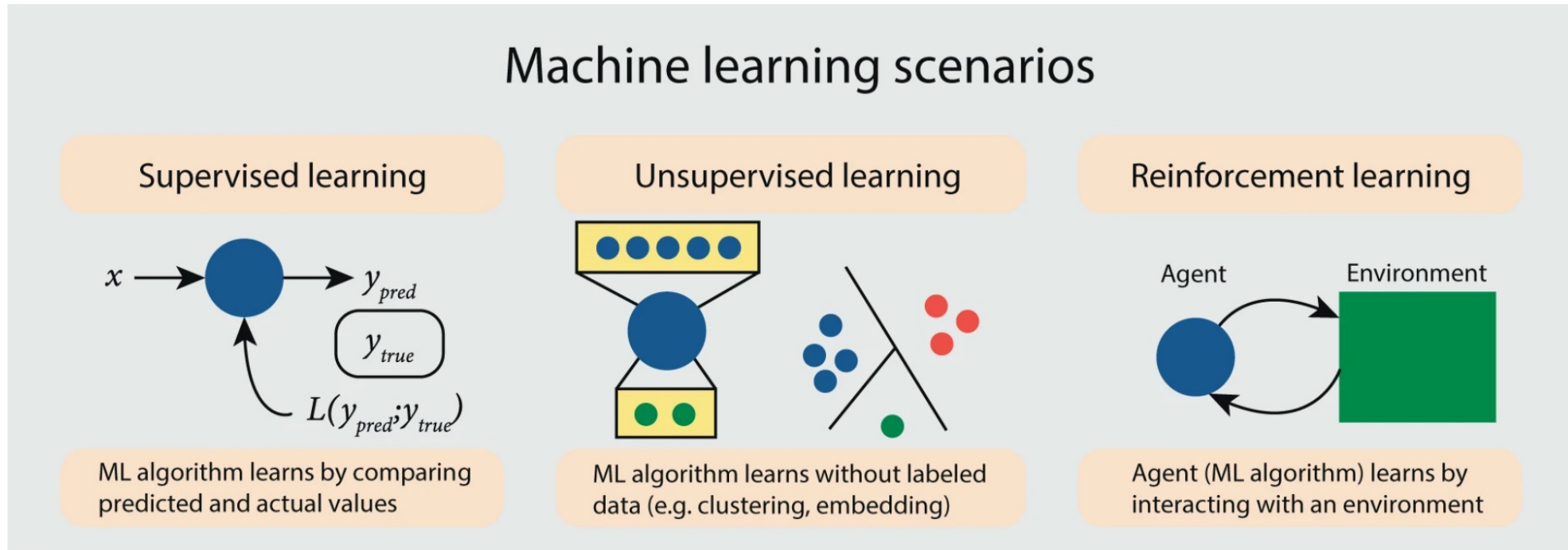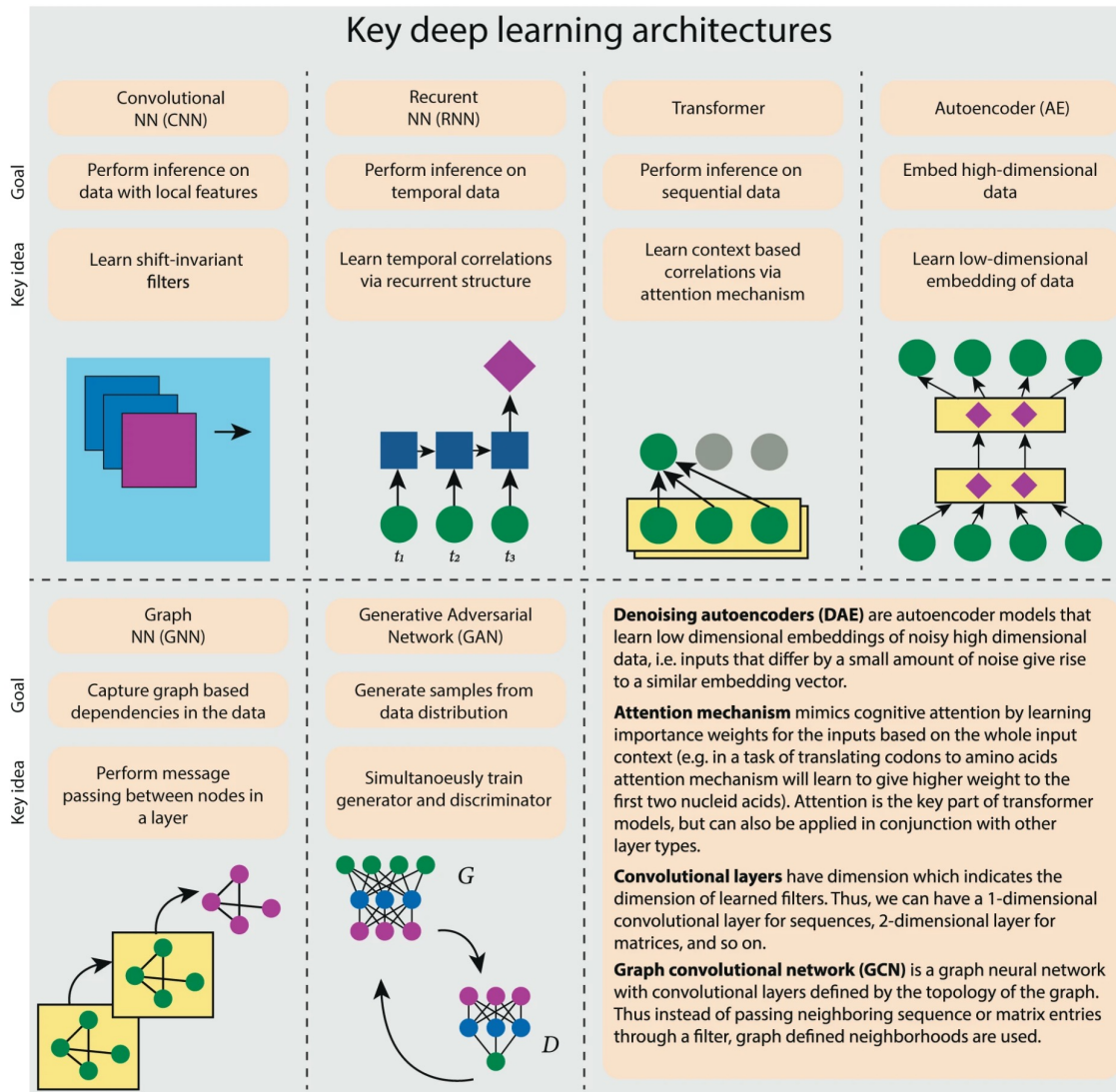# Applications for Deep Learning in Computational Biology

