

Advancing bioinformatics with large language models: components, applications and perspectives

Jiajia Liu^{1,†}, Mengyuan Yang^{2,†}, Yankai Yu^{3,†}, Haixia Xu^{1,†}, Tiangang Wang¹, Kang Li⁴, Xiaobo Zhou^{1,5,6,*}

¹Center for Computational Systems Medicine, McWilliams School of Biomedical Informatics, The University of Texas Health Science Center at Houston, Houston, Texas, 77030, USA

²Department of Cell Biology and Genetics, School of Basic Medical Sciences, Xi'an Jiaotong University Health Science Center, Xi'an, China

³School of Computing and Artificial Intelligence, Southwest Jiaotong University, Chengdu, Sichuan 611756, China

⁴West China Biomedical Big Data Center, West China Hospital, Sichuan University, Chengdu, Sichuan 610041, China

⁵McGovern Medical School, The University of Texas Health Science Center at Houston, Houston, TX 77030, USA

⁶School of Dentistry, The University of Texas Health Science Center at Houston, Houston, TX 77030, USA

[†]These authors have contributed equally to this work.

* Address correspondence to:

