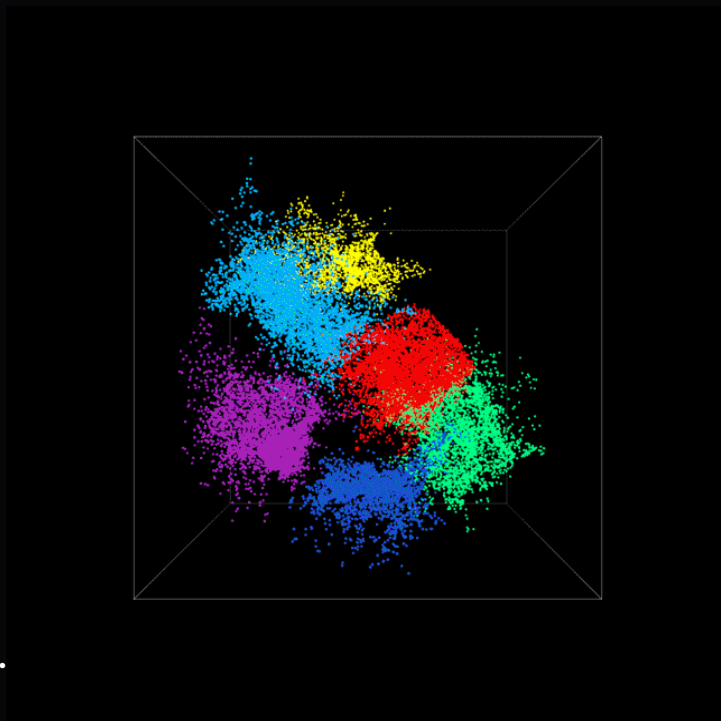
Whole Slide Image Segmentation Tool

From Classification To Segmentation

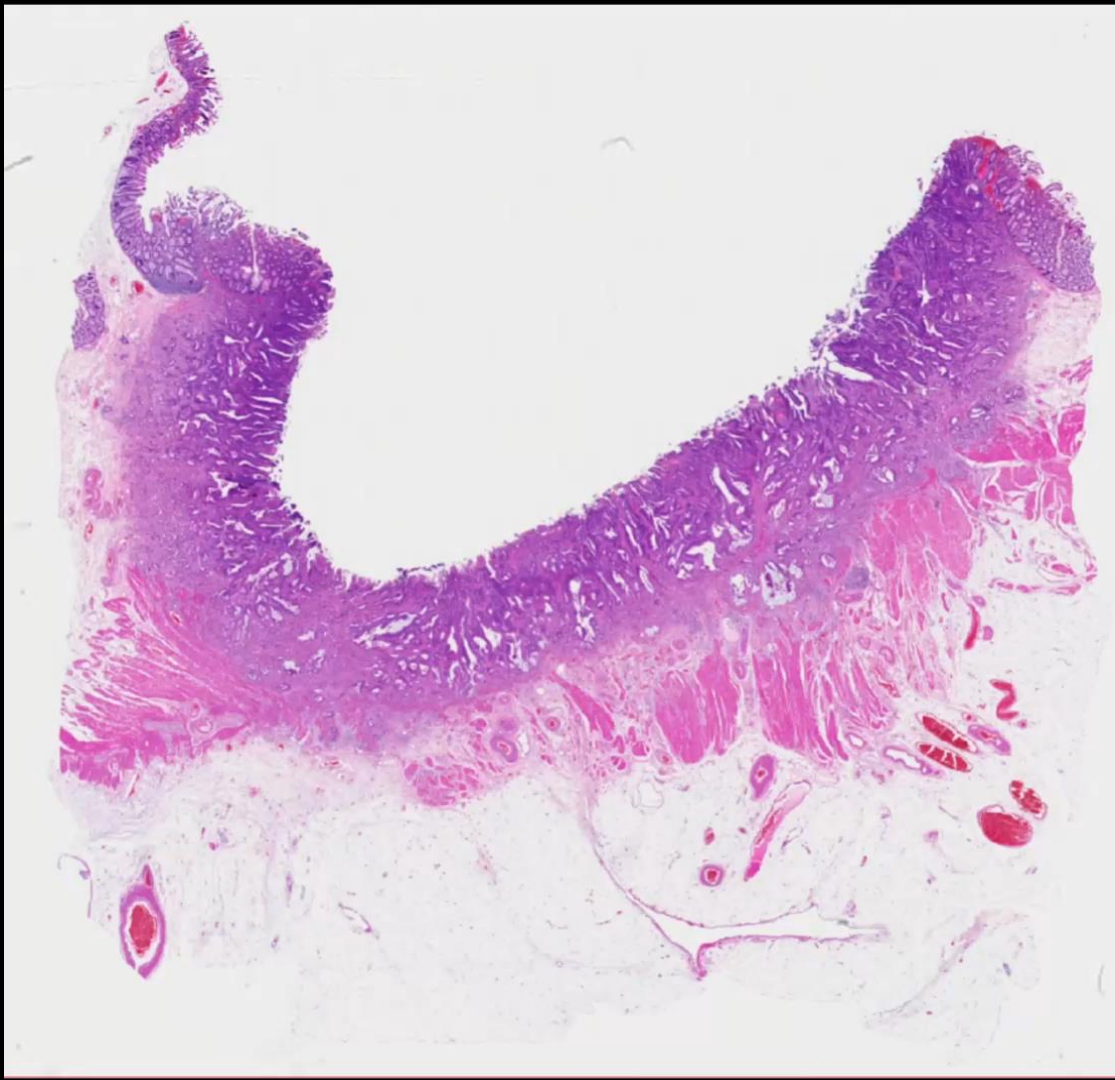
Use a **state-of-the-art vision foundation model** that we train specifically on pathology images.

Decompose the image into **patches**.

Transform the process from image **segmentation** to **patch classification**.



Source: Own Data



Scalable and Fast

By processing thousands of patches in parallel,
We deliver **more than a 100 times speed-up** compared to manual
annotation.

Step 2:

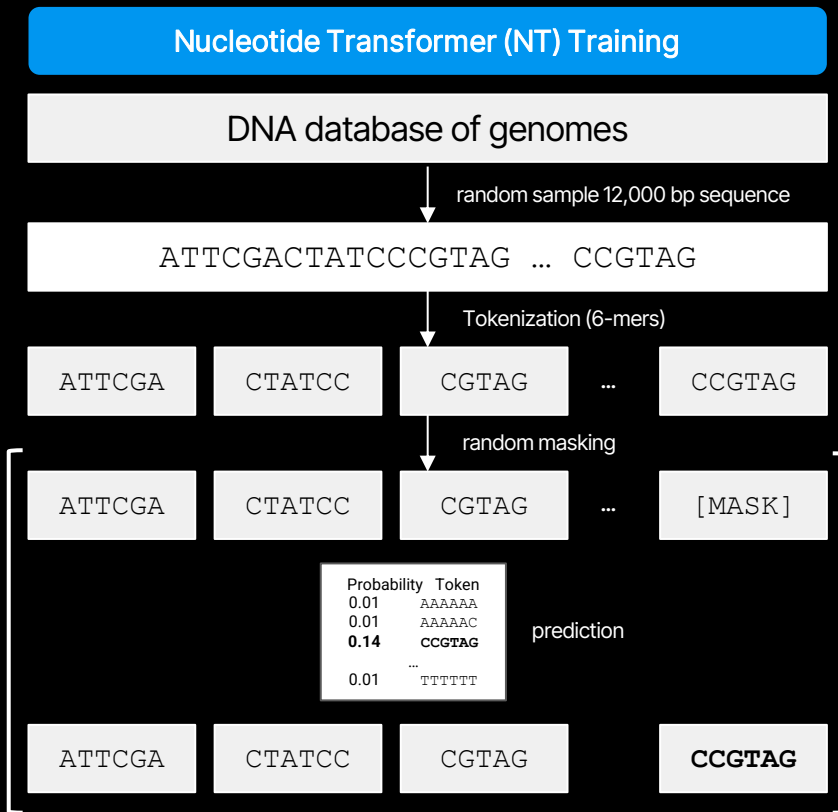
DNA Foundation Models at Nucleotide Resolution

Nucleotide Transformer: Self-Supervised Learning on Genomes

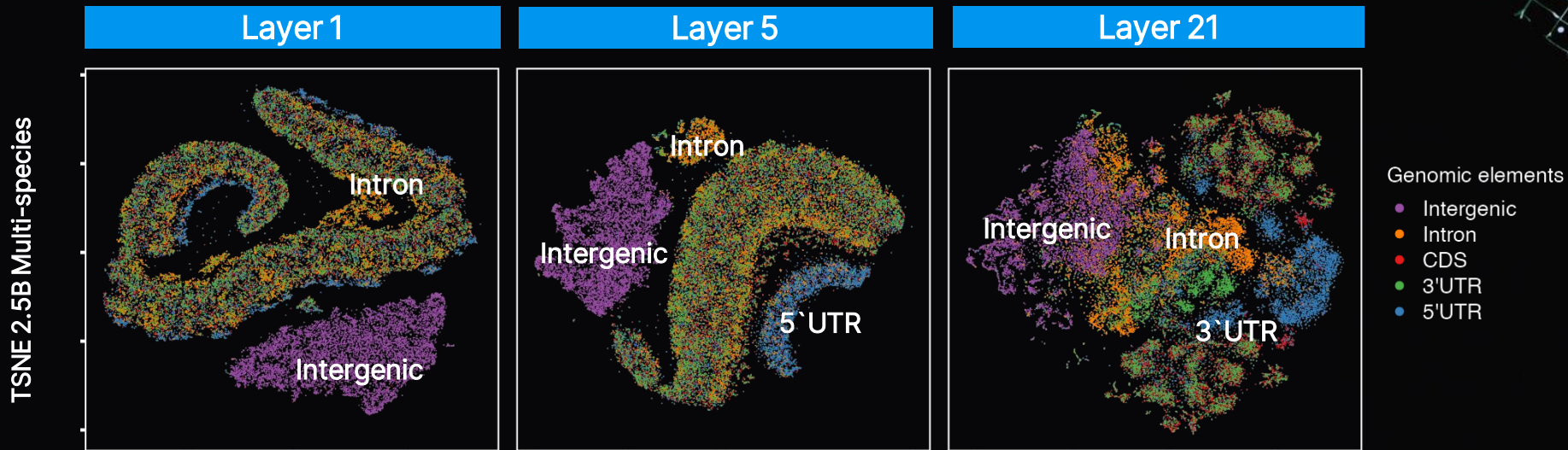
We believe automated analysis and predictions from genome sequences have the potential to transform tomorrow's health care and agriculture.

InstaDeep's Nucleotide Transformer Models

- *Architecture*: Masked Language Models (Bert-style training).
- *Datasets*: Trained on 5 datasets from different sizes with inter and intra-species variability from the whole tree of life.
- *Nucleotide Transformers (NT)*
 - V1: 500M, 1B, 2.5B parameters (2022)
 - V2: 50M, 100M, 250M parameters (2023)
- *Hardware*
 - Cambridge-1 Datacenter (collaboration with Nvidia)
 - TPUv4-1024 Pod (collaboration with Google Cloud)



NT Acquires Genomics Knowledge During Pre-Training

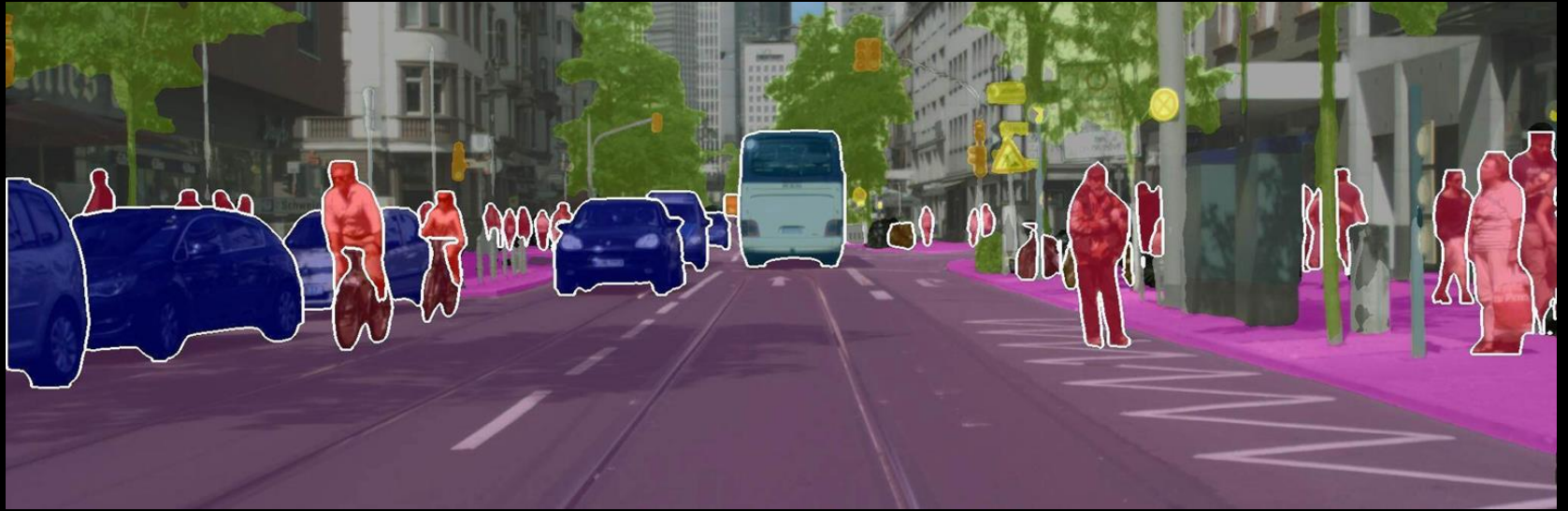


Even without supervision, information about genomic sequence features is learned in the "sequence representation"

Source: Dalla-Torre et al., Nat Methods 2024 (in press)

SegmentNT: Inspiration from Computer Vision Segmentation Models

Computer Vision

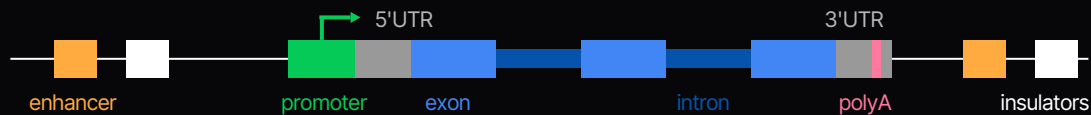


Genomics



SegmentNT: Annotating the Genome at Nucleotide Resolution

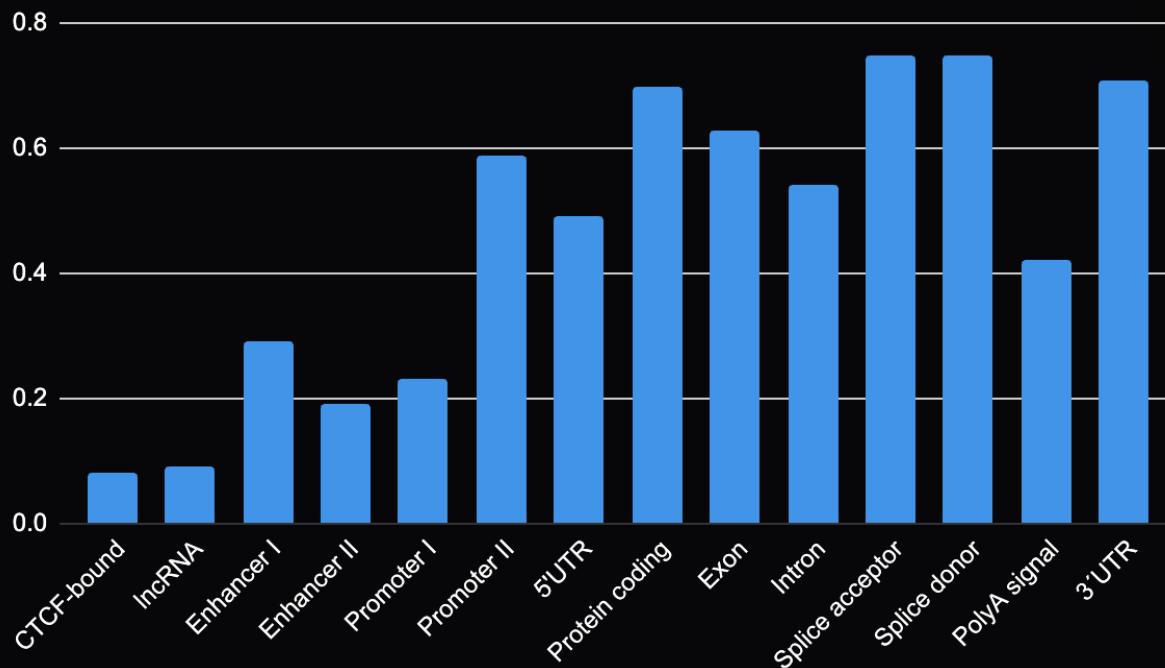
We fine-tuned the Nucleotide Transformers on 2.5M high-quality gene and regulatory elements annotations.



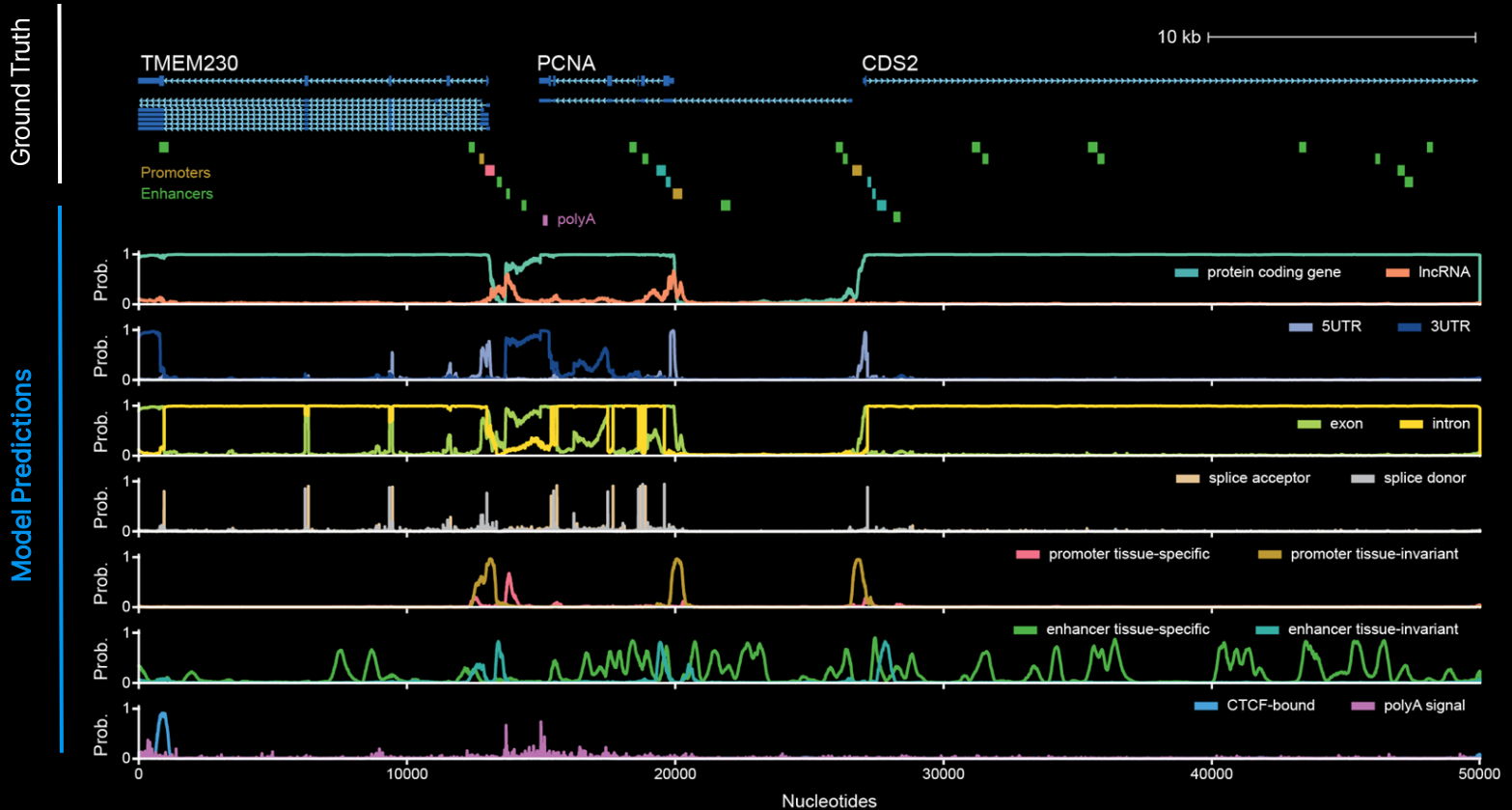
14 annotations per nucleotide
(e.g. 700,000 predictions at
50kbp)

SegmentNT: Annotating the Genome at Nucleotide Resolution

The resulting SegmentNT model segments sequences up to 50kbp with state-of-the-art performance for splicing, gene finding and regulatory element detection.

