# Workshop in Systems Immunology



*Emanuele de Rinaldis*
*Magnus Fontes*
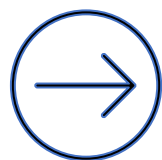*Shameer Khader*
*Giorgio Gaglia*

*June 19th 2023*

# Today's Plan

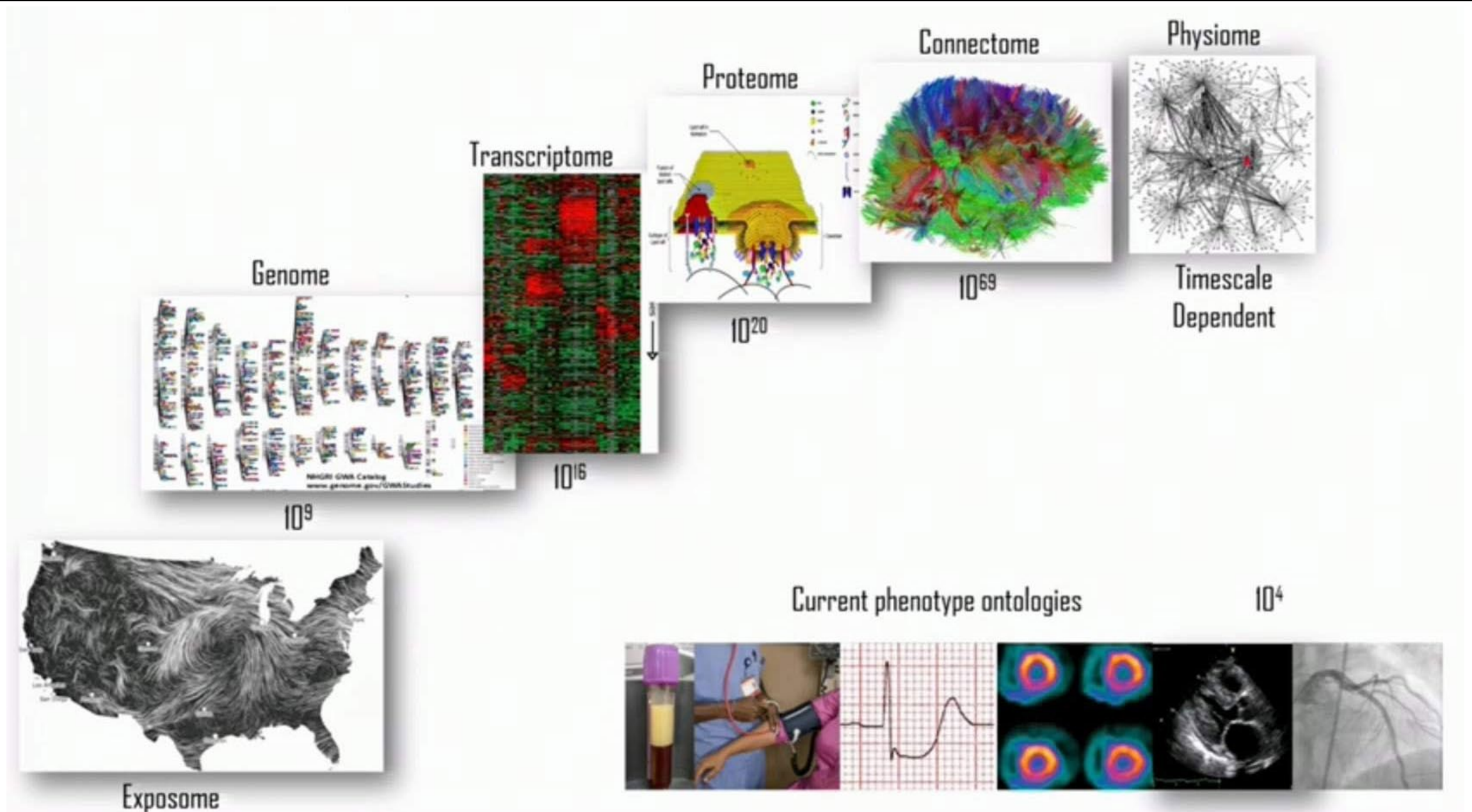| Time | Session |
|---|---|
| 8:00-8:10 am | Course Overview and Objectives – Emanuele de Rinaldis, PhD |
| 8:10-9:00 am | Introduction to Systems Immunology – Emanuele de Rinaldis, PhD |
| 9:00-9:45 am | Systems Immunology & Immune Oncology: A Data-Centric View – Magnus Fontes, PhD |
| 9:45-10:00 am | Break |
| 10:00-11:30 am | Deep Dive Into Selected Scientific Case Studies: From Systems Immunology to Novel Therapeutic Insights – Emanuele de Rinaldis, PhD |
| 11:30 am-12:00 pm | Q/A and Panel Discussion |
| 12:00-1:00 pm | Break for Lunch |
| 1:00-2:00 pm | Spatial Biology Methods and Analytics for Immunology & Oncology – Giorgio Gaglia, PhD |
| 2:00-2:15 pm | Break |
| 2:15-3:30 pm | Artificial Intelligence – A Primer for Immunologists – Shameer Khader, PhD, MPH |
| 3:30-3:45 pm | Break |
| 3:45-4:45 pm | Interactive Data Analysis Session – Magnus Fontes, PhD |
| 4:45-5:00 pm | Wrap Up Notes & Final Remarks |

# AI for Immunologists
# – An Introduction

- Background

- Data boom in biology and need for AI

- Examples of AI in Immunology
  - **Classical ML and Predictive models – Open Targets / Target Immune Engine**
  - **Graph ML – AsthmaGraph (Poster at FOCIS!)**
  - **Emerging themes in AI: Encoders, Embedding, Transformers, GANs, and LLMs**

- Future outlook

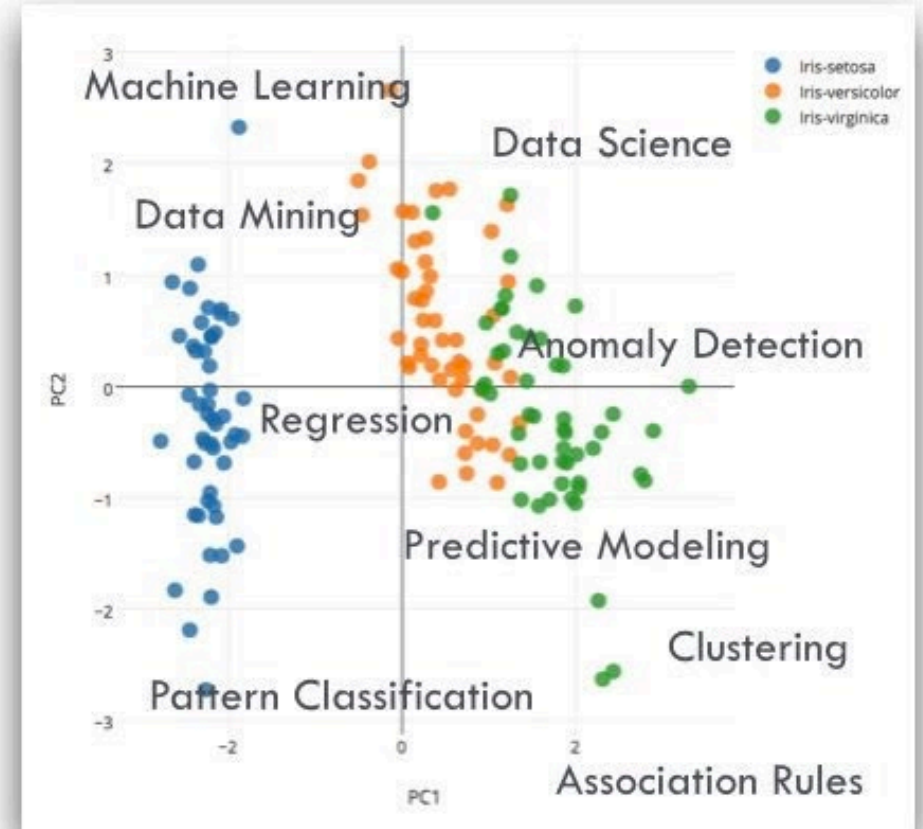*From Calium MacRae / https://twitter.com/daniel_kraft/status/1011692279445123072*

***Big data*** *(noun) extremely large data sets that may be analyzed computationally to reveal patterns, trends, and associations, especially relating to human behavior and interactions.*

***Predictive analytics*** *is the area of data mining concerned with forecasting probabilities and trends.*

***Data science*** *is an interdisciplinary field about processes and systems to extract knowledge or insights from data.*