Yuhang Chen

TF Lecture for CBB 752

Biomedical Data Science: Mining and Modeling

Spring, 2023

Yale University

# An overview of common machine learning paradigms



Machine learning scenarios

**Supervised learning**

$x \rightarrow \bullet \rightarrow y_{pred}$

$y_{true}$

$L(y_{pred}; y_{true})$

ML algorithm learns by comparing predicted and actual values

**Unsupervised learning**

ML algorithm learns without labeled data (e.g. clustering, embedding)

**Reinforcement learning**

Agent        Environment

Agent (ML algorithm) learns by interacting with an environment

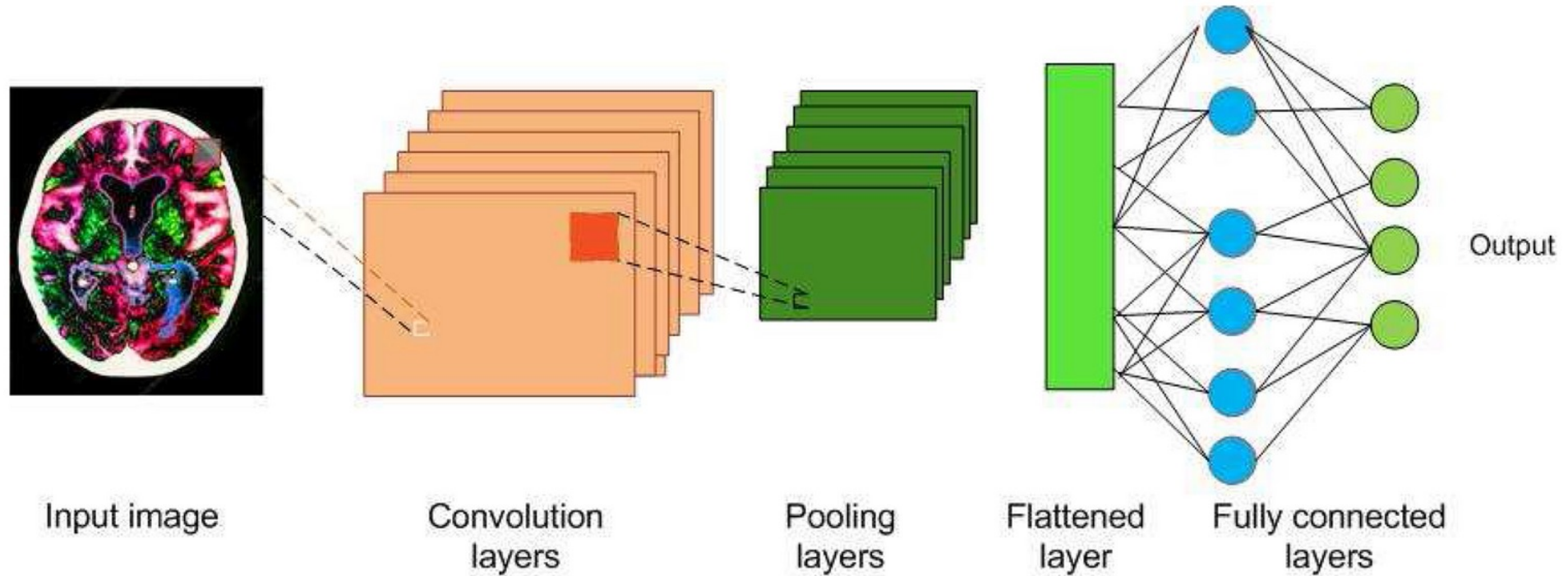# An overview of commonly used DL architecture



Key deep learning architectures

Follow up with Martin's lecture structure:

➢ Deep Supervised Learning: Deep CNN/RNN for image classification/sequence classification

➢ Deep Unsupervised Learning: Deep Autoencoder, Deep Generative Models

➢ Deep Reinforcement Learning: AlphaGO, AlphaZero

➢ Explainable AI (XAI) in biology
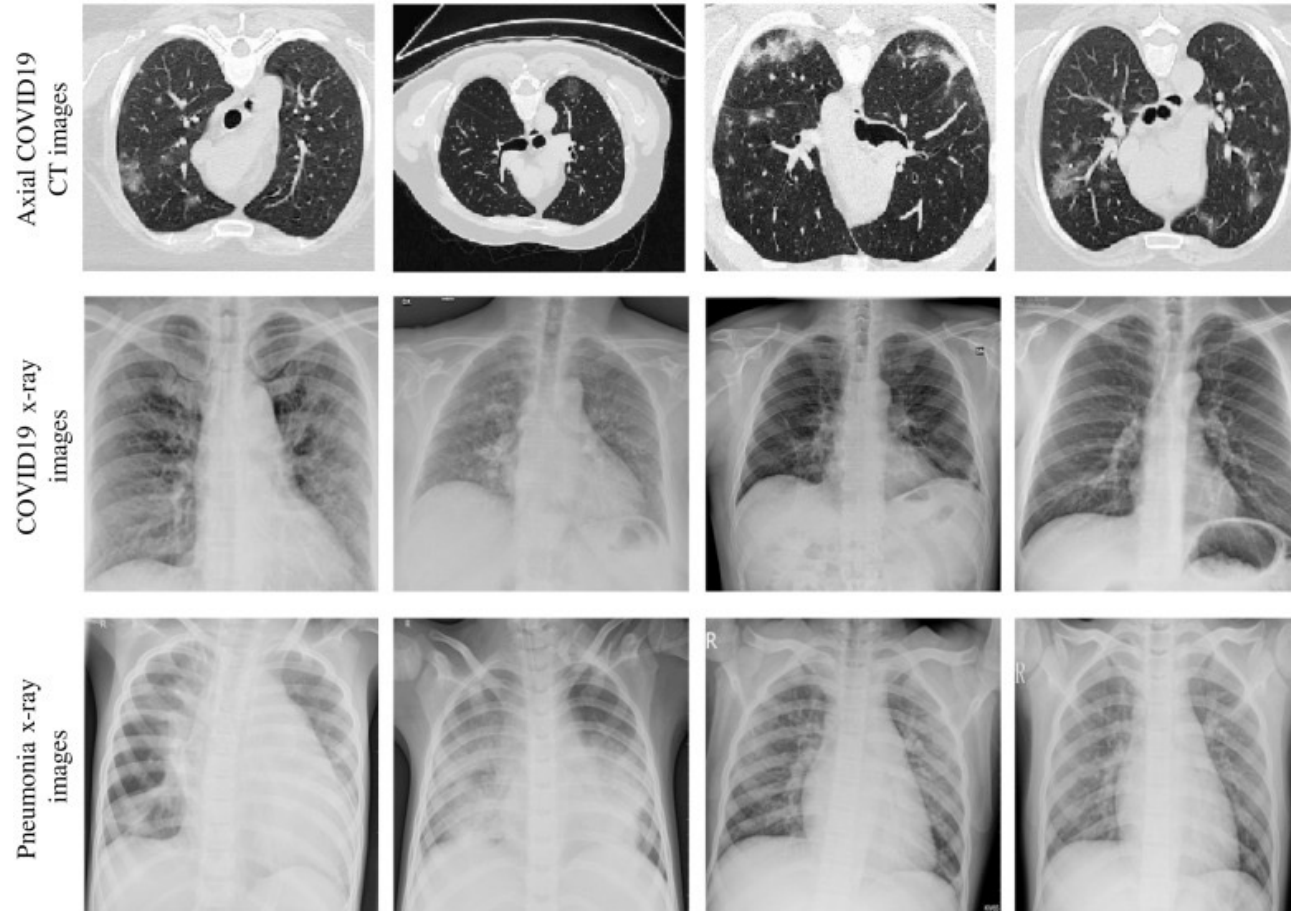
# Deep Supervised Learning