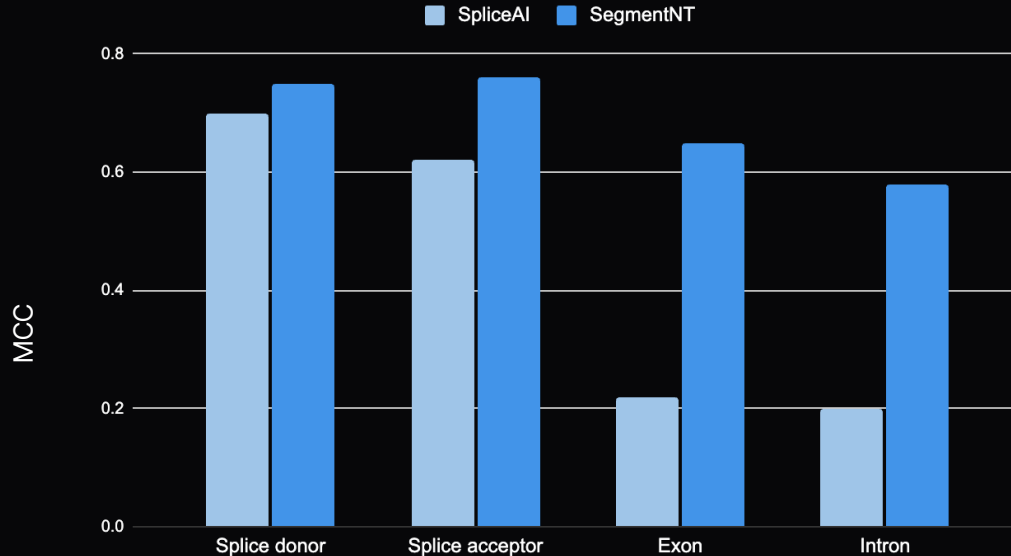


SegmentNT: 700,000 Accurate Predictions over 50kbp in less than a second



SegmentNT is State-of-the-Art for Canonical Splicing Detection

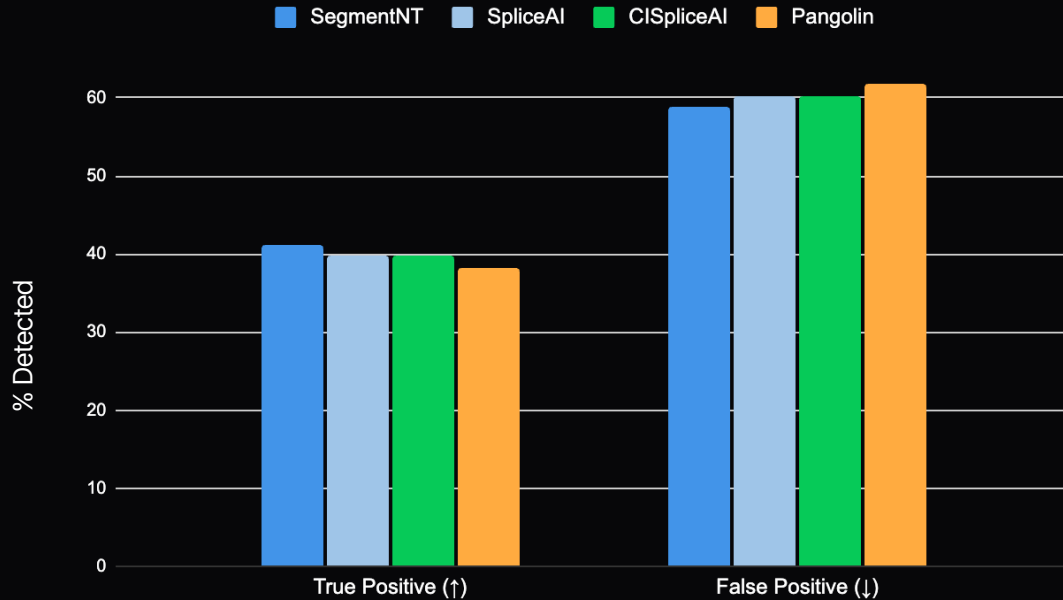
Splicing is a biological process that **removes non-coding sequences** (introns) from a primary messenger RNA (mRNA) transcript and joins the coding sequences (exons) together to create a mature mRNA. **Dysregulated splicing** can be a **vulnerability in cancers**. SegmentNT **outperforms state-of-the-art SpliceAI** for splicing event detection.



Test performance for splicing detection over full chromosomes on human reference genome

Alternative Splicing Events Detection with SegmentNT

Alternative splicing events can disrupt protein production and cancer pathways and is associated with cancer development. We finetuned segmentNT to identify tumor antigen candidates from alternative splicing events, which represent potential targets for personalized cancer immunotherapies. After finetuning, SegmentNT can accurately predict alternative splicing events in cancer data.

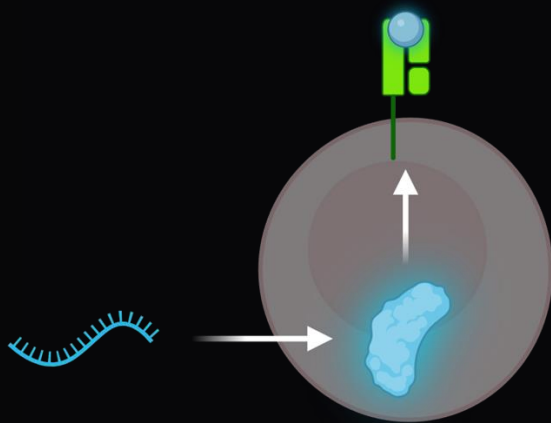


Test performance for 2,000 alternative splicing detection over cancer data (TCGA LUAD data)

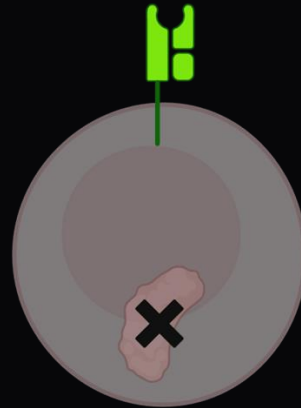
Step 3:

AI-Enhanced Proteomics for Target Discovery

Intracellular Proteins are Processed and Presented on MHC Complexes



T cell-targeting RNA vaccine



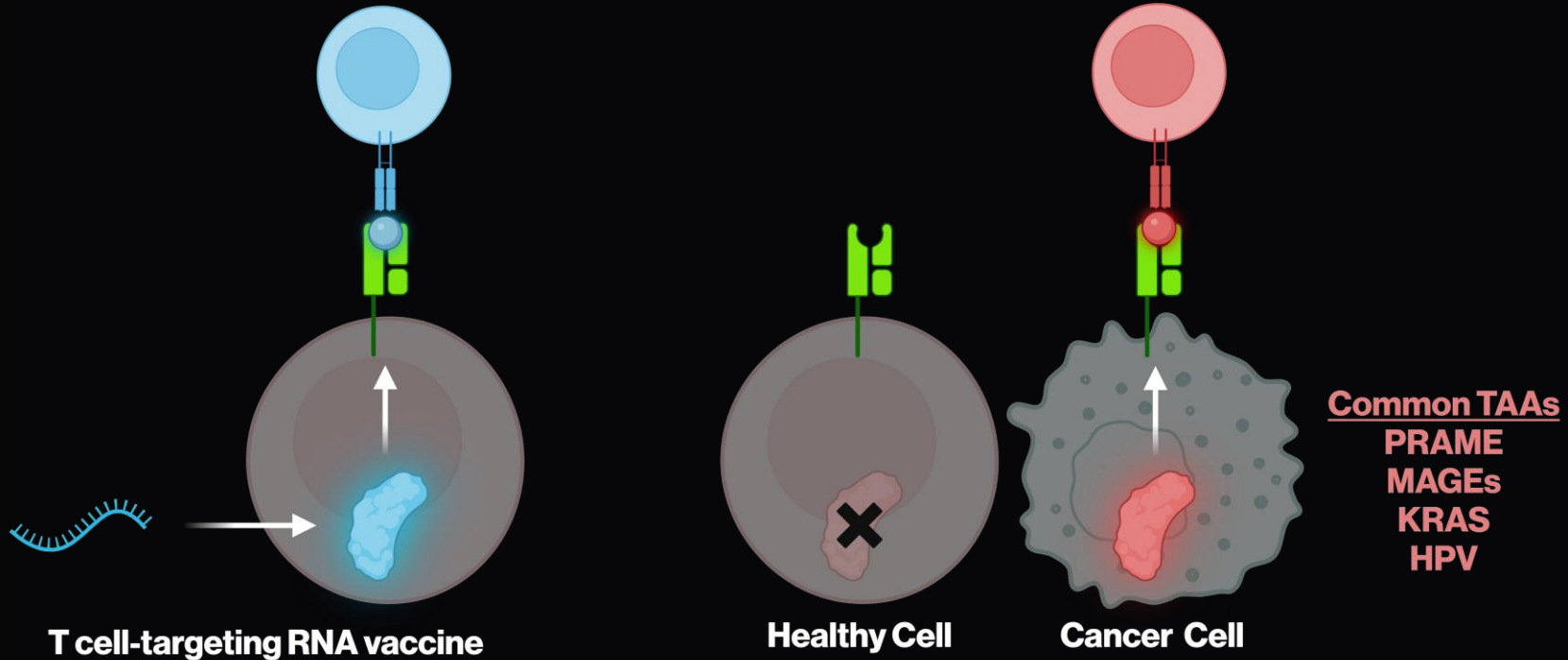
Healthy Cell



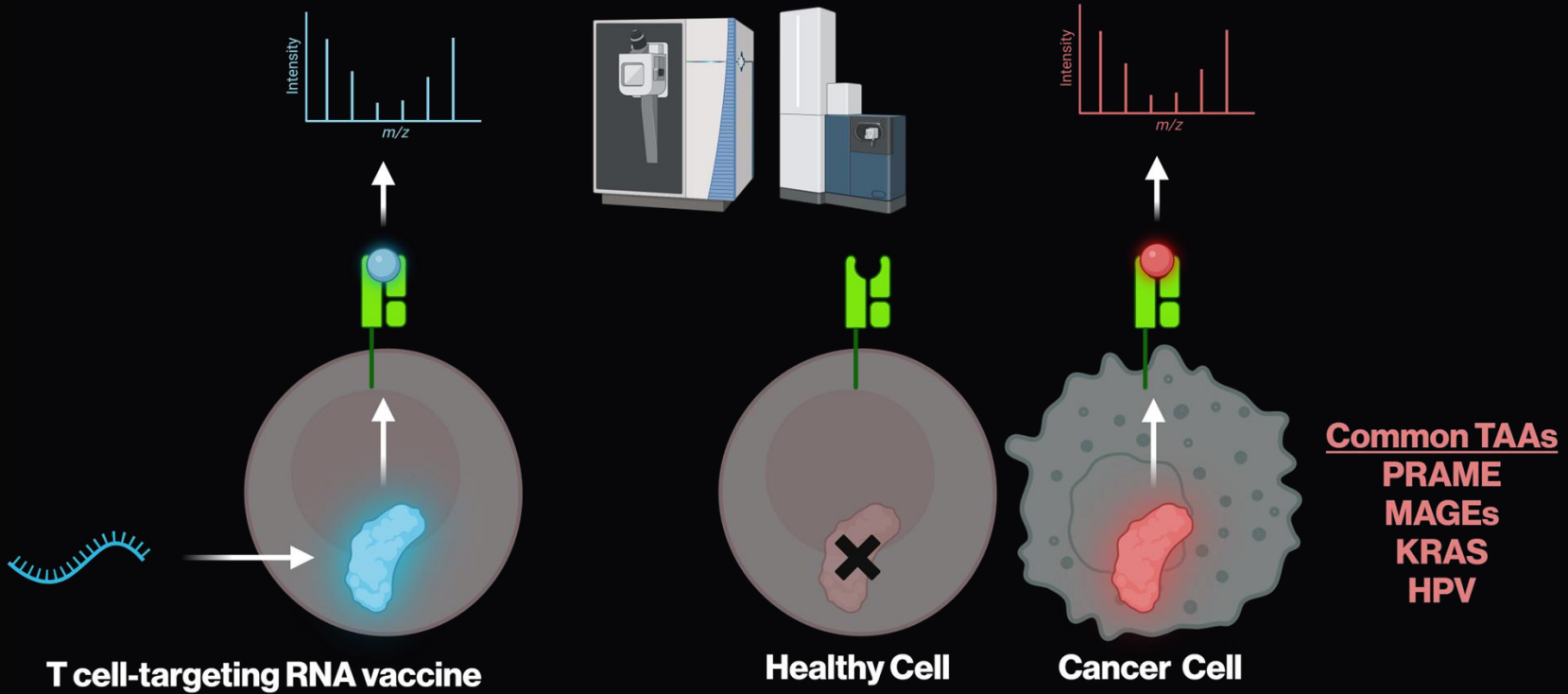
Cancer Cell

Common TAAs
PRAME
MAGEs
KRAS
HPV

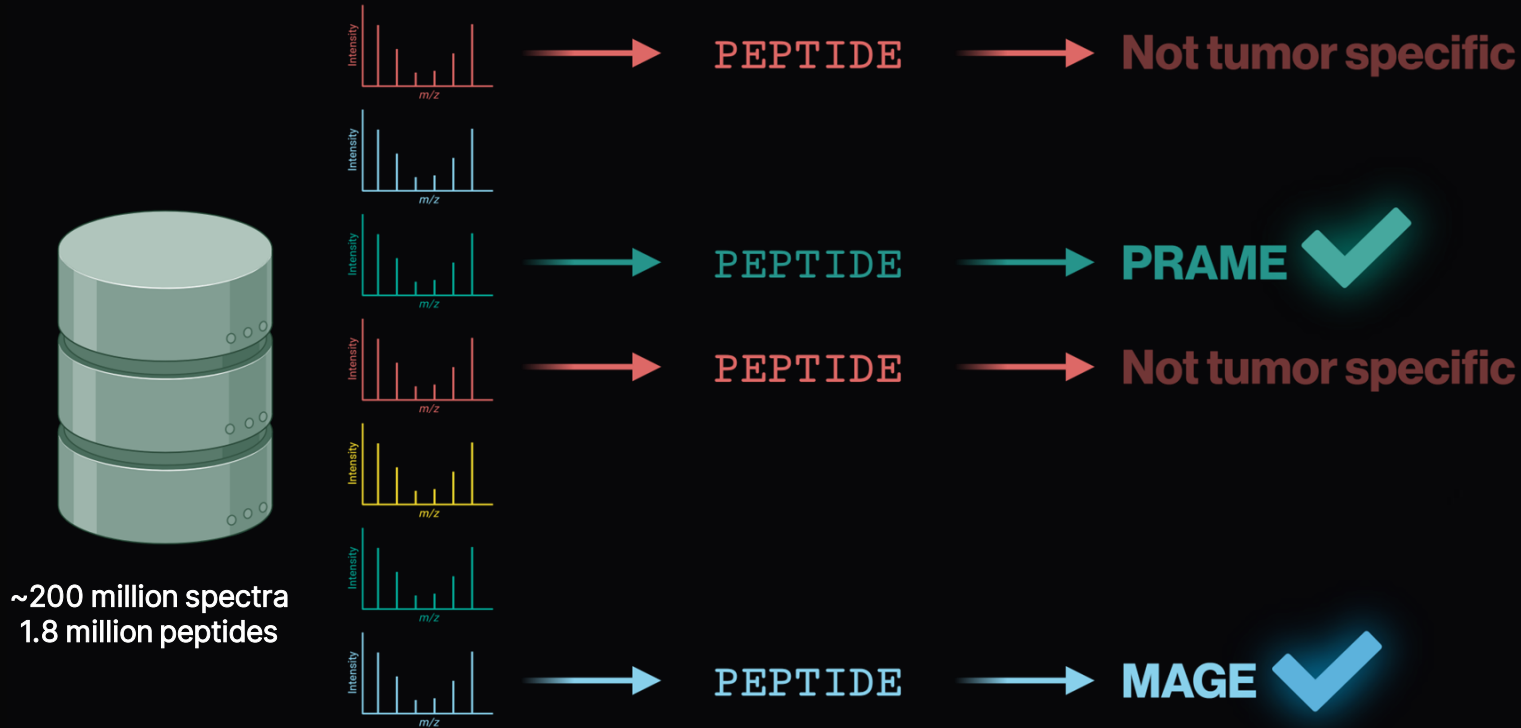
MHC-presented Epitopes are Immune System's Window into the Intracellular Proteome



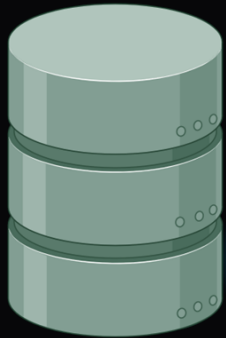
Mass Spectrometry is the Current State-of-the-art for Detecting, Identifying, and Quantifying MHC-presented Epitopes