***Artificial intelligence (AI)*** *is wide-ranging branch of computer science concerned with building smart machines capable of performing tasks that typically require human intelligence.*
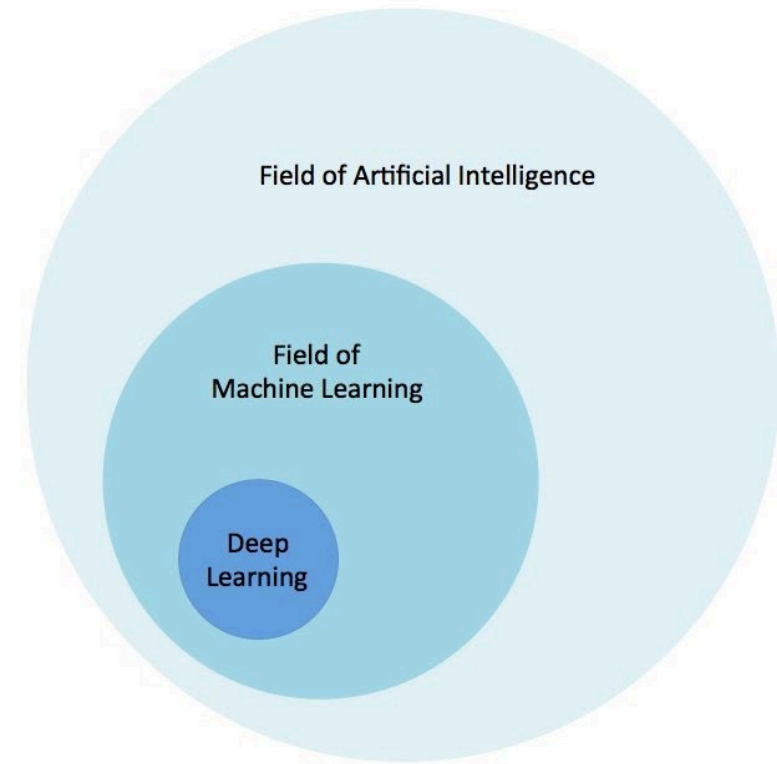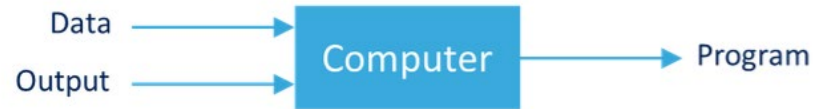
Disease networks

*No disease is an island!*
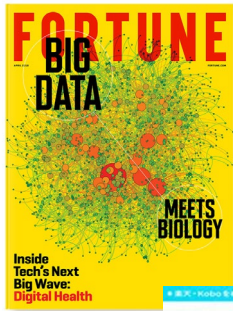
Protein-Nucleic acid

Protein-Protein

Protein-Small molecule

**AI in Biomedicine: Convergence of Big Data, Predictive Modeling, Data science and AI to design, develop and deliver Precision Medicine solutions**

# From Biology to Therapy to Healthcare via Data & AI

## Medicine

- Traditional data types
- Centralized
- GBs or TBs in size
- Structured
- Stable data model
- Low-dimensional
- Statistical approaches
- Cohort size (~10K)
- Hypothesis-driven

## Data-driven Precision Medicine

- Evolving data types
- Decentralized
- Petabytes, exabytes…
- Semi or unstructured
- Evolving, flat data model
- High-dimensional
- Machine or deep learning
- Large cohort size (>10K)
- Data-driven

Healthcare

Biology

Chemistry

Pharmacology

Medicine

*K. Shameer et. al(2018) Machine learning in Cardiovascular Medicine, Are we there yet? BMJ Heart pii: heartjnl-2017-311198. doi: 10.1136/heartjnl-2017-311198.*

# Enablers of AI



**Data from nursing flow sheets!** → **Text-based description of neurological status of patients!** → **Mapping to AVPU Schema**

| | | |
|---|---|---|
| Alert | | 339,899 |
| Voice | 46,588 | |
| Pain | 16,963 | |
| Unresponsive | 3,448 | |
| Missing | 7,616 | |

## PheKB

- **Access Validated Phenotype Algorithms**
- **Collaborate on Phenotype Algorithms**
- **Share Validated Phenotype Algorithms**

One-stop documentation and versioning of validated phenotype algorithms
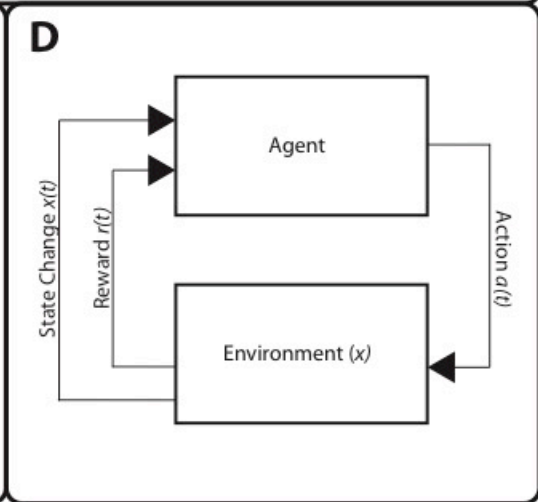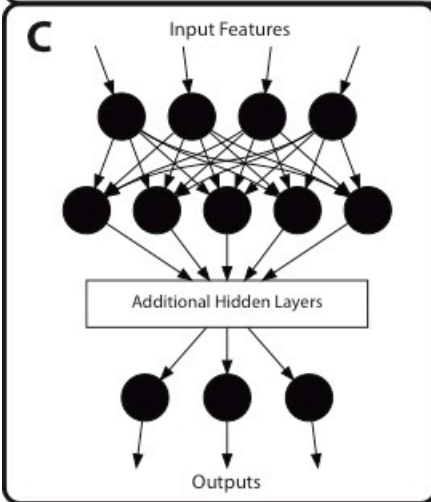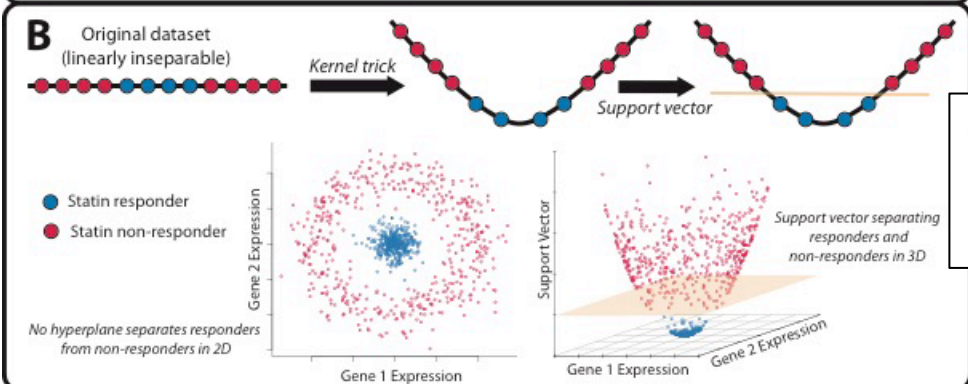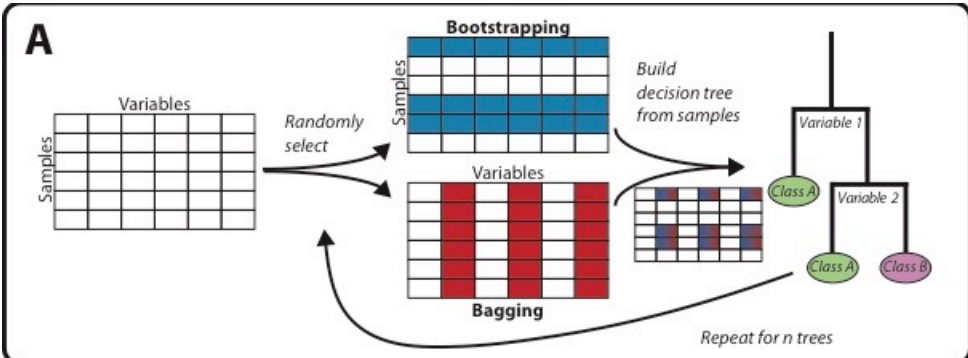
Tailored searches for algorithms applicable to your EMR system

Validate existing phenotype algorithms on your EMR

Receive feedback and additional validation

Publicize your work to better find collaborators

Receive feedback and validation of your algorithm

### Location specific data
- Climate, allergies
- Pollution, crime
- Social networks & characteristics
- Food, walkability

**Environment Repository**

### Health monitoring devices
- Fitness, sleep and heart rate monitoring
- Blood pressure and glucose monitoring
- Diet
- Current prescription and OTC medications

**Personal Health Repository**

### Electronic Health Record
- Vital signs, labs
- Family history and medications
- Adverse events, Biomarkers, imaging and Biopsy
- Multi-omic data

**Clinical Data Repository**

Data interoperability using health and clinical data application program interface

Data transfer and exchange according to PHI, HIPAA, HL-7 and other regulatory codes

| Diagnostic alerts | Predictive models | Data-driven clinical trials |
|---|---|---|

**Actionable recommendations**

## Concept Classes

| | |
|---|---|
| 🟥 Emergent care event | 🟪 Symptoms/Disorders | ⬜ Suppressors |
| 🟦 Emergent care facilities | 🟩 Unplanned event triggers | 🟧 Negation / Family History |

### Step 1: Text Processing

Over 300,000 Clinical Notes → Tagger → **Terminology**

$T_1$ emergency dept;
$T_2$ emergent care
$T_3$ urgent care
$T_4$ ER
$T_5$ ED
...
$T_n$ Mercy Hospital
$T$ Mercy

### Retrieved Snippets

S-1: "Mr. X presented to the Santa Clara Valley Medical Center Emergency Room with a complaint of left-sided chest pain."