# Building blocks of a CNN



Input image | Convolution layers | Pooling layers | Flattened layer | Fully connected layers | Output

CNN is mainly used for applications in image and speech recognition.

What makes CNNs so effective is their ability to learn a sequence of filters to extract more and more complex patterns. In particular, these convolutional filters are characterized by their compact support, and by the property of being translation-invariant.

Sarvamangala, D.R., Kulkarni, R.V. Convolutional neural networks in medical image understanding: a survey. Evol. Intel. 15, 1–22 (2022).
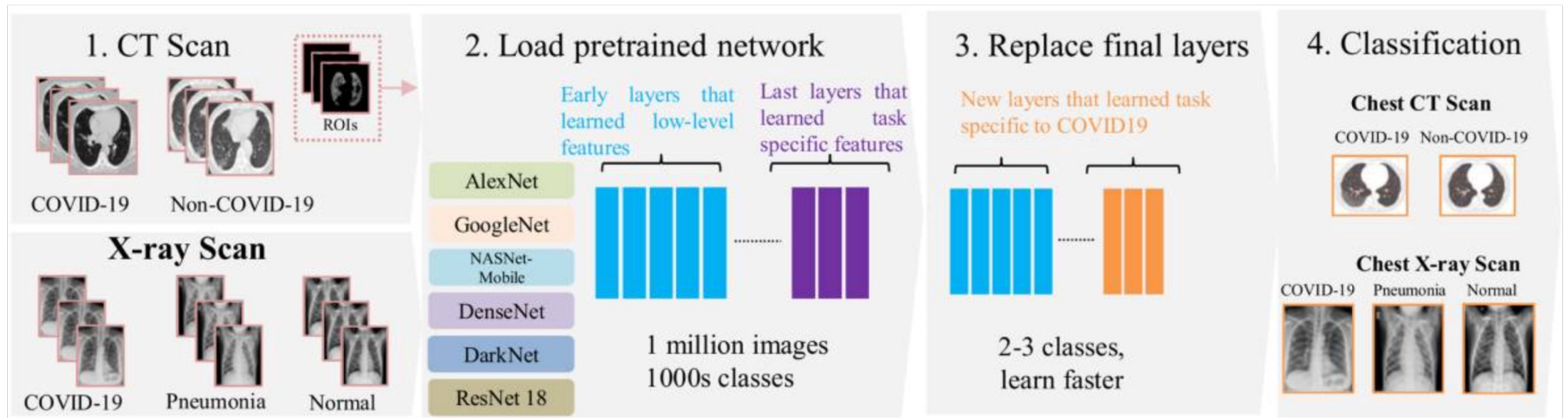
# Example: Deep CNN models for predicting COVID-19 in CT and x-ray images



Examples of COVID-19 in CT and x-ray images. First row: axial COVID-19 CT images with lesions in different positions and sizes. Second row: COVID-19 x-ray images. Third row: pneumonia x-ray images.

Chaddad A, Hassan L, Desrosiers C. Deep CNN models for predicting COVID-19 in CT and x-ray images. J Med Imaging (Bellingham). 2021 Jan;8(Suppl 1):014502.