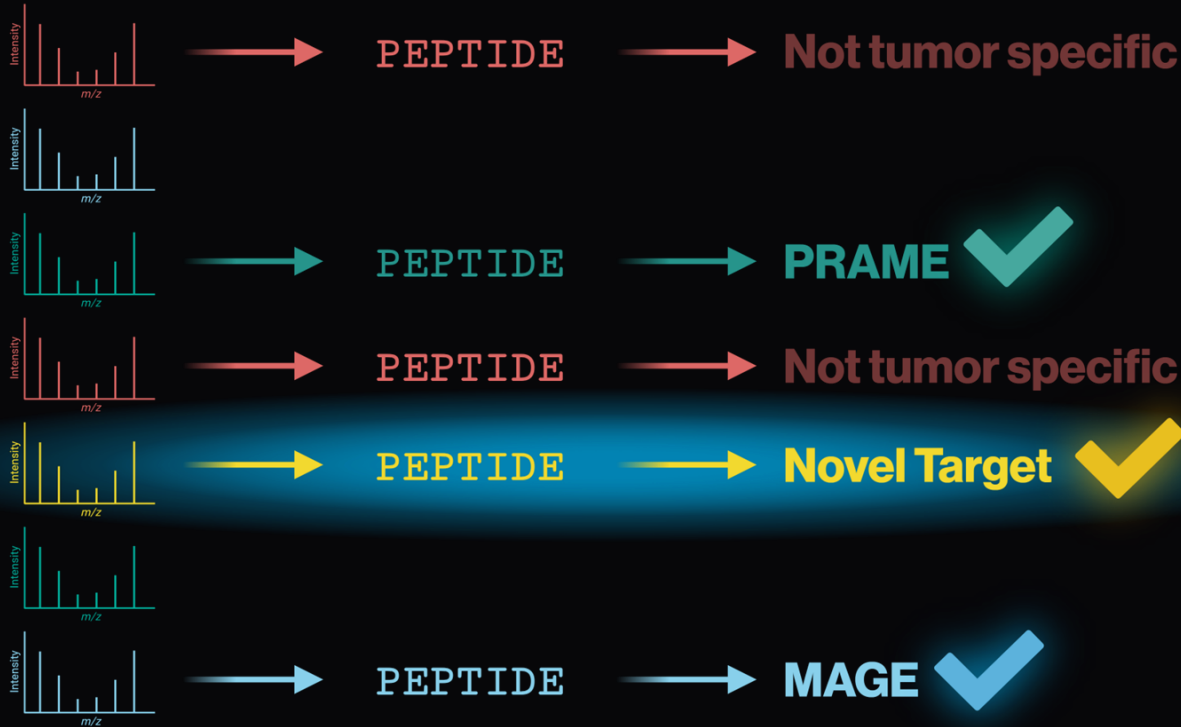
BioNTech has a Massive Database of MS Validated MHC-bound Epitope Peptides from Studies Performed Internally and Externally



With the Power of AI, we can Dig Deeper into our Data to Identify Novel Therapeutic Targets



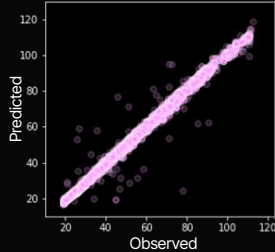
~200 million spectra
>>1.8 million peptides



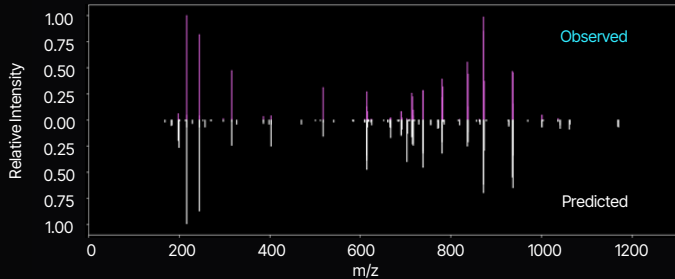
AI Maximizes our Ability to Discover Novel Cancer Targets

Up to 200% increase in recovered peptide IDs

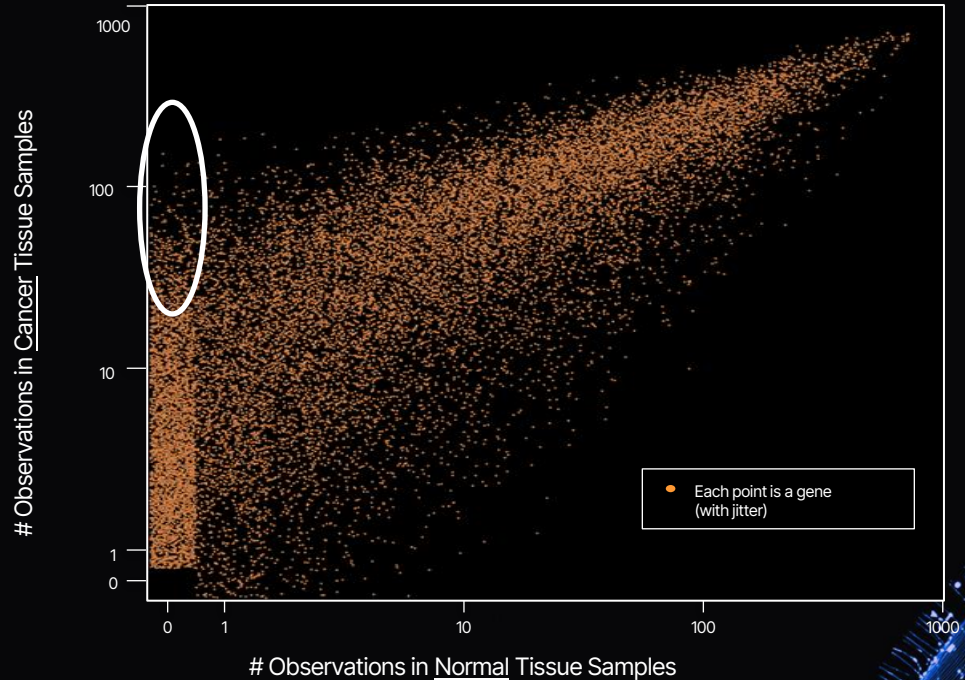
Chromatographic Retention Time



Signal intensities of peptide fragments

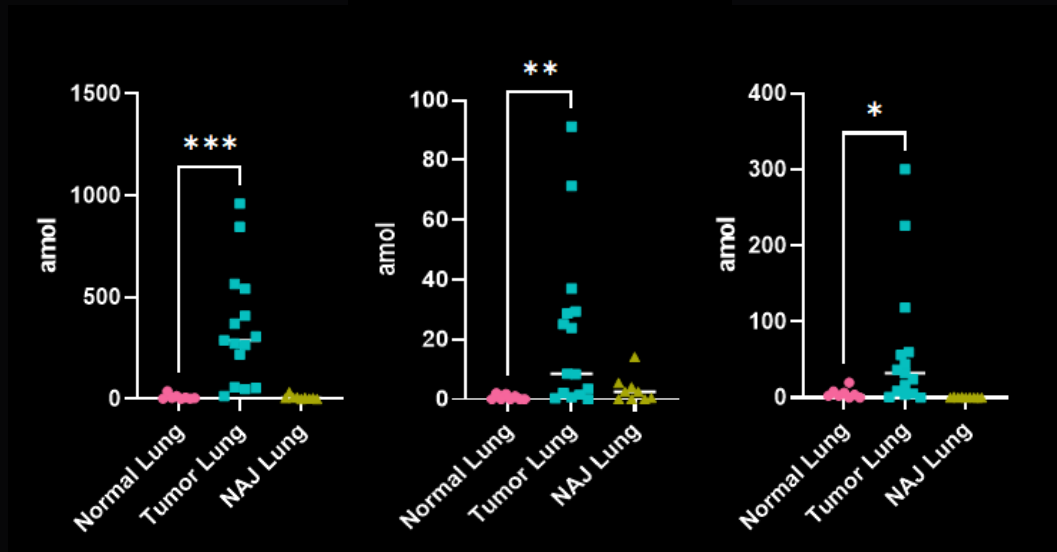


Identification of novel, tumor-specific peptides



Targets are Validated with High-precision Mass Spectrometry

Lung Squamous Cell Carcinoma:
Example peptides validated using synthetic controls



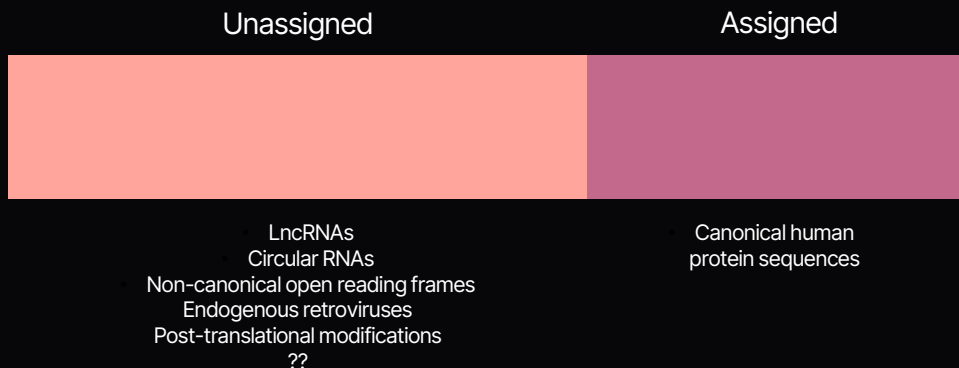
Validated targets are candidates for TCR-based therapies