S-2: "Ms. X had coffee ground emesis and went to the ER of CPMC at which time she claims was almost pulseless required …"

S-3: "Repeat ER/PR and HER2 testing was done on the cellblock from the UCSF FNA and showed that ER was positive, but PR …"

First 30 days through one year post-Dx

### Step 2: Candidate Event Matrix

### Step 3: Filtering

## Mount Sinai Data Flow

EPIC & other data capturing systems

Mount Sinai Data Warehouse | BioME & other molecular data

| Structured data | Semi-structured | Unstructured data |
|---|---|---|
| Normalization ! | Text-mining ! | NLP! |

**Mount Sinai Health Base**

**Mount Sinai HealthIO**
Big data analytics and predictive modeling

**EHDViz**
High-dimensional and real-time data visualization

## Input / Output

Electronic Health Record (EHR) Data | Biomedical Monitor Data | Wearable Technology Data

Read Data Into Memory — De-identify Data

- Clinical Phenotypes (diagnoses and procedure codes)
- Medications (normalized using RxNorm and NDC)
- Corpus of Clinical Notes (compiled using 86 note-types)
- Laboratory Measurements (normalized and QC'd data)
- Genomic Profiling (genetic and structural variants)
- Wellness Data (biomedical devices and technologies)

Machine Learning / Risk Algorithms

ggplot2 / gridExtra

Select Visualization Style

Merge and Normalize Data Streams — reshape2

EHR-agnostic visualization

Population Health Management

*Badgeley MA & Shameer K, et.al; . BMJ Open. 2016 Mar 24;6(3):e010579. doi: 10.1136/bmjopen-2015-010579*
*Brief Bioinform. 2016 Feb 14. pii: bbv118. PubMed PMID: 26876889.*

# Key AI methods and Applications in Immunology



A) Random Forests
B) Support Vector Machines
C) Convolutional Neural Network
D) Reinforcement leaning

DeepImmuno: deep learning-empowered prediction and generation of immunogenic peptides for T-cell immunity

Guangyuan Li, Balaji Iyer, V.B. Surya Prasath, Yizhao Ni and Nathan Salomonis

Corresponding author. Guangyuan Li, University of Cincinnati, 3333 Burnet Ave, MLC7024, Cincinnati, OH 45267, USA. Tel: 15138031584; E-mail: li2g2@mail.uc.edu

**Significance**

Applying artificial intelligence tools to a highly complex question of immunology, we show that a deep neural network can learn to predict the patterns of chromatin opening across 81 stem and differentiated cells across the immune system, solely from the DNA sequence of regulatory regions. It does so by discovering ab initio the binding motifs for known master regulators, along with some unknown ones, and their combinatorial operation. These predictions validated biochemically, and a mouse-trained neural network predicts human enhancer/promoter activity much better than sequence comparisons would. Beyond serving as a trove of testable functional frameworks, this work is important in showing how massively complex integrated questions of immunology can be handled with such tools.
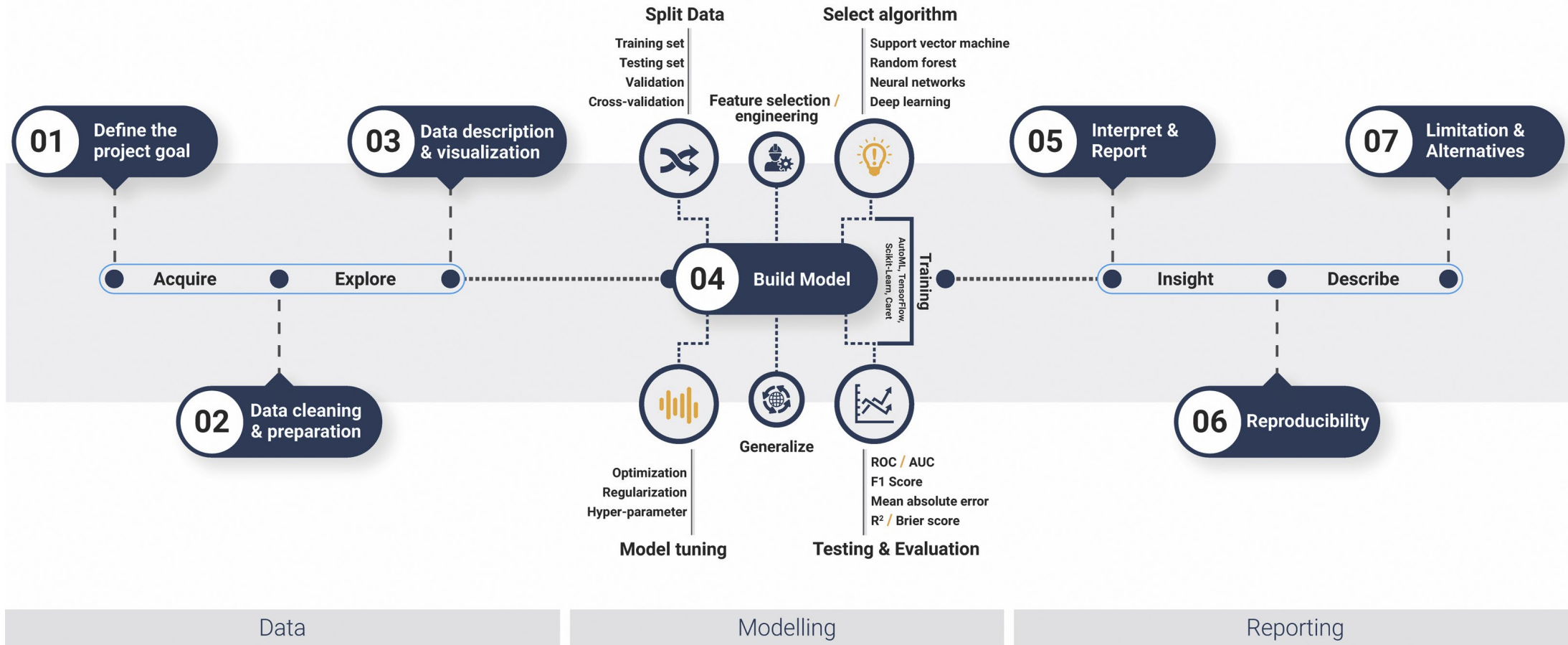
https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7924715/
https://journals.plos.org/plosone/article?id=10.1371/journal.pone.0011267
https://www.pnas.org/doi/10.1073/pnas.2011795117

# Integration of mechanistic immunological knowledge into a machine learning pipeline improves predictions

Anthony Culos[1,2,10], Amy S. Tsai[1,10], Natalie Stanley[1,2], Martin Becker[1,2], Mohammad S. Ghaemi[1,2,3], David R. McIlwain[4], Ramin Fallahzadeh[1,2], Athena Tanada[1,2], Huda Nassar[1,2], Camilo Espinosa[1,2] Maria Xenochristou[1,2], Edward Ganio[1], Laura Peterson[1,5], Xiaoyuan Han[1], Ina A. Stelzer[1], Kazuo Ando[1], Dyani Gaudilliere[1], Thanaphong Phongpreecha[1,2,6], Ivana Marić[1,5], Alan L. Chang[1,2], Gary M. Shaw[5], David K. Stevenson[5], Sean Bendall[6], Kara L. Davis[5], Wendy Fantl[4,7,8], Garry P. Nolan[6], Trevor Hastie[2,9], Robert Tibshirani[2,9], Martin S. Angst[1,11], Brice Gaudilliere[1,5,11] and Nima Aghaeepour[1,2,5,11] ✉

The dense network of interconnected cellular signalling responses that are quantifiable in peripheral immune cells provides a wealth of actionable immunological insights. Although high-throughput single-cell profiling techniques, including polychromatic flow and mass cytometry, have matured to a point that enables detailed immune profiling of patients in numerous clinical settings, the limited cohort size and high dimensionality of data increase the possibility of false-positive discoveries and model overfitting. We introduce a generalizable machine learning platform, the immunological Elastic-Net (iEN), which incorporates immunological knowledge directly into the predictive models. Importantly, the algorithm maintains the exploratory nature of the high-dimensional dataset, allowing for the inclusion of immune features with strong predictive capabilities even if not consistent with prior knowledge. In three independent studies our method demonstrates improved predictions for clinically relevant outcomes from mass cytometry data generated from whole blood, as well as a large simulated dataset. The iEN is available under an open-source licence.