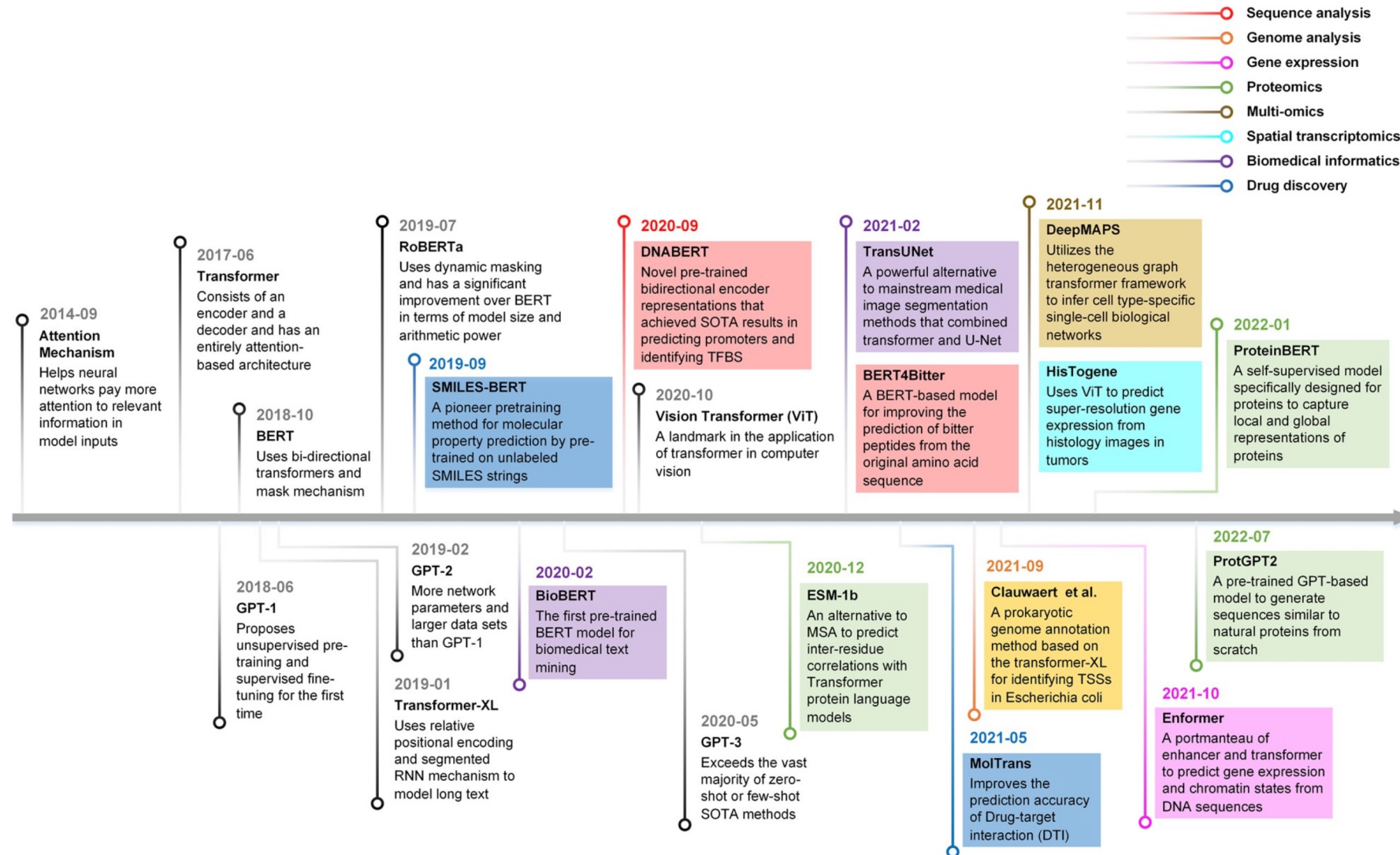# Example: Deep CNN models for predicting COVID-19 in CT and X-ray images



Regions of interest (ROI) corresponding to ground-glass opacities (GGO), consolidations, and pleural effusions were labeled in 100 axial lung CT images from 60 COVID-19-infected subjects. These segmented regions were then employed as an additional input to six deep convolutional neural network (CNN) architectures, pretrained on natural images, to differentiate between COVID-19 and normal CT images.
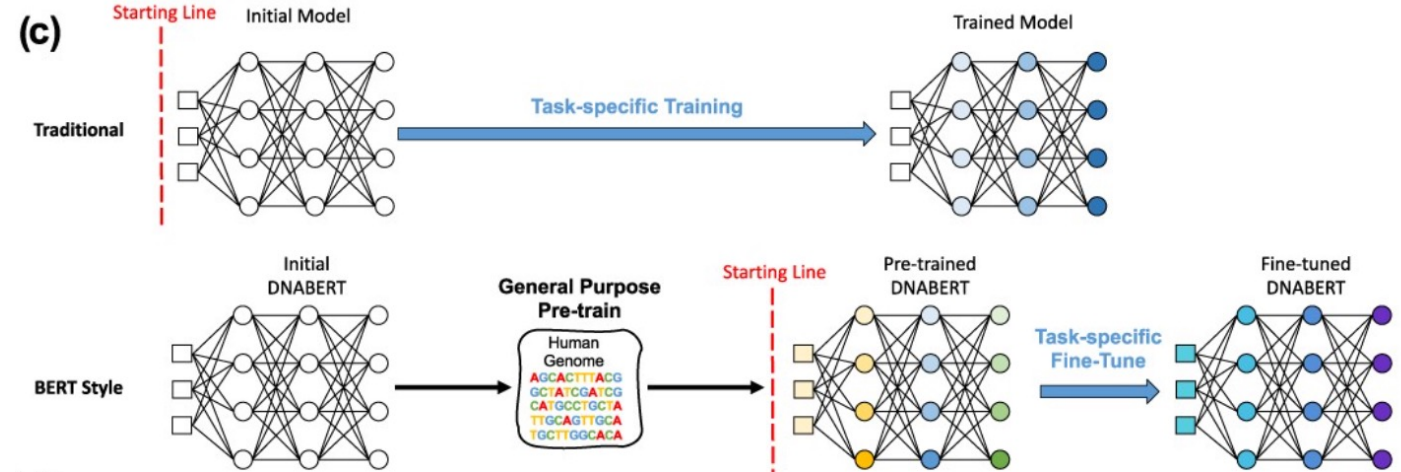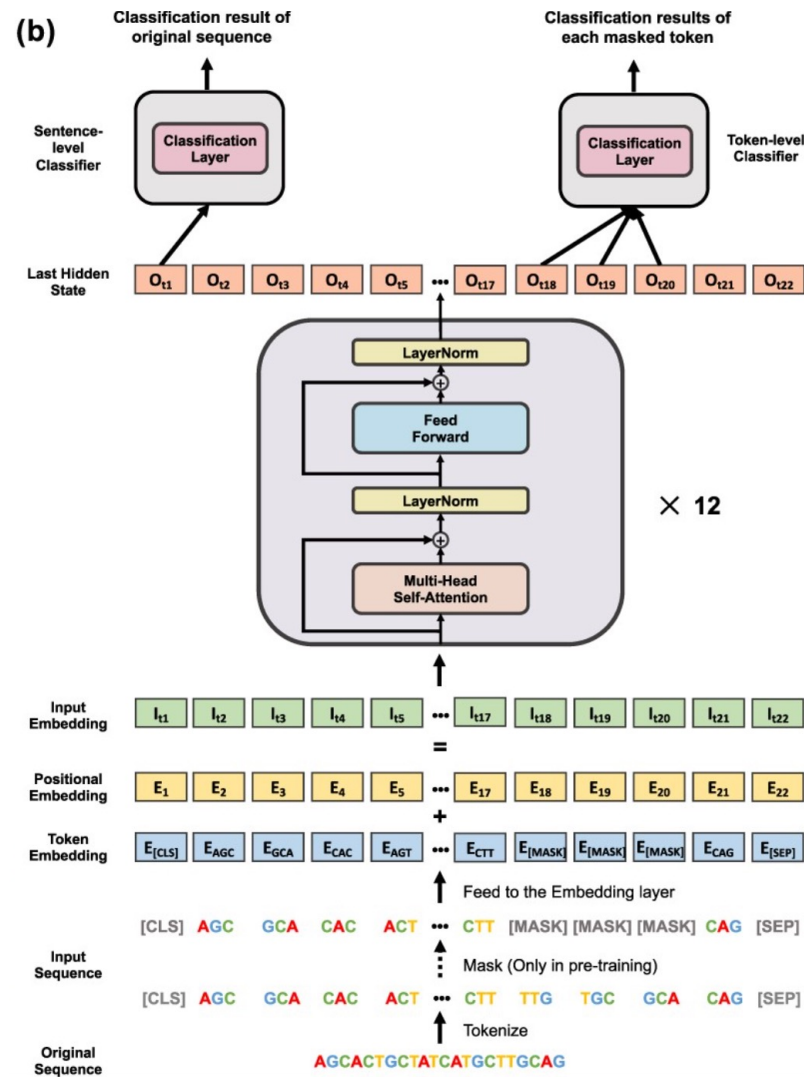Also explored the model's ability to classify x-ray images as COVID-19, non-COVID-19 pneumonia, or normal.