Underway: Novel *in silico* approaches to discover and enhance TCRs

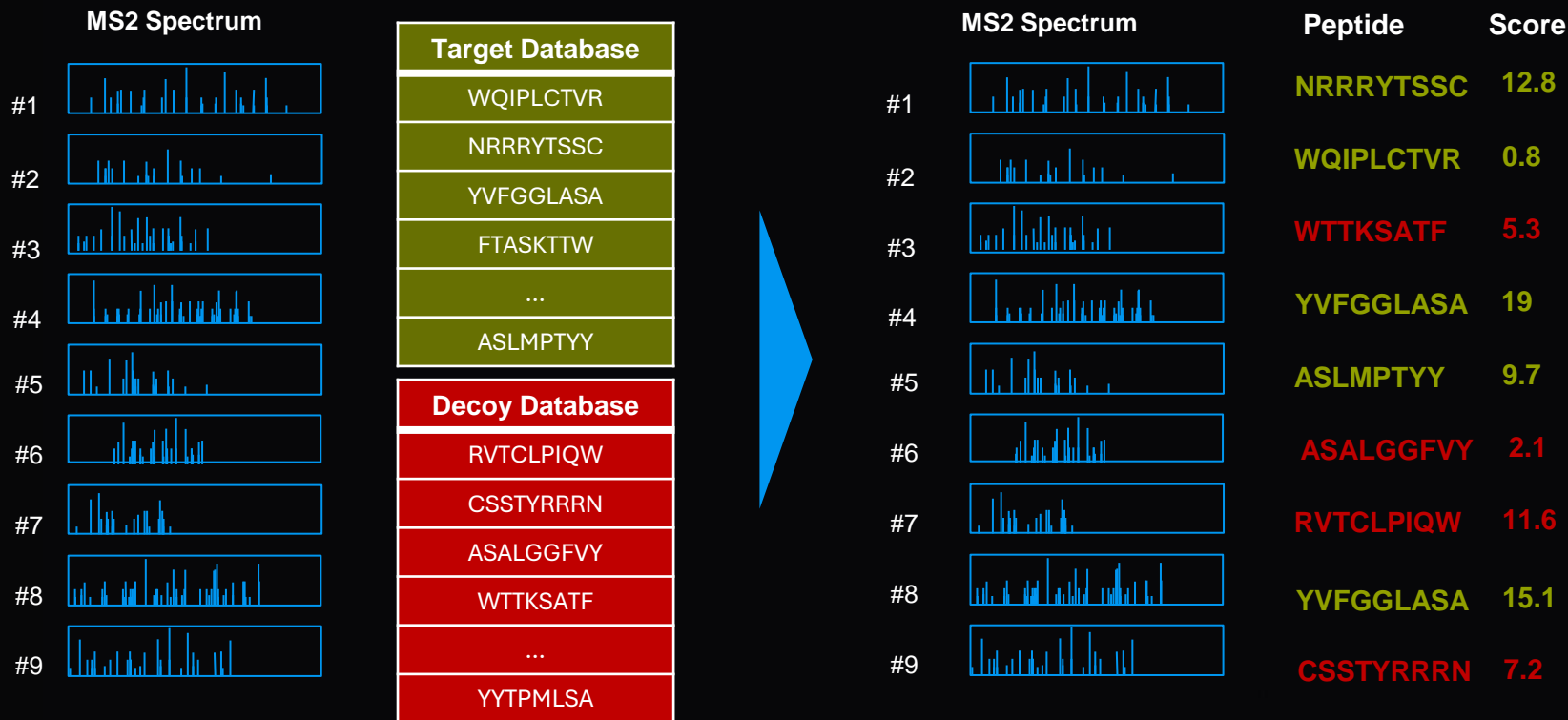
Source: Internal

The Challenge of HLA Mass Spectrometry

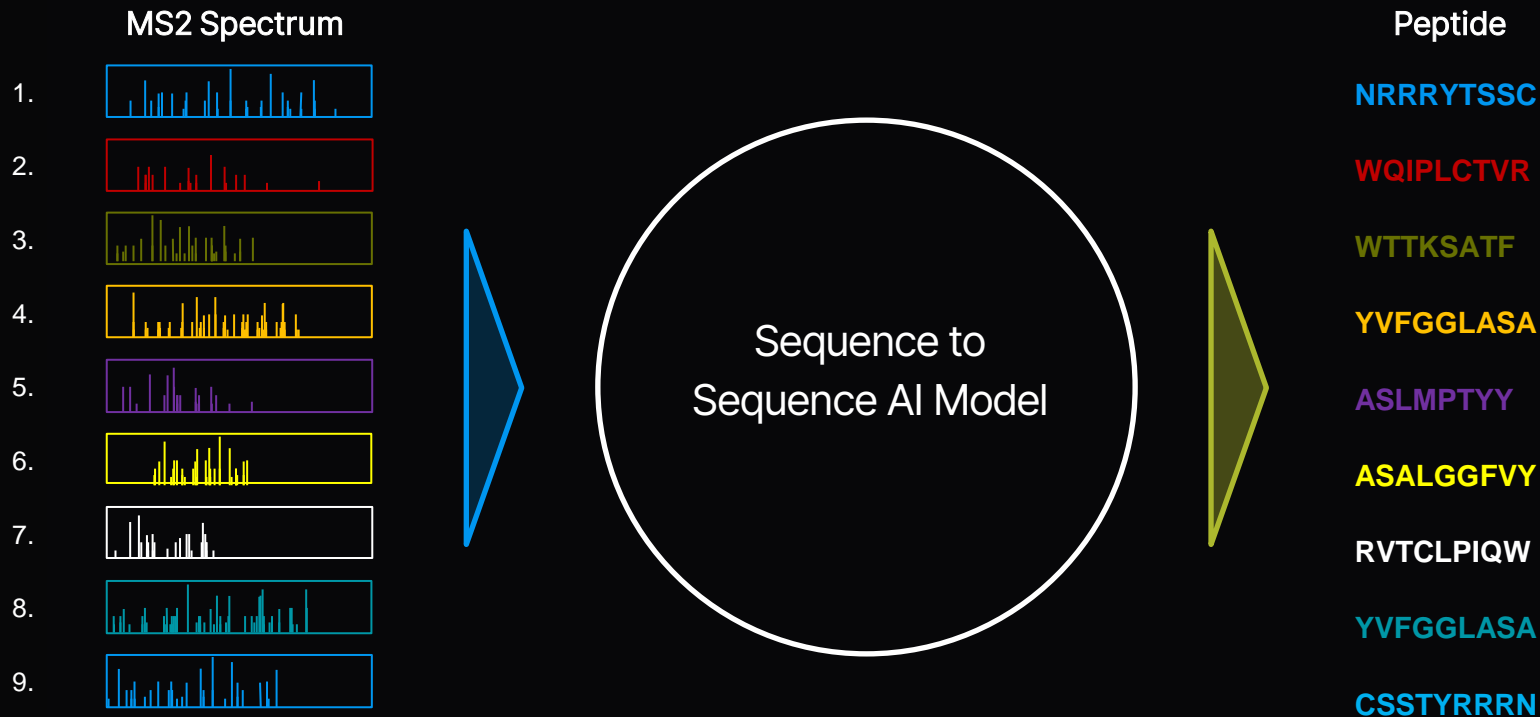
55-75% of data cannot be mapped to a known human peptide



Traditional Mass Spectrometry Target Decoy Search



De Novo Peptide Sequencing



InstaNovo - *De Novo* Peptide Sequencing with Deep Learning

The approach

De novo peptide sequencing using deep learning. No database needed.

The dataset

Model trained on **28 million labeled spectra** matched to **742 thousand human peptides** from ProteomeTools project.

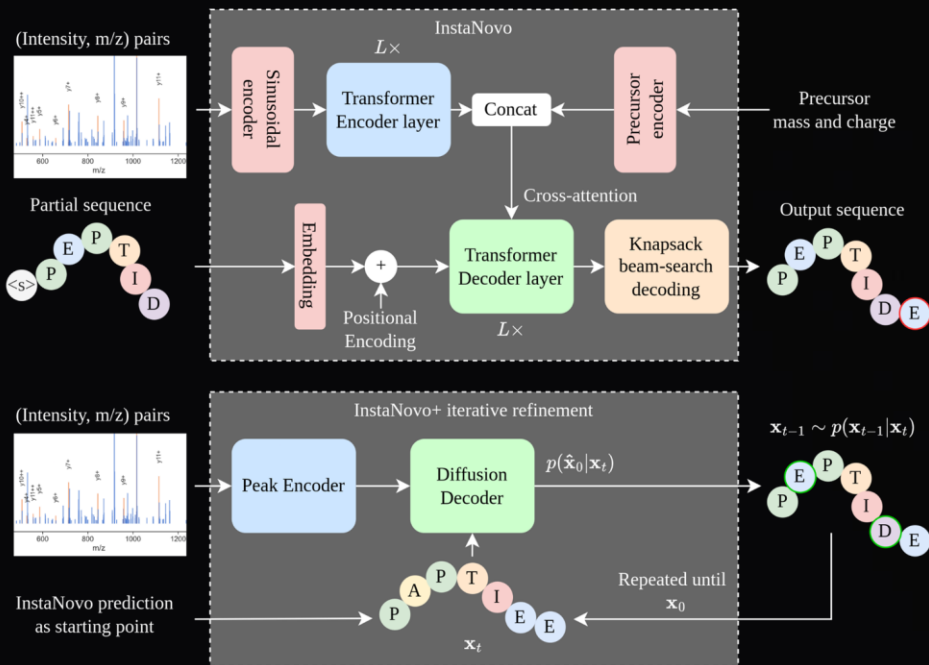
The models

InstaNovo

Autoregressive encoder-decoder transformer model with special MS2 spectrum encoder

InstaNovo+

Multinomial diffusion model to further improve performance using iterative refinement



Source: *De novo peptide sequencing with InstaNovo: Accurate, database-free peptide identification for large scale proteomics experiments* (<https://www.biorxiv.org/content/10.1101/2023.08.30.555055v3>)



InstaNovo - *De Novo* Peptide Sequencing with Deep Learning

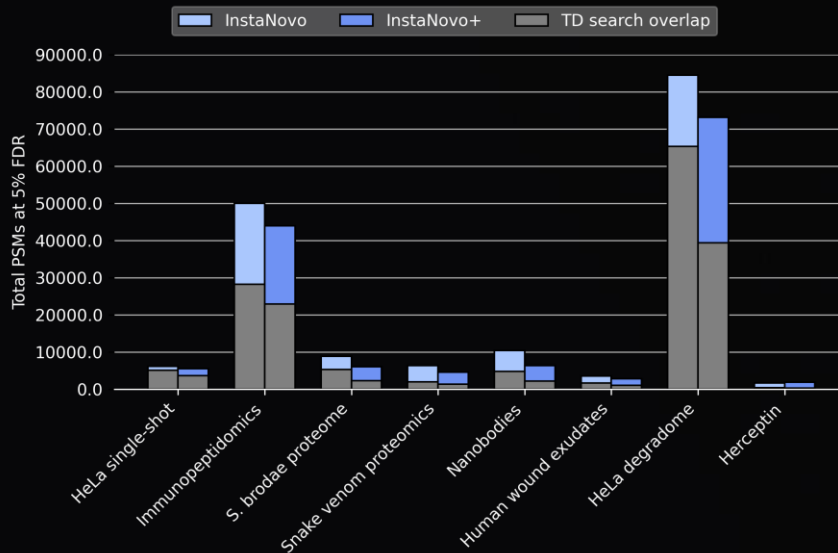
The results

- InstaNovo has performed well across most datasets
- Increases PSM rate in HeLa proteomes
- Expanded an immunopeptidomics dataset by 42%
- Found peptides from individual-specific mutations, splice variants, and post-translational modifications.
- Discovered new HLA peptides in immunopeptidome experiments

Preprint and code available

Preprint is available on **BioRxiv**
<https://www.biorxiv.org/content/10.1101/2023.08.30.555055v3>

Code is available on **GitHub**
<https://github.com/instadeepai/instanovo>



Source: *De novo peptide sequencing with InstaNovo: Accurate, database-free peptide identification for large scale proteomics experiments* (<https://www.biorxiv.org/content/10.1101/2023.08.30.555055v3>)



Step 4:

Protein Design: RiboMab™ Platform

Press Release

BioNTech and InstaDeep Announce Strategic Collaboration and Form AI Innovation Lab to Develop Novel Immunotherapies

25 November 2020

The strategic collaboration will focus on three core areas:

- **Novel Drug Design:** BioNTech is advancing a pipeline of novel mRNA-based vaccines and therapeutics and will apply InstaDeep's DeepChain™ protein design platform to engineer new mRNA sequences for protein targets, including for its RiboMab™ and RiboCytokine™ platforms, which use messenger RNA to encode antibodies and cytokines *in vivo*.