# Designing a classical machine learning project: key steps
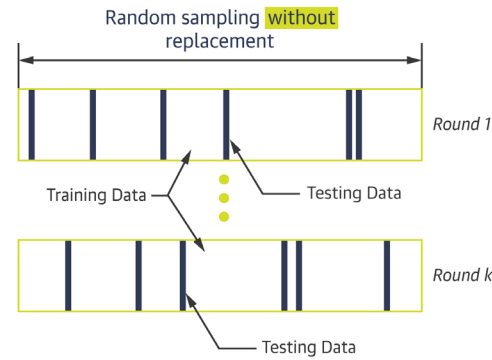
# Feature engineering: Key concepts

Feature Engineering &
Learning Approaches

Hand-Engineering
(e.g., LBP, SIFT)

Dimensionality Reduction
(e.g., PCA, ICA, ISOMAP)

Deep Learning
(e.g., Deep Convolutional auto-encoder)

# Supervised vs. Unsupervised Learning

**Best Practice Strategies**

- Choose Tall/Thin Data
- Remove Redundant Features
- Deal with Outliers
- Impute Missing Values
- Address Class Imbalance

- Optimize Bias-Variance Tradeoff
- Employ Bootstrapping or k-fold CV for Large Datasets
- Use LOOCV for Small Datasets

- Release Code under Open-source GPL License
- Release a Code Book
- Provide Code & Data as Supplements
- Use Docker or Sphinx to Generate Code-embedded S/W Manuals

**Checklist**

- Format/Describe Data
- Ensure Data is Clean
  - Normalize Variables
  - Impute Missing Data
  - Remove Outliers
  - Balance Class
- Describe Feature Selection

- Optimize Model Parameters
- Train Model using CV
- Define Ensemble Methods (if used)
- Identify Model Interpretability

- Consider Sharing Code or Scripts on Public repositories
- Provide a Data Dictionary
- Document Detailing Software & Libraries

**ML Pipeline**

Data Standardization

Model Assessment

Model Replicability

1  2  3  4  5  6  7

Designing Study Plan

Selecting ML Models

Model Evaluation

Reporting Limitations

Describe the Study Plan
- Need for ML Models
- Goals/Objectives
- Summary Statistics
- ML Workflow

- Define Analysis Goal
- Identify ML Task
- Define use of Simple /Complex Models
- Benchmark Complex Models

- Clearly Define Training Test & Validation Sets
- Provide Summary of Model Parameter's
- Report Class Balancing Measures for Classification Tasks

- Report all Assumptions, Biases & Limitations
- Provide Performance Metrics on Hold-out or External Validation Set

**Checklist**

- Identify the Need for Using ML
- Identify Input/Output Variables
- Identify Biases in the Data
- Define Steps in ML Pipeline

- Select Simple Models First
- Benchmark Complex Models
- Optimize Model by Tuning Hyperparameters

- Include Bland-Altman Plots
- Report Inter-/Intra-Observer Variability
- Convey Misclassification Risk
- Report Balanced Class Accuracies

- Check & Report all Assumptions
- Evaluate the Model using External Validation Dataset
- Justify the Use of More Complex Models

**Best Practice Strategies**

BMC Bioinformatics

http://snapshots.cell.com/

# Application of Knowledge Graphs in Immunology



Franck Rapaport, Travis Ahn-Horst, Emanuele de Rinaldis & Shameer Khader

| Approach | Domain | Model | Prediction Task | Entites | Relations | Entity Types | Relation Types | Num Datasets in Graph |
|---|---|---|---|---|---|---|---|---|
| Decagon [167] | Drug-Drug Interactions | Relational GCN with tensor factorisation decoder | Link Prediction | 19.6K | 5.3M | 2 | 964 | ≈7 |
| TriModel [102] | Drug-Target Interactions | Tensor factorisation | Link Prediction | 5K | 12K | 2 | 1 | 1 |
| Rosalind [113] | Disease-Gene Prioritisation | Tensor factorisation | Link Prediction | 319K | 2.6M | 5 | 11 | ≈15 |

**Preclinical validation of therapeutic targets predicted by tensor factorization on heterogeneous graphs**

Saee Paliwal[1], Alex de Giorgio[2], Daniel Neil[1], Jean-Baptiste Michel[1] & Alix MB Lacoste[1]

**Link prediction, embeddings & algorithms**

**Graph Neural Networks for Multirelational Link Prediction**

Decagon is a graph convolutional neural network for multirelational link prediction in heterogeneous graphs.

### Target identification

### Drug repositioning

### Target validation

**Graph integration**

## Pre-processing & filtration

*General*

*Disease-specific*

Filter

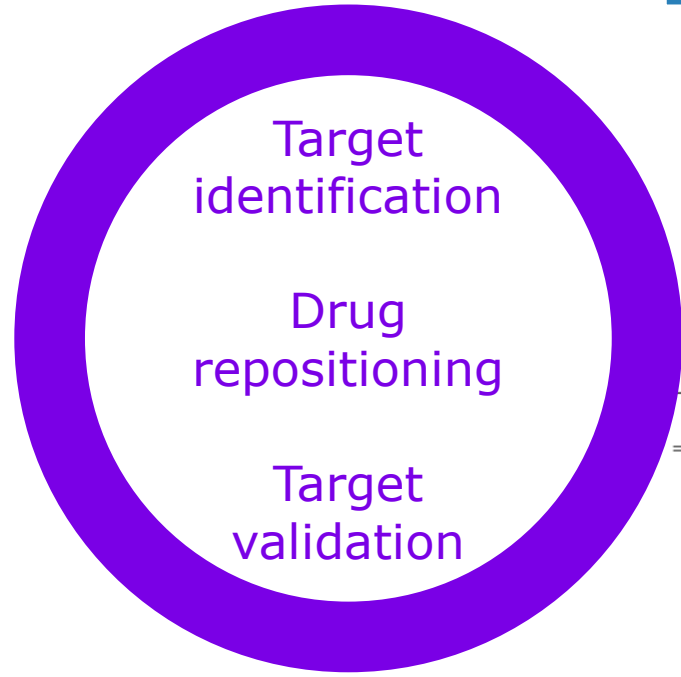An integrative knowledge graph for rare diseases, derived from the Genetic and Rare Diseases Information Center (GARD)

Qian Zhu[1*+], Dac-Trung Nguyen[1+], Ivan Grishagin[1], Noel Southall[1], Eric Sid[2] and Anne Pariser[2]

*Public knowledge graphs*

| KG Dataset | Link | Entities | Triples |
|---|---|---|---|
| Hetionet [56] | https://het.io/ | 47K | 2.2M |
| DRKG [65] | https://github.com/gnn4dr/DRKG | 97K | 5.7M |
| BioKG [151] | https://github.com/dsi-bdi/biokg | 105K | 2M |
| PharmKG [164] | https://github.com/MindRank-Biotech/PharmKG | 7.6K | 500K |
| OpenBioLink [14] | https://zenodo.org/record/3834052 | 184K | 4.7M |
| Clinical Knowledge Graph [124] | https://data.mendeley.com/datasets/mrcf7f4tc2/1 | 16M | 220M |

*Multi-omics inferred networks*

Phenotype — Phenotypic similarities
Scale — Functional similarities
Pathway co-membership
Protein interactions
Tissue-specific co-expression
Genotype — Genetic interactions

ARTICLE

Network analysis reveals rare disease signatures across multiple levels of biological organization

Pisanu Buphamalai[1,2], Tomislav Kokotovic[1,3,4], Vanja Nagy[1,3,4] & Jörg Menche[1,2,5,6]

## Generative AI Timeline

**Model Type**
- Variational Autoencoder
- Generative Adversarial Network
- Autoregressive / Transformer
- Normalizing Flow
- Energy-Based / Diffusion Model
- Multimodal Model

**2014**
- VAE

**2015**
- GAN
- CGAN
- GRU

**2016**
- VAE-GAN
- DCGAN
- PixelRNN
- PixelCNN
- RealNVP

**2017**
- pix2pix
- WGAN
- CycleGAN
- PixelCNN++
- Transformers

**2018**
- VQ-VAE
- MuseGAN
- ProGAN
- World Models
- SAGAN
- GPT
- GLOW

**2019**
- BigGAN
- StyleGAN
- BERT
- Music Transformer
- GPT-2
- MuseNet
- FFJORD
- VQ-VAE-2
- NCSN

**2020**
- StyleGAN2
- T5
- GPT-3
- DDPM
- DDIM

**2021**
- VQ-GAN
- Vision Transformer
- CLIP
- DALL.E
- GPT-Neo
- GPT-J
- Codex
- StyleGAN3

**2022**
- ViT VQ-GAN
- Megatron-Turing NLG
- Gopher
- StyleGAN-XL
- LaMDA  GPT-NeoX
- PaLM  Chinchilla
- OPT
- BLOOM
- Latent Diffusion
- GLIDE
- DALL.E 2
- Flamingo
- Imagen
- Parti
- Stable Diffusion

**2023**
- ChatGPT
- Toolformer
- LLaMA
- MusicLM
- ControlNet
- Visual ChatGPT
- MUSE
- Dreamix
- PaLM-E
- GPT-4

**(Selected) Emerging themes in AI**
- Neural network
- Pre-training, fine-tuning and transfer learning
- Attention, Embedding, Autoencoders, Transformers

A mostly complete chart of

# Neural Networks

©2019 Fjodor van Veen & Stefan Leijnen    asimovinstitute.org

Input Cell

Backfed Input Cell

Noisy Input Cell

Hidden Cell

Probablistic Hidden Cell

Spiking Hidden Cell

Capsule Cell

Output Cell

Match Input Output Cell

Recurrent Cell

Memory Cell

Gated Memory Cell

Kernel

Convolution or Pool

Perceptron (P)