# An overview of important works related to TRANSFORMER in computational biology regime
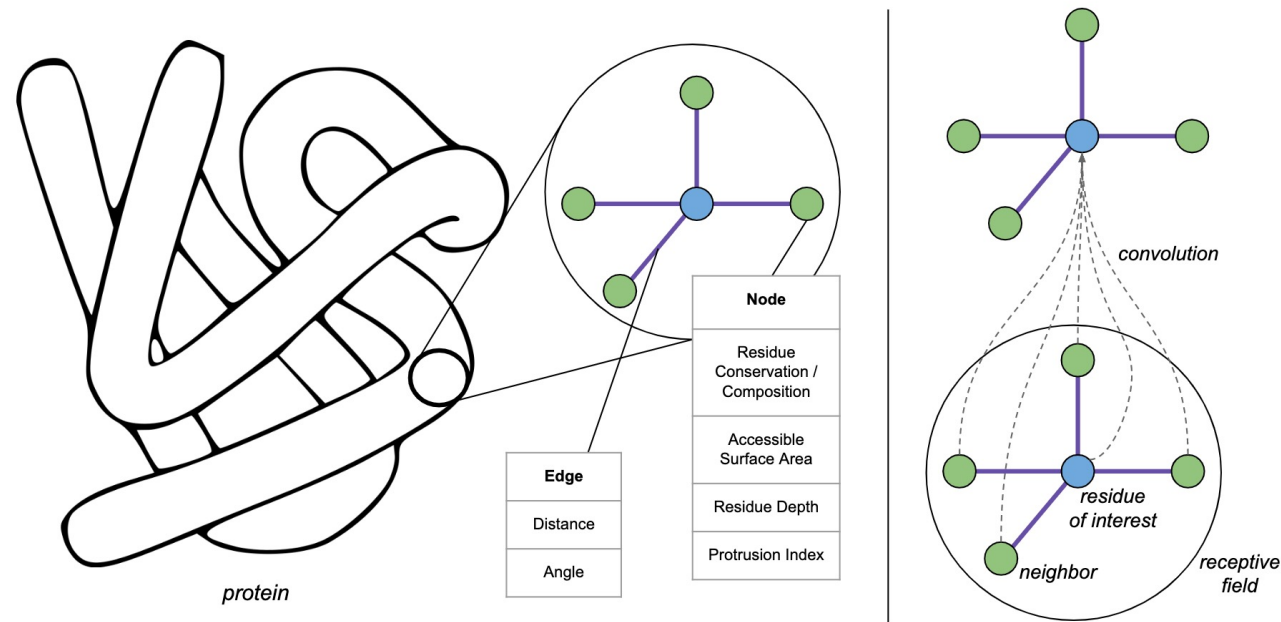
# Example: DNABERT – Transformer model for predicting promoters and identifying TFBSs



- Left: DNABERT uses tokenized k-mer sequences as input, which also contains a CLS token (a tag representing meaning of entire sentence), a SEP token (sentence separator) and MASK tokens (to represent masked k-mers in pre-training). The input passes an embedding layer and is fed to 12 Transformer blocks.

- Top:  DNABERT adopts general-purpose pre-training which can then be fine-tuned for multiple purposes using various task-specific data.

Yanrong Ji, Zhihan Zhou, Han Liu, Ramana V Davuluri, DNABERT: pre-trained Bidirectional Encoder Representations from Transformers model for DNA-language in genome, Bioinformatics

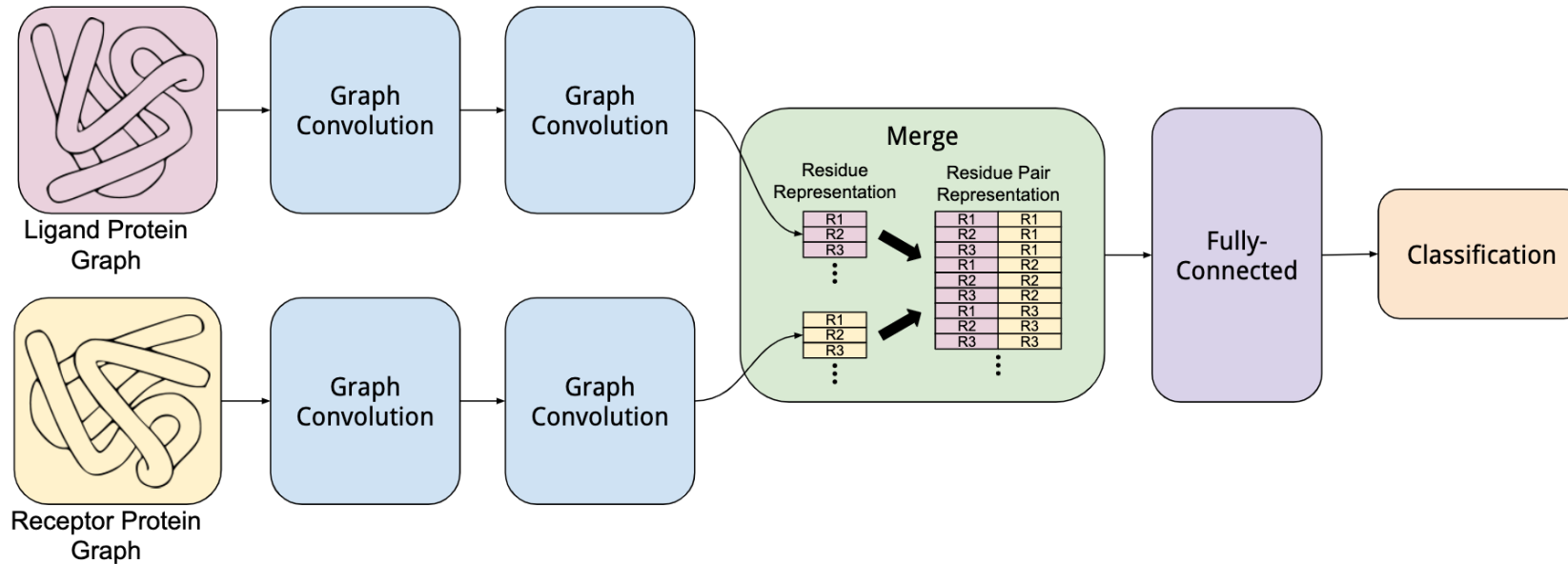# Example: Protein Interface Prediction using Graph Convolutional Networks