Source: <https://investors.biontech.de/news-releases/news-release-details/biontech-and-instadeep-announce-strategic-collaboration-and-form/>

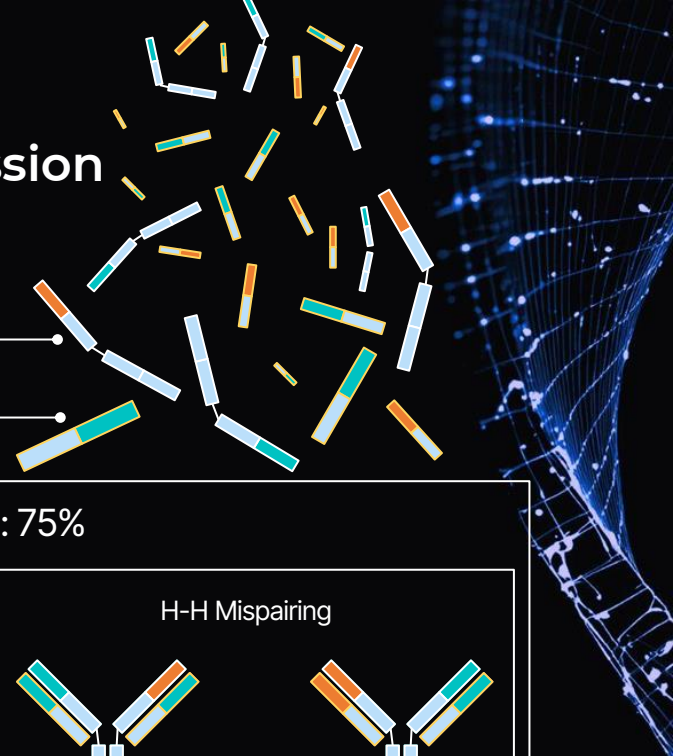
Enabling Antibody Co-Formulation / Co-Expression

Co-expressed and bi-specific antibodies hold significant therapeutic interest.

However, these require precise pairing of heavy and light chains.

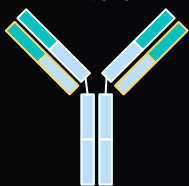
Heavy chain

Light chain



Correct Pairing

Antibody A
12.5%



Antibody B
12.5%



Functional form of co-expressed antibodies.

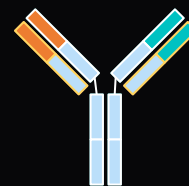
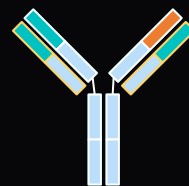
Mispairing: 75%

H-L Mispairing



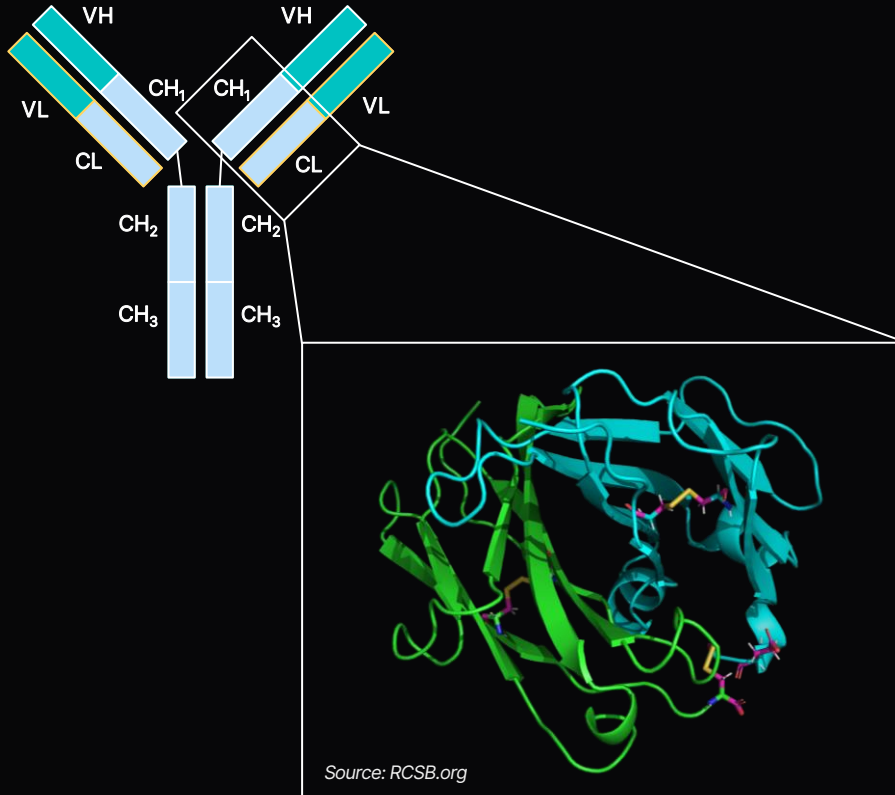
Our goal: engineer the H-L interface to prevent mispairing.

H-H Mispairing

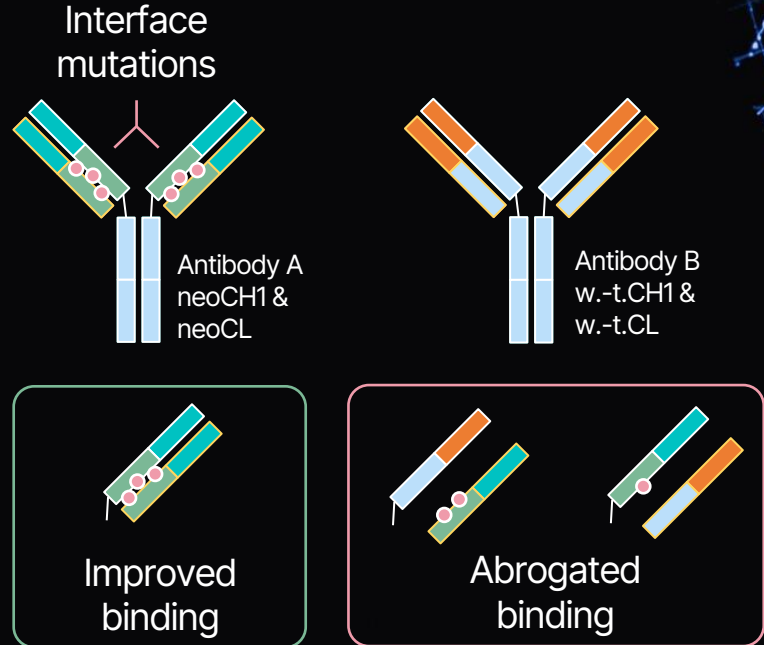


H-H mispairing out of scope.

Our Protein Engineering Approach



We set out to introduce mutations to enforce orthogonality between neoCH1—neoCL v.s. wild-typeCH1—wild-typeCL.



Finding Optimal Orthogonal Mutations

A combinatorial multi-objective optimisation problem we could solve thanks to our DeepChain™ platform and an efficient *in silico* – *in vitro* feedback loop.

Binding energy estimations

For all correctly paired and mispaired complexes.

Mitigation of thermostability changes

For all both heavy and light chains.

Structural modelling

Of each interface mutation.

Key interaction understanding

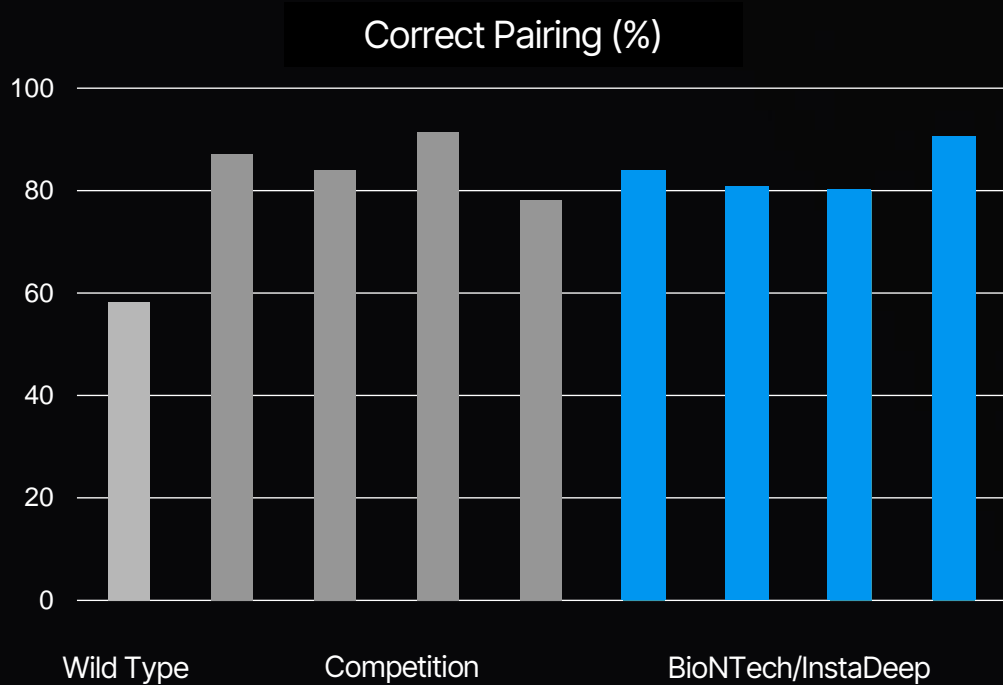
In depth knowledge of the physics of the interface.