Feed Forward (FF)

Radial Basis Network (RBF)

Deep Feed Forward (DFF)

Recurrent Neural Network (RNN)

Long / Short Term Memory (LSTM)

Gated Recurrent Unit (GRU)

Auto Encoder (AE)

Variational AE (VAE)

Denoising AE (DAE)

Sparse AE (SAE)

Markov Chain (MC)

Hopfield Network (HN)

Boltzmann Machine (BM)

Restricted BM (RBM)

Deep Belief Network (DBN)

Deep Convolutional Network (DCN)

Deconvolutional Network (DN)

Deep Convolutional Inverse Graphics Network (DCIGN)

Generative Adversarial Network (GAN)

Liquid State Machine (LSM)

Extreme Learning Machine (ELM)

Echo State Network (ESN)

Deep Residual Network (DRN)

Differentiable Neural Computer (DNC)

Neural Turing Machine (NTM)

Capsule Network (CN)

Kohonen Network (KN)

Attention Network (AN)

1. **Input Layer**

2. **Hidden Layers**

3. **Neurons (Nodes)**

4. **Activation Functions**

5. **Parameters (Weights and Biases)**

6. **Loss Function**

7. **Optimization Algorithm**

8. **Output Layer**

9. **Training**

10. **Inference**

# Pre-training, Fine-tuning & Foundation models: Concept & Example

- Pre-training: Pre-training refers to the initial training of a model on a large, diverse dataset to learn general representations of the input data.

- The pre-training process involves training a model on a self-supervised or unsupervised task, where the model learns to predict missing or masked parts of the input data.

- Fine-tuning is the process of taking a pre-trained model and further training it on a specific task or dataset.

- A foundation model is a pre-trained model that serves as the base for further development or fine-tuning in the context of large language models

**Traditional ML** vs **Transfer Learning**

- Isolated, single task learning:
  - Knowledge is not retained or accumulated. Learning is performed w.o. considering past learned knowledge in other tasks
- Learning of a new tasks relies on the previous learned tasks:
  - Learning process can be faster, more accurate and/or need less training data

# Transfer learning framework for cell segmentation with incorporation of geometric features

Yinuo Jin[1,*], Alexandre Toberoff[1,*], Elham Azizi[2,†]
[1]Department of Computer Science, Columbia University, New York, NY, USA
[2]Department of Biomedical Engineering and Irving Institute for Cancer Dynamics, Columbia University, New York, NY, USA
{yj2589, aat2167, ea2690}@columbia.edu

## Abstract

With recent advances in multiplexed imaging and spatial transcriptomic and proteomic technologies, cell segmentation is becoming a crucial step in biomedical image analysis. In recent years, Fully Convolutional Networks (FCN) have achieved great success in nuclei segmentation in *in vitro* imaging. Nevertheless, it remains challenging to perform similar tasks on *in situ* tissue images with more cluttered cells of diverse shapes. To address this issue, we propose a novel transfer learning, cell segmentation framework incorporating shape-aware features in a deep learning model, with multi-level watershed and morphological post-processing steps. Our results show that incorporation of geometric features improves generalizability to segmenting cells in *in situ* tissue images, using solely *in vitro* images as training data.

- Transfer learning is a machine learning technique that allows knowledge gained from one task to be applied to another related task.

- It involves leveraging pre-trained models that have been trained on large datasets for a specific task and then reusing and adapting them for a different but related task.

https://www.biorxiv.org/content/10.1101/2021.02.28.433289v2

- Embedding refers to a technique used in neural networks to represent categorical or discrete variables as continuous, dense vectors.

- Embeddings are commonly used in various domains, including natural language processing (NLP), recommender systems, and computer vision.

- The main idea behind embedding is to map high-dimensional categorical data to lower-dimensional continuous representations, where similar or related categories are closer together in the embedding space.

- Examples of embeddings: Word Embeddings, Sentence Embeddings, Document Embeddings, Image Embeddings, Graph Embeddings, Knowledge Graph Embeddings



https://pubmed.ncbi.nlm.nih.gov/29218877/

- Encoding refers to the process of representing data in a specific format or representation.

- It involves transforming data from its original representation into a different format that is suitable for a particular purpose or task.

- The goal of encoding is often to make the data compatible with a specific learning algorithm or to facilitate efficient processing.

- Examples of encoding: One-Hot Encoding, Label Encoding, Binary Encoding, Ordinal Encoding, Hash Encoding, Target Encoding, Feature Hashing

## Benzy

- 1 year old
- From MD
- Likes to play with water
- Friendly with kids
- Comfortable with other pets
- Unconditional hugs

RESEARCH ARTICLE

## Autoencoder based local T cell repertoire density can be used to classify samples and T cell receptors

Shirit Dvorkin, Reut Levi, Yoram Louzoun*

Department of Mathematics, Bar Ilan University, Ramat Gan, Israel

* louzouy@math.biu.ac.il

## Abstract

Recent advances in T cell repertoire (TCR) sequencing allow for the characterization of repertoire properties, as well as the frequency and sharing of specific TCR. However, there is no efficient measure for the local density of a given TCR. TCRs are often described either through their Complementary Determining region 3 (CDR3) sequences, or theirV/J usage, or their clone size. We here show that the local repertoire density can be estimated using a combined representation of these components through distance conserving autoencoders and Kernel Density Estimates (KDE). We present ELATE–an Encoder-based LocAl Tcr dEnsity and show that the resulting density of a sample can be used as a novel measure to study repertoire properties. The cross-density between two samples can be used as a similarity matrix to fully characterize samples from the same host. Finally, the same projection in combination with machine learning algorithms can be used to predict TCR-peptide binding through the local density of known TCRs binding a specific target.

## Author summary

T cell repertoires contain a vast amount of information on the donors, and can be used to characterize the donor, and apply machine learning algorithms on such repertoires. A limiting factor in the analysis of such repertoire is the lack of a good representation of the T cell receptors. We here propose an autoencoder, named ELATE to present receptors as real vectors. We show that this encoder can be used to characterize both full donors and specific receptors using either supervised or unsupervised methods.

https://www.biorxiv.org/content/10.1101/2021.02.28.433289v2

Michael Widrich*    Bernhard Schäfl*    Milena Pavlović[†,‡]    Hubert Ramsauer*

Lukas Gruber*    Markus Holzleitner*    Johannes Brandstetter*    Geir Kjetil Sandve[‡]

Victor Greiff[†]                          Sepp Hochreiter*,[§]

Günter Klambauer*
*ELLIS Unit Linz and LIT AI Lab,
Institute for Machine Learning,
Johannes Kepler University Linz, Austria
[†]Department of Immunology, University of Oslo, Norway
[‡]Department of Informatics, University of Oslo, Norway
[§]Institute of Advanced Research in Artificial Intelligence (IARAI)

## Modern Hopfield Networks and Attention for Immune Repertoire Classification

### Abstract

A central mechanism in machine learning is to identify, store, and recognize patterns. How to learn, access, and retrieve such patterns is crucial in Hopfield networks and the more recent transformer architectures. We show that the attention mechanism of transformer architectures is actually the update rule of modern Hopfield networks that can store exponentially many patterns. We exploit this high storage capacity of modern Hopfield networks to solve a challenging multiple instance learning (MIL) problem in computational biology: immune repertoire classification. In immune repertoire classification, a vast number of immune receptors are used to predict the immune status of an individual. This constitutes a MIL problem with an unprecedentedly massive number of instances, two orders of magnitude larger than currently considered problems, and with an extremely low witness rate. Accurate and interpretable machine learning methods solving this problem could pave the way towards new vaccines and therapies, which is currently a very relevant research topic intensified by the COVID-19 crisis. In this work, we present our novel method DeepRC that integrates transformer-like attention, or equivalently modern Hopfield networks, into deep learning architectures for massive MIL such as immune repertoire classification. We demonstrate that DeepRC outperforms all other methods with respect to predictive performance on large-scale experiments including simulated and real-world virus infection data and enables the extraction of sequence motifs that are connected to a given disease class. Source code and datasets: *https://github.com/ml-jku/DeepRC*

- Attention in AI refers to a mechanism that enables models to focus on specific parts of input data while performing a task.

- It mimics the selective attention mechanism observed in human cognition, where we prioritize certain information over others.

Figure 1: The Transformer - model architecture.

https://proceedings.neurips.cc/paper/2020/file/da4902cb0bc38210839714ebdcf0efc3-Paper.pdf

- Transformer refers to a type of neural network architecture that has gained significant popularity, particularly in the field of natural language processing (NLP)

- The transformer architecture is designed to process sequential data efficiently, such as sentences, paragraphs, or time series data.

- Transformers employs a mechanism called self-attention or scaled dot-product attention to capture relationships and dependencies between different elements of the input sequence.

- Transformers have achieved remarkable success in various NLP tasks, including machine translation, language generation, sentiment analysis, and text classification.