Graph convolution on protein structures.
Left: Each residue in a protein is a node in a graph where the neighborhood of a node is the set of neighboring nodes in the protein structure; each node has features computed from its amino acid sequence and structure, and edges have features describing the relative distance and angle between residues.
Right: Schematic description of the convolution operator which has as its receptive field a set of neighboring residues, and produces an activation which is associated with the center residue

Alex Fout, Jonathon Byrd, Basir Shariat, and Asa Ben-Hur. 2017. Protein interface prediction using graph convolutional networks. NIPS'17.

# Example: Protein Interface Prediction using Graph Convolutional Networks



An overview of the pairwise classification architecture.
Each neighborhood of a residue in the two proteins is processed using one or more graph convolution layers, with weight sharing between legs of the network. The activations generated by the convolutional layers are merged by concatenating them, followed by one or more regular dense layers.

Alex Fout, Jonathon Byrd, Basir Shariat, and Asa Ben-Hur. 2017. Protein interface prediction using graph convolutional networks. NIPS'17.

# Deep Unsupervised Learning

# RNA Velocity

Single-cell RNA-seq provides only static snapshots of cellular states at the moment of the measurement.
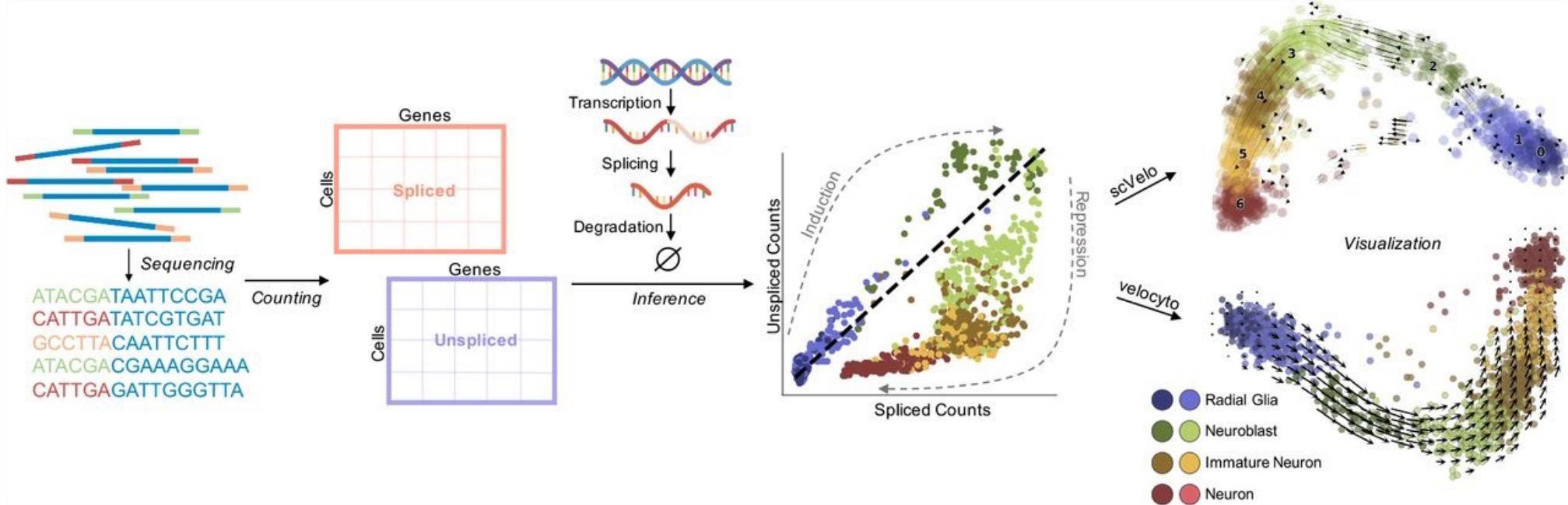
RNA velocity ([La Manno et al, 2018](#); [Bergen et al. 2020](#)) can predict the direction and speed of movement of cells in transcriptome space.

Application: analysis of cell dynamics → developmental biology, tissue regeneration, disease progression