RiboMab™ Results obtained with DeepChain™

>90% correct pairing

Matching the best patented designs on the market.

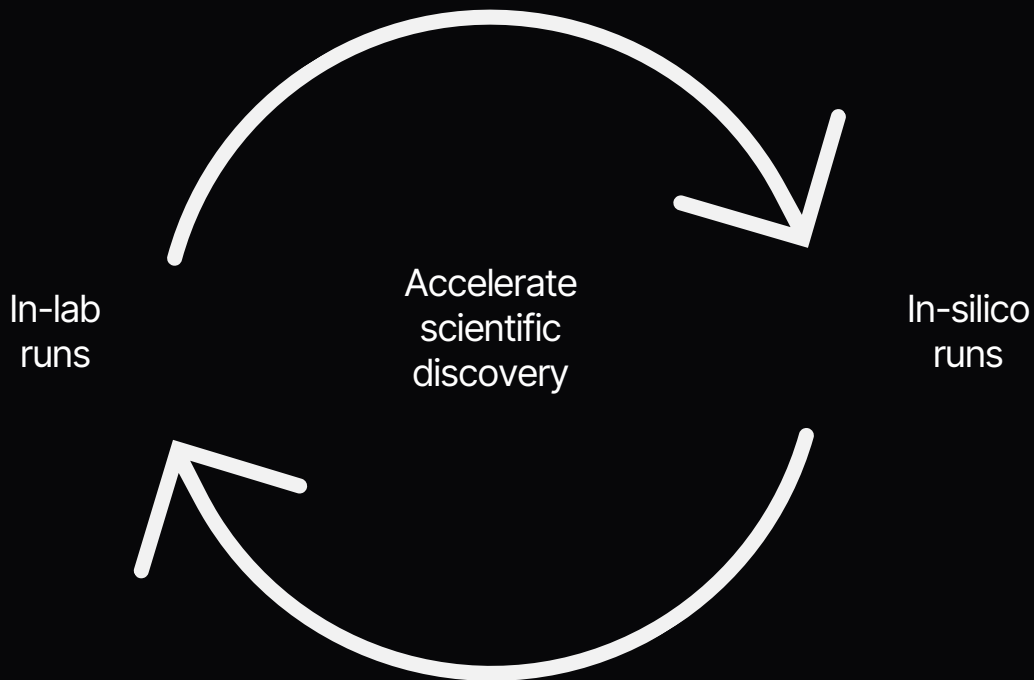
Functional activity of antibodies confirmed.



Source: Internal Data.

Step 5: Lab Automation

Lab Automation Could Transform Research & Development



Challenges

Change

R&D is changing constantly

Complexity

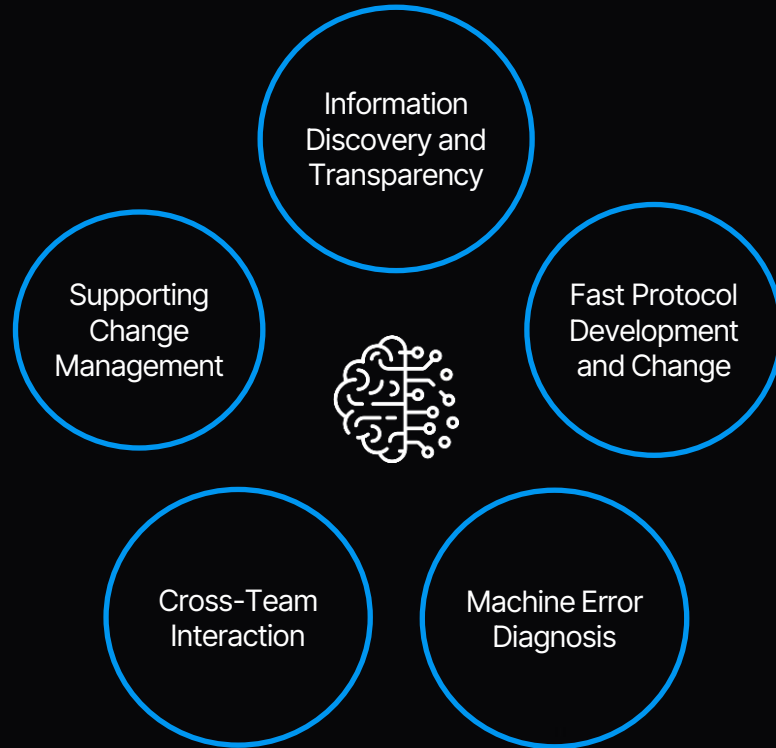
High complexity of science and automation

Transparency

Scientists need transparency and trust

Opportunities to Unlock Full Lab Automation with AI

With the assistance of Artificial Intelligence, we see opportunities to **overcome these challenges** and unlock the **full potential of laboratory automation**

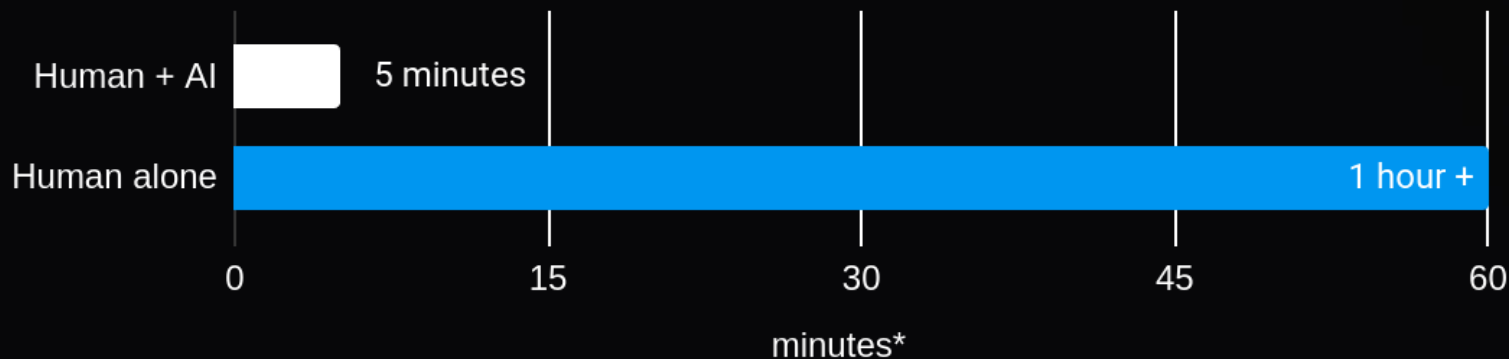


Demo



Increased Efficiency Across Laboratory Research Activities

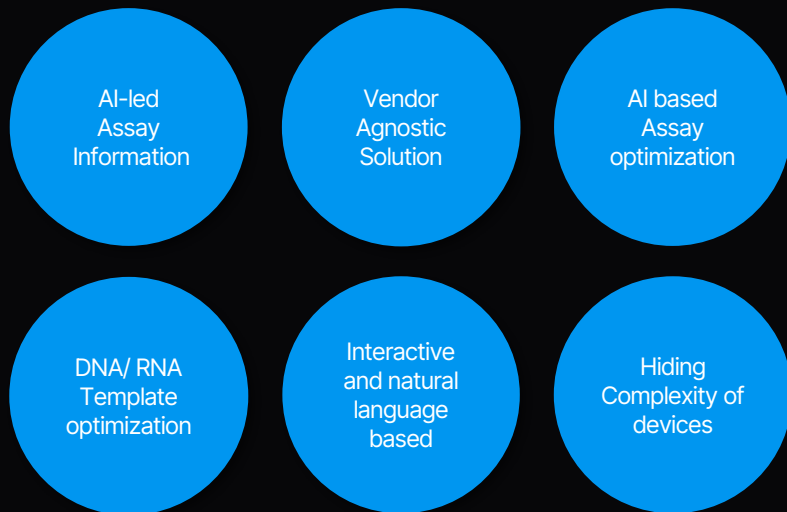
Case Study: liquid handler error diagnosis



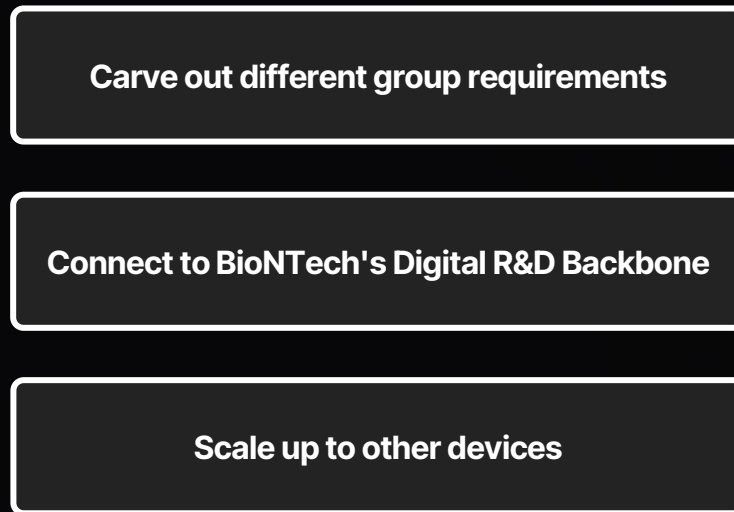
*Stated times are for an unexpected error

Lab Automation Demo – Future Outlook

Technology capability established



Next Steps



One more thing...



AI Day Executive Summary



Ugur Sahin
Founder & CEO
BioNTech



Ryan Richardson
Chief Strategy Officer
BioNTech



Karim Beguir
CEO
InstaDeep

Thank You!
THE END