GPT (Generative Pre-trained Transformers)
- A well-known transformer models is the Generative Pre-trained Transformer (GPT)
- Developed by OpenAI

https://www.biorxiv.org/content/10.1101/2021.02.28.433289v2

arXiv > cs > arXiv:1706.03762

**Computer Science > Computation and Language**

[Submitted on 12 Jun 2017 (v1), last revised 6 Dec 2017 (this version, v5)]

**Attention Is All You Need**

Output
Probabilities

Softmax

Linear

Add & Norm

Feed
Forward

Add & Norm

Multi-Head
Attention

Nx

Add & Norm

Feed
Forward

Add & Norm

Multi-Head
Attention

Nx

Add & Norm

Masked
Multi-Head
Attention

Positional
Encoding

Positional
Encoding

Input
Embedding

Output
Embedding

Inputs

Outputs
(shifted right)

Figure 1: The Transformer - model architecture.

**BERT**
Bidirectional encoder
representations from transformers

**OpenAI**
**GPT-4**
Generative pretrained transformer

**Many others**
T5, LLaMA, Bard,
open source models

| Model | Training | Parameters | Year |
|-------|----------|------------|------|
| BERT | 3.3B words | 340M | 2018 |
| GPT-3 | 500B tokens | 175B | 2020 |
| ChatGPT | 300B words | 1.5B | 2022 |
| LLaMA | 1.4T tokens | 65B | 2023 |

**GLUE (General Language Understanding Evaluation) Benchmark Tasks:**
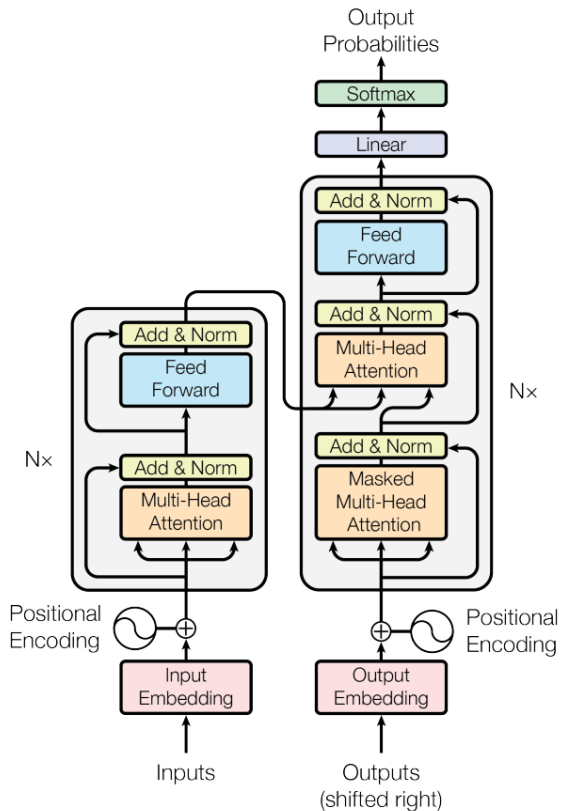
| Task | Example | Dataset | Metric |
|------|---------|---------|--------|
| Grammatical | "This toast is than that one." = Ungrammatical | CoLA | Matthews |
| Sentiment Analysis | "Toy Story 2 was okay." = .543291 (neutral) | SST-2 | Accuracy |
| Similarity | a.) A pride of lions surrounded a monkey. b.) Lions encompassed a monkey. = 4.7 (Very Similar) | STS-B | Person / Spearman |
| Paraphrase | A. Last week, Seattle reported 12 new earthquakes. B. Seattle reported another 12 earthquakes yesterday. = A Paraphrase | MRPC | Accuracy / F1 |
| Question Similarity | a.) How can I cook noodles over a campfire? b.) How do you make Mac & Cheese? = Not Similar | QQP | Accuracy / F1 |
| Contradiction | a.) Glossier products are the best! b.) Glossier products are overpriced. = Contradiction | MNLI-mm | Accuracy |
| Answerable | a.) How does the Dyson Airwrap work? b.) The Airwarp uses the Coanda effect to create a vortex pulling the hair towards the attachments. = Answerable | QNLI | Accuracy |
| Entail | a.) In 2006, Paul David bought a Microprocessing center to create 30,000 jobs in Northern Minnesota. b.) Paul David created 30,000 jobs in MN. = Entail | RTE | Accuracy |
| Ambiguous pronouns | a.) Federico spoke to Marie, breaking her focus. b.) Federico spoke to Marie, breaking Federico's focus. = Incorrect Referent | WNLI | Accuracy |

# Large Language Models are large and expensive!

| Optimal LLM Training Cost | | | | |
|---|---|---|---|---|
| **Model** | **Size (# Parameters)** | **Tokens** | **GPU** | **Optimal Training Compute Cost** |
| MosaicML GPT-30B | 30 Billion | 610 Billion | A100 | $ 325,855 |
| Google LaMDA | 137 Billion | 168 Billion | A100 | $ 368,846 |
| Yandex YaLM | 100 Billion | 300 Billion | A100 | $ 480,769 |
| Tsinghua University Zhipu.AI GLM | 130 Billion | 400 Billion | A100 | $ 833,333 |
| Open AI GPT-3 | 175 Billion | 300 Billion | A100 | $ 841,346 |
| AI21 Jurassic | 178 Billion | 300 Billion | A100 | $ 855,769 |
| Bloom | 176 Billion | 366 Billion | A100 | $ 1,033,756 |
| DeepMind Gopher | 280 Billion | 300 Billion | A100 | $ 1,346,154 |
| DeepMind Chinchilla | 70 Billion | 1,400 Billion | A100 | $ 1,745,014 |
| MosaicML GPT-70B | 70 Billion | 1,400 Billion | A100 | $ 1,745,014 |
| Nvidia Microsoft MT-NLG | 530 Billion | 270 Billion | A100 | $ 2,293,269 |
| Google PaLM | 540 Billion | 780 Billion | A100 | $ 6,750,000 |

*Source: semianalysis.com (calculated using Chinchilla pricing)*

# Architecture of GPT

GPT (Generative Pre-trained Transformers)

- A well-known transformer models is the gpt-3.5-turbo

- Pre-trained Transformer (GPT)

- Developed by OpenAI

**LLM**



| GPT | ChatGPT |
|---|---|
| Transformer Encoder | |
| ↓ | |
| Pre-training | Transformer Encoder |
| ↓ | ↓ |
| Transformer Decoder | Context Window |
| ↓ | ↓ |
| Attention Mechanism | Language Generation |
| ↓ | ↓ |
| Positional Encoding | Fine-tuning |
| ↓ | |
| Fine-tuning | |

# Large Language Models in Biomedicine

## Single cell

## Small molecule



nature machine intelligence
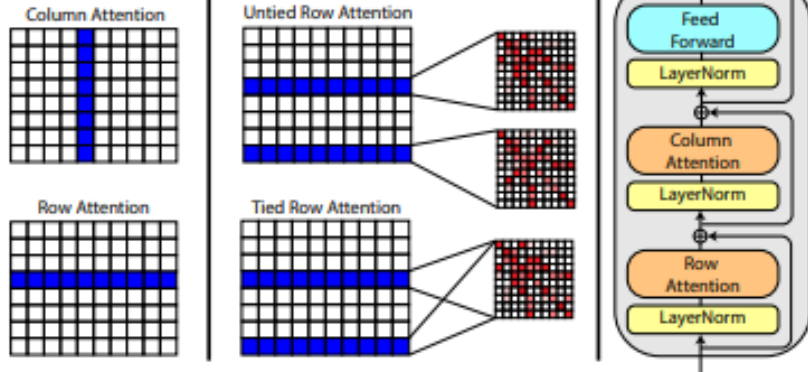
## scBERT as a large-scale pretrained deep language model for cell type annotation of single-cell RNA-seq data

## Protein

## PubMed

OXFORD

## BioGPT: generative pre-trained transformer for biomedical text generation and mining

Renqian Luo [ID], Liai Sun, Yingce Xia [ID], Tao Qin [ID], Sheng Zhang [ID], Hoifung Poon [ID] and Tie-Yan Liu