https://www.youtube.com/watch?v=ODEP3JhyZq4

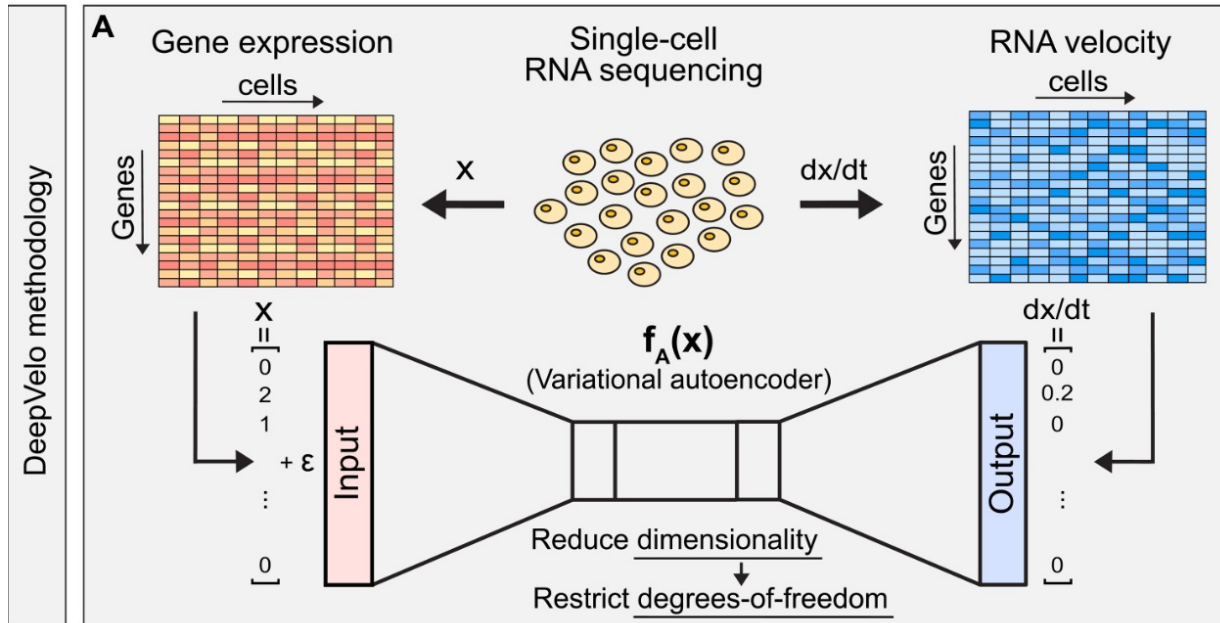# RNA velocity workflow



- Based on the relative abundance of mature (spliced) RNA and unspliced RNA to estimate the rate of RNA splicing and degradation

- Use the 2 count matrices to infer the directionality of transcription events within cells

- Phase plots describing the dynamical transcription process -> convert into embeddings showing with top 2 PCs
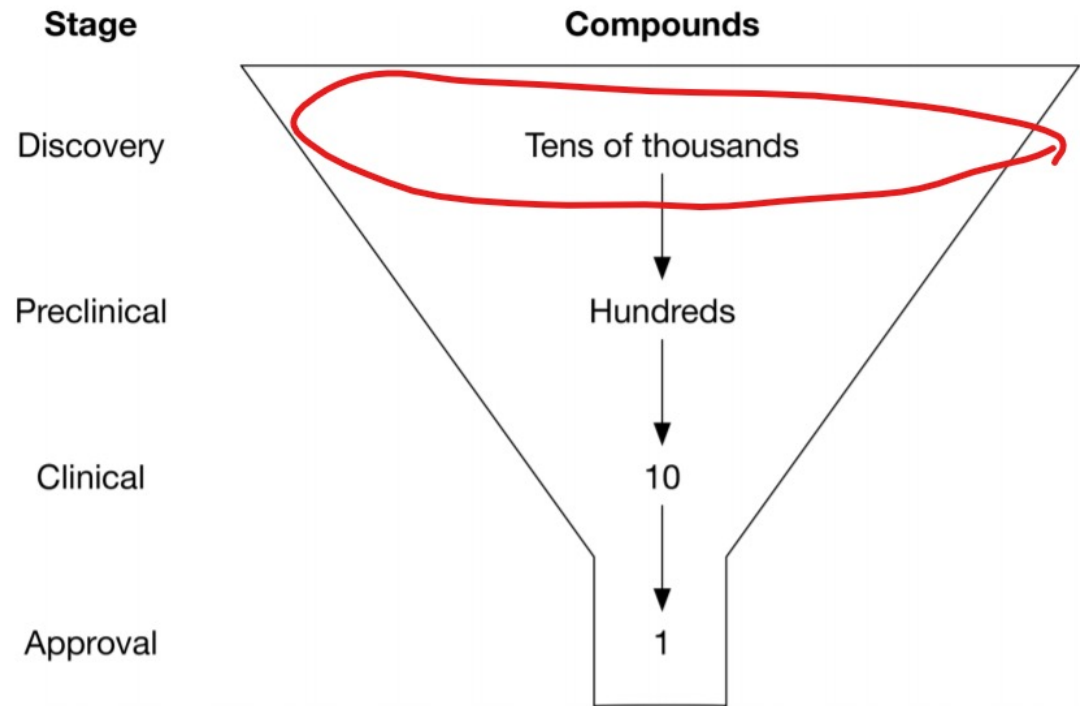
https://www.youtube.com/watch?v=ODEP3JhyZq4

# Example: DeepVelo – Model Single-cell transcriptomic velocity using VAE



- Gene expression profile of an individual cell ($x$)
- RNA velocity $\partial x/\partial t$
- Existing methods that assume linear gene interactions (i.e., $\partial x/\partial t = Ax$ with matrix A)
- Train a VAE $f_A$ to capture the nonlinear gene regulatory relationships (e.g., multiple TFs coactivating gene transcription) and map gene expression state to the RNA velocity, expressed by $\partial x/\partial t = \mathbf{f}_A(x)$

Chen Z et. al. DeepVelo: Single-cell transcriptomic deep velocity field learning with neural ordinary differential equations. Sci Adv. 2022

# GANs for Biological Research

- GANs is particularly useful for establishing potential directions in scientific study: we can generate molecules or try out potential protein structures using GANs.

- Molecules that GANs output are rarely stable or potentially useful, but we can subsequently use other deep learning models to screen the few promising molecules in a dataset.



The red circle represents the phase of drug discovery GANs will impact

# GANs architecture



The 'generator' produces a specific type of data (e.g., an image, text, or a protein sequence). The 'discriminator" tries to distinguish between the artificial data created by the 'generator' and authentic or real data.

Subsequently, the generator uses the feedback provided by the discriminator to generate new data. The generator never processes or analyzes real data and the data it produces. Therefore, its learning relies solely on the outcome of the analyses carried out by the discriminator.

# Example: ProteinGAN - A generative adversarial network that generates functional protein sequences



- Given a random input vector, the Generator network produces a protein sequence, which is scored by the Discriminator network by comparing it to the natural protein sequences. The generator tries to fool the discriminator by generating sequences that will eventually look like real ones (the generator never actually sees real enzyme sequences).

- ProteinGAN learns the evolutionary relationships of protein sequences directly from the amino-acid sequence space and creates new, highly diverse sequence variants with natural-like physical properties.

- 24% those new proteins are experimentally tested to be functional in vitro

Repecka, D., Jauniskis, V., Karpus, L. et al. Expanding functional protein sequence spaces using generative adversarial networks. Nat Mach Intell 3, 324–333 (2021).