Corresponding authors: Tao Qin, Microsoft Research AI4Science, Beijing, China, E-mail: taoqin@microsoft.com; Renqian Luo, Microsoft Research AI4Science, Beijing, China, E-mail: renqianluo@microsoft.com; Yingce Xia, Microsoft Research AI4Science, Beijing, China, E-mail: yinxia@microsoft.com
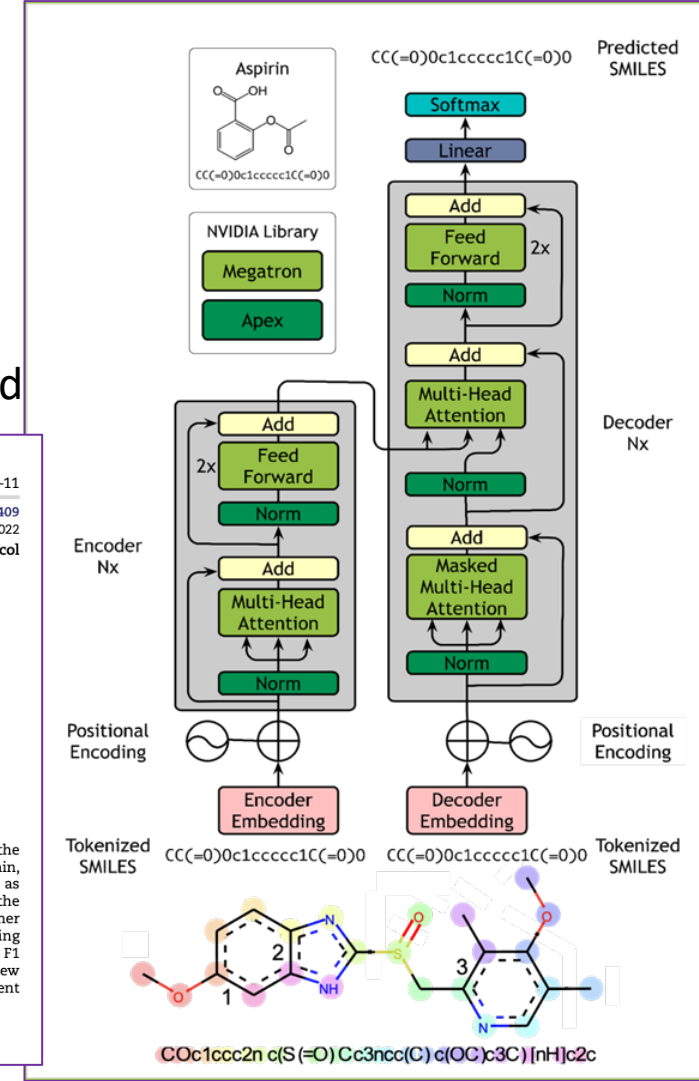
**Abstract**

Pre-trained language models have attracted increasing attention in the biomedical domain, inspired by their great success in the general natural language domain. Among the two main branches of pre-trained language models in the general language domain, i.e. BERT (and its variants) and GPT (and its variants), the first one has been extensively studied in the biomedical domain, such as BioBERT and PubMedBERT. While they have achieved great success on a variety of discriminative downstream biomedical tasks, the lack of generation ability constrains their application scope. In this paper, we propose BioGPT, a domain-specific generative Transformer language model pre-trained on large-scale biomedical literature. We evaluate BioGPT on six biomedical natural language processing tasks and demonstrate that our model outperforms previous models on most tasks. Especially, we get 44.98%, 38.42% and 40.76% F1 score on BC5CDR, KD-DTI and DDI end-to-end relation extraction tasks, respectively, and 78.2% accuracy on PubMedQA, creating a new record. Our case study on text generation further demonstrates the advantage of BioGPT on biomedical literature to generate fluent descriptions for biomedical terms.

**Keywords:** biomedical literature, generative pre-trained language model, text generation, text mining

## Language generation

## Relation extraction

## Q&A style interface

## nature biotechnology

# Efficient evolution of human antibodies from general protein language models

**Fig. 1 | Guiding evolution with protein language models. a,b,** Two possible models for relating the space of mutations with high evolutionary plausibility (for example, mutations seen in antibodies) to the space with high fitness under specific selection pressures (for example, mutations that result in high binding affinity to a specific antigen). Both models assume that mutations with high fitness make up a rare subset of the full mutational space and that, in general, high-fitness mutations are also evolutionarily plausible. Under the first model (a), mutations with high fitness are rare within the subset of mutations that are evolutionarily plausible. Under the second model (b), when restricted to the regime of plausible mutations, improvements to fitness become much more common. c, Protein language models, trained on millions of natural protein sequences learn amino acid patterns that are likely to be seen in nature. We hypothesized that most mutations with high language model likelihood would also be evolutionarily plausible. Assuming that this is true, and if the second model (b) better describes nature, then a language model with no information about specific selection pressures can still efficiently guide evolution.

15. Bepler, T. & Berger, B. Learning the protein language: evolution, structure and function. *Cell Syst.* **12**, 654–669 (2021).

Article   CAS   PubMed   PubMed Central   Google Scholar

16. Bepler, T. & Berger, B. Learning protein sequence embeddings using information from structure. *International Conference on Learning Representations*. Preprint at *arXiv* https://doi.org/10.48550/arXiv.1902.08661 (2019).

17. Hie, B., Zhong, E., Berger, B. & Bryson, B. Learning the language of viral evolution and escape. *Science* **371**, 284–288 (2021).

Article   CAS   PubMed   Google Scholar

18. Alley, E. C., Khimulya, G., Biswas, S., AlQuraishi, M. & Church, G. M. Unified rational protein engineering with sequence-based deep representation learning. *Nat. Methods* **16**, 1315–1322 (2019).

Article   CAS   PubMed   PubMed Central   Google Scholar

To select evolutionarily plausible mutations, we use
(Fig. 1c) to learn patterns that are likely to occur in r
we used general language models[19,20], trained on n
meant to represent variation across all natural prot
general evolutionary rules than could a model train
sequences[24,25,26,27] or a model directly supervised v
starting sequence, we used these language models
substitutions that we then experimentally screened
algorithm requires only a single wild-type sequence
knowledge of the antigen, task-specific supervision
structure information.

Fig. 4 | Guiding evolution without explicitly modeling fitness. a, The same strategy and language models that we use to affinity mature antibodies can also recommend high-fitness changes across a diversity of selection pressures and protein families, as identified experimentally using high-throughput scanning mutagenesis assays[8,45] (described in Supplementary Table 13). 'Fraction positive' indicates the percentage of high-fitness amino acid substitutions within either the set of substitutions recommended by the language model (LM guided) or the set of all single-residue substitutions (Background). A large portion of language-model-guided substitutions have high fitness, which, in many cases, is significantly enriched compared to the background percentage; also see Extended Data Figs. 4–6, and see Supplementary Table 13 for the exact one-sided hypergeometric P values and sample sizes. ADRB2, adrenoreceptor beta 2; β-la.,

β-lactamase; Env, envelope glycoprotein; infA, translation initiation factor 1; MAPK1, mitogen-activated protein kinase 1; PafA, phosphate-irrepressible alkaline phosphatase. b, Conceptually, the prior information encoded by evolutionary plausibility is represented in this cartoon by the rainbow road, where ascending corresponds to improving fitness and descending corresponds to lowering fitness. Moving in any direction (for example, via random or brute force mutagenesis) would most likely decrease fitness or have a high chance of being a detrimental change (represented by the green ball). However, if evolutionary plausibility is an efficient prior (Fig. 1b), then movement that is constrained to the plausible regime (for example, when guided by a language model) substantially increases the chance of improving fitness (represented by the red ball).

Trials

# The role of machine learning in clinical research: transforming the future of evidence generation

Check for updates

E. Hope Weissler[1*], Tristan Naumann[2], Tomas Andersson[3], Rajesh Ranganath[4], Olivier Elemento[5], Yuan Luo[6], Daniel F. Freitag[7], James Benoit[8], Michael C. Hughes[9], Faisal Khan[3], Paul Slater[10], Khader Shameer[3], Matthew Roe[11], Emmette Hutchison[3], Scott H. Kollins[1], Uli Broedl[12], Zhaoling Meng[13], Jennifer L. Wong[14], Lesley Curtis[1], Erich Huang[1,15] and Marzyeh Ghassemi[16,17,18,19]
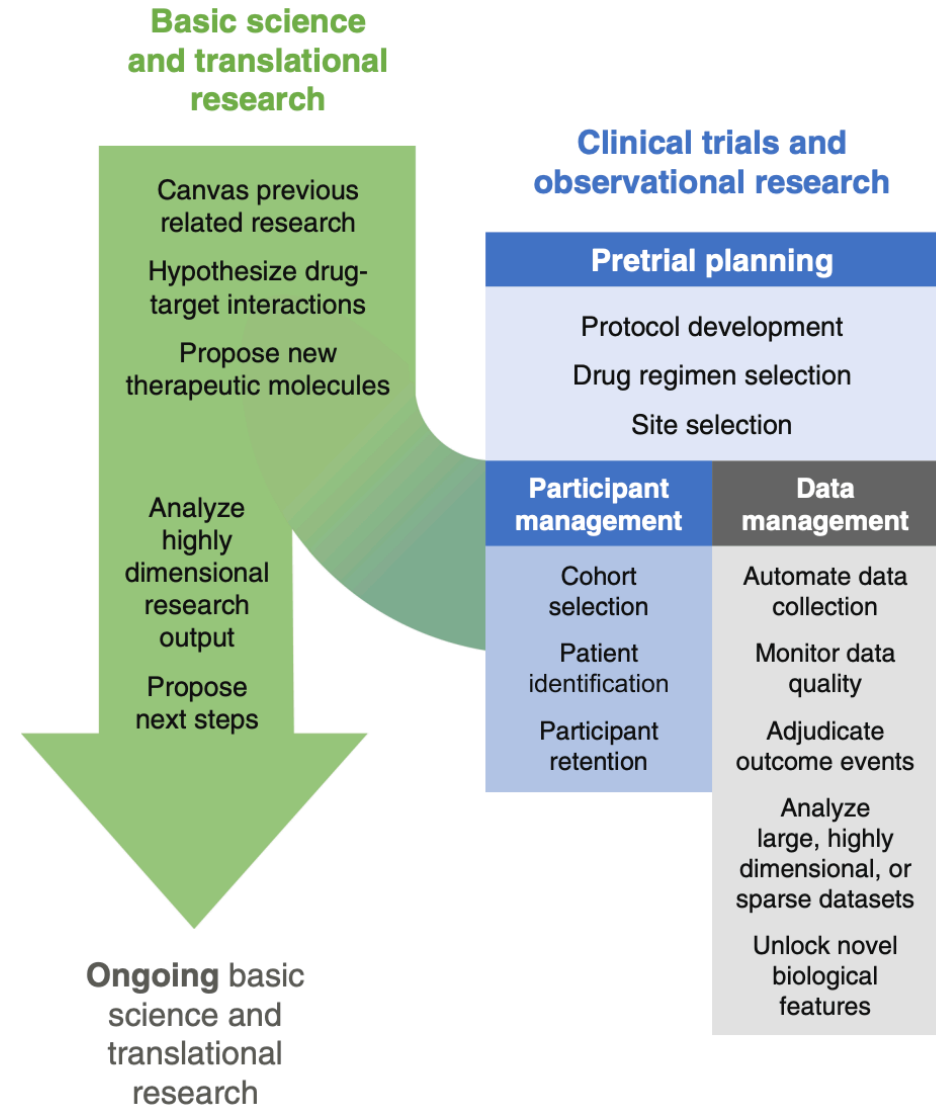
**Abstract**

**Background:** Interest in the application of machine learning (ML) to the design, conduct, and analysis of clinical trials has grown, but the evidence base for such applications has not been surveyed. This manuscript reviews the proceedings of a multi-stakeholder conference to discuss the current and future state of ML for clinical research. Key areas of clinical trial methodology in which ML holds particular promise and priority areas for further investigation are presented alongside a narrative review of evidence supporting the use of ML across the clinical trial spectrum.
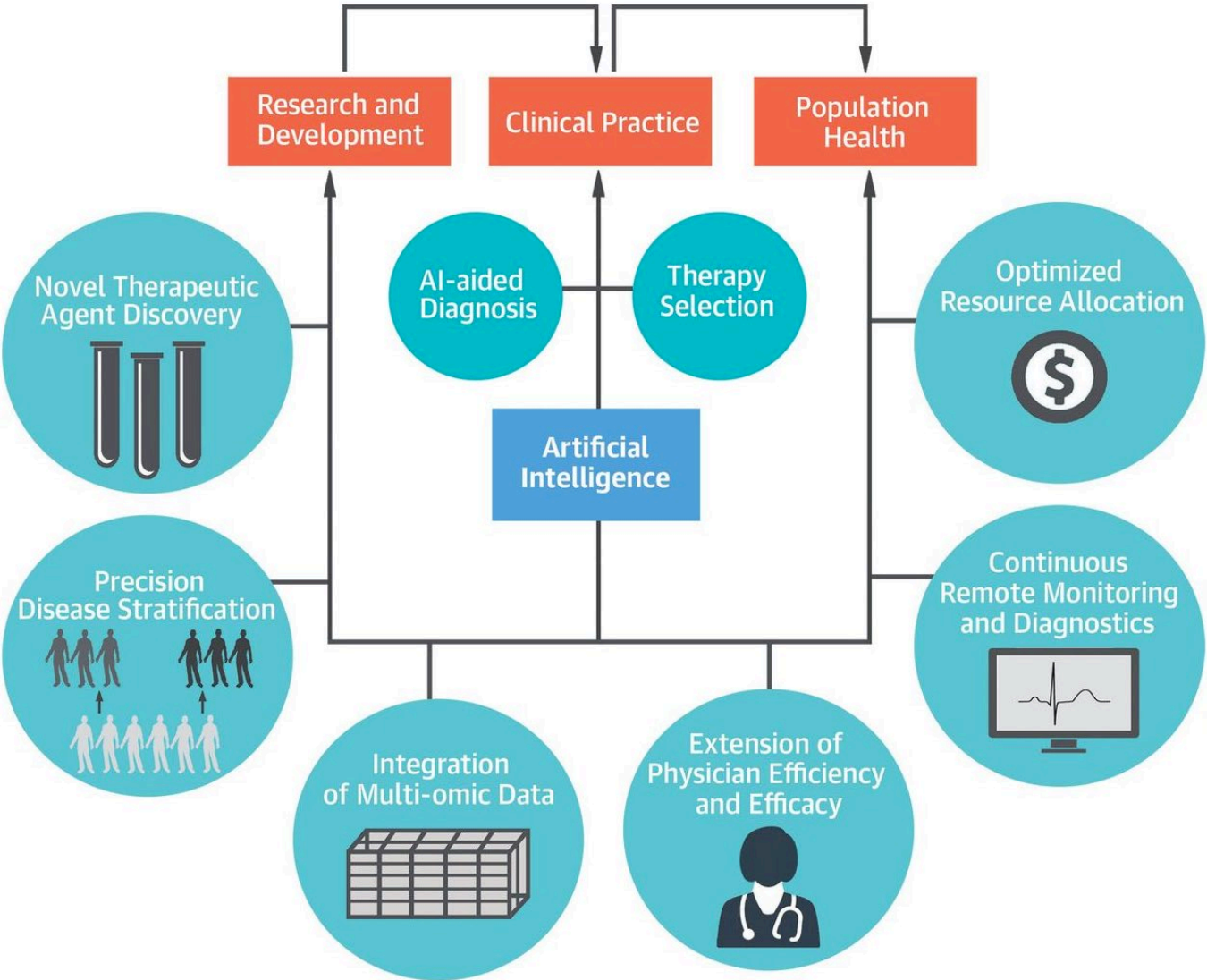
**Results:** Conference attendees included stakeholders, such as biomedical and ML researchers, representatives from the US Food and Drug Administration (FDA), artificial intelligence technology and data analytics companies, non-profit organizations, patient advocacy groups, and pharmaceutical companies. ML contributions to clinical research were highlighted in the pre-trial phase, cohort selection and participant management, and data collection and analysis. A particular focus was paid to the operational and philosophical barriers to ML in clinical research. Peer-reviewed evidence was noted to be lacking in several areas.

**Conclusions:** ML holds great promise for improving the efficiency and quality of clinical research, but substantial barriers remain, the surmounting of which will require addressing significant gaps in evidence.

**Keywords:** Clinical trials as topic; Machine learning, Artificial intelligence, Research design, Research ethics

**Basic science and translational research**

**Clinical trials and observational research**

Canvas previous related research

Hypothesize drug-target interactions

Propose new therapeutic molecules

Analyze highly dimensional research output

Propose next steps

**Ongoing** basic science and translational research

**Pretrial planning**

Protocol development

Drug regimen selection

Site selection

| **Participant management** | **Data management** |
|---|---|
| Cohort selection | Automate data collection |
| Patient identification | Monitor data quality |
| Participant retention | Adjudicate outcome events |
| | Analyze large, highly dimensional, or sparse datasets |
| | Unlock novel biological features |

# AI is deemed to play a significant role across the biomedical verticals

- Data availability is **growing** in biomedicine and healthcare

- Implementing data-driven methods that use AI-algorithms real-time variables in a hypothesis-drive/hypothesis-free approach could help us to find new targets, therapies and indications

- **Evolving** platforms including EMRs, integration engines, data mining systems and phenotyping approaches are growing

- **Integrating** novel, scalable and low-cost molecular profiling technologies with AI approaches would accelerate precision medicine development in immunology

- **Standardization** in AI, Bioinformatics and Advanced analytics would lead to develop computational medicine standards

- **Advances in AI** (including AGI) will further improve the application of AI and its impact in biomedicine and healthcare

# References

- Efficient evolution of human antibodies from general protein language models
https://pubmed.ncbi.nlm.nih.gov/37095349/


- The role of machine learning in clinical research: transforming the future of evidence generation
https://pubmed.ncbi.nlm.nih.gov/34399832/


- Sepsis in the era of data-driven medicine: personalizing risks, diagnoses, treatments and prognoses
https://pubmed.ncbi.nlm.nih.gov/31190075/


- Additional reading:
  - Shameer K, et. al; Machine learning in cardiovascular medicine: are we there yet? Heart. 2018 Jan 19. pii: heartjnl-2017-311198. doi: 10.1136/heartjnl-2017-311198. [Epub ahead of print] Review. PubMed PMID: 29352006.
  - Peters LA, et. al; Functional genomics predictive network model identifies regulators of inflammatory bowel disease. Nat Genet. 2017 Oct;49(10):1437-1449. doi: 10.1038/ng.3947. Epub 2017 Sep 11. PubMed PMID: 28892060; PubMed Central PMCID: PMC5660607.
  - Shameer K et. al; Translational bioinformatics in the era of real-time biomedical, health care and  wellness data streams. Brief Bioinform. 2017 Jan;18(1):105-124. doi: 10.1093/bib/bbv118. Epub 2016 Feb 14. PubMed PMID: 26876889; PubMed Central PMCID: PMC5221424.