# Explainable AI (XAI) in biology

# A tradeoff between accuracy vs. interpretability



When FPR=0.1, TPR changes from 0.58 to 0.73.

TPR (% of desats correctly predicted)

FPR (% of non-desats incorrectly predicted)

- GBM trees (AUC 0.90)
- Linear lasso (AUC 0.86)
- SaO2 linear SVM (AUC 0.76)
- Parzen window (AUC 0.70)
- Baseline (AUC 0.70)
- Random (AUC 0.5)

- Receiver operating characteristic (ROC) curve from 10-fold cross validation tests.
  - Simple vs. complex models

**Complex model f (.)**

**Generalized linear model**

**Black Box**

$X_1$
$X_2$
$\vdots$
$X_p$
?
→ Y

**X:** Features **Y:** Outcome

$X_1$ — $w_1$
$X_2$ — $w_2$
$\vdots$
$X_p$ — $w_p$
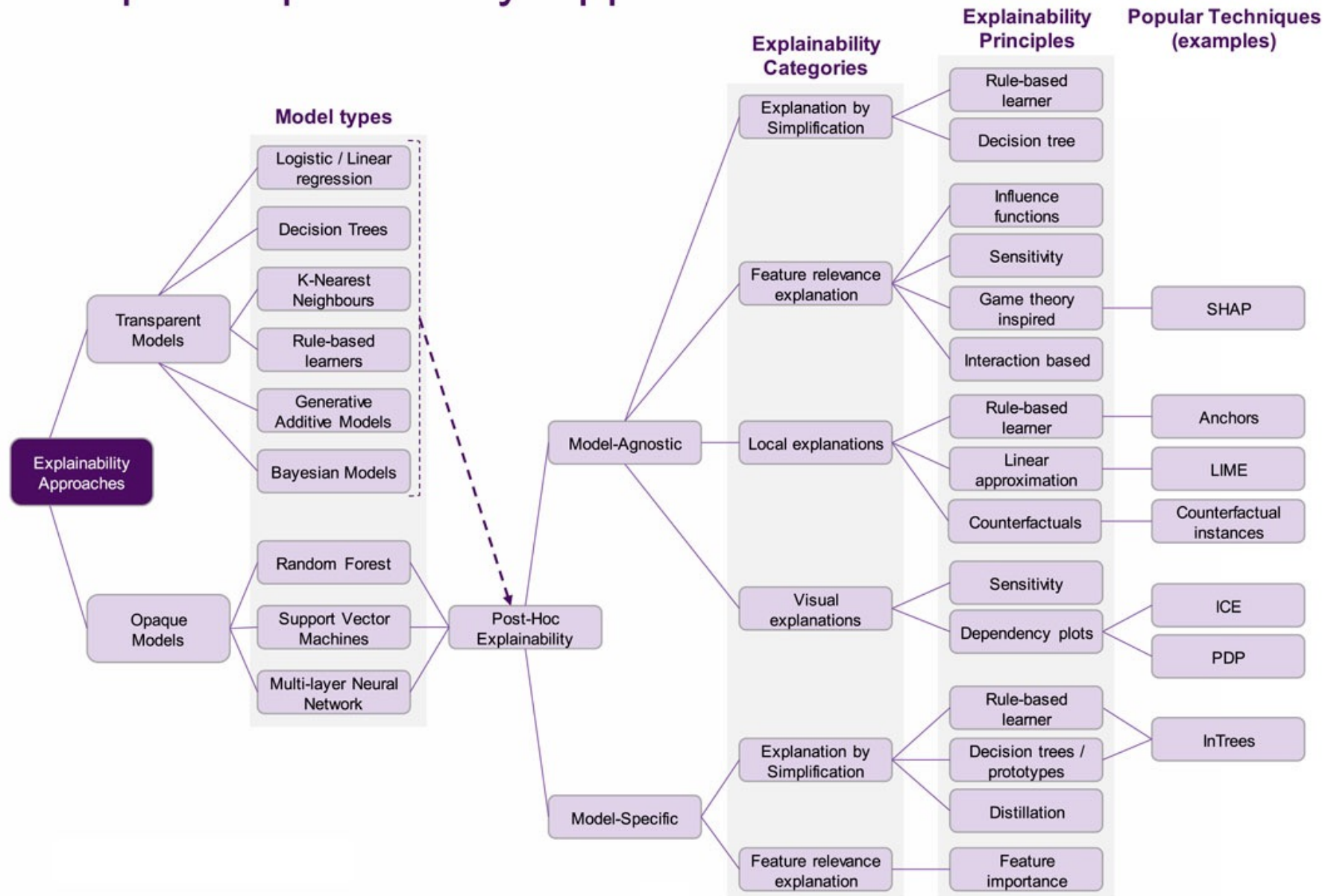$\Sigma$ → Y

# Some dimensions to evaluating explainability of a model

- **Comprehensibility:** The extent to which extracted representations are humanly comprehensible, and thus touching on the dimensions of transparency considered earlier.

- **Fidelity:** The extent to which extracted representations accurately capture the opaque models from which they were extracted.

- **Accuracy:** The ability of extracted representations to accurately predict unseen examples.

- **Scalability:** The ability of the method to scale to opaque models with large input spaces and large numbers of weighted connections.

- **Generality:** The extent to which the method requires special training regimes or restrictions on opaque models.

# Map of Explainability Approaches

# Example: An explainable artificial intelligence approach for decoding the enhancer histone modifications code and identification of novel enhancers in Drosophila



ChIP-seq data for histone modifications and STARR-seq enhancer annotations are combined and tiled into bins covering the Drosophila genome. Using these bins, traditional machine learning models (ML) and explainable AI models (XAI) can be trained to predict enhancer locations.

Wolfe, J.C., Mikheeva, L.A., Hagras, H. et al. An explainable artificial intelligence approach for decoding the enhancer histone modifications code and identification of novel enhancers in Drosophila. Genome Biol 22, 308 (2021).

# Example: An explainable artificial intelligence approach for decoding the enhancer histone modifications code and identification of novel enhancers in Drosophila



Illustration of rules identified by explainable AI model to classify regions as either enhancers or non-enhancers in Drosophila. The rules were determined to be the most effective while remaining explainable when constrained to a maximum of three epigenetic modifications per rule, and a maximum of 50 rules. These parameters were chosen to ensure that the model was explainable while maintaining a high degree of predictive power.

# Thanks for attention!
## Q&A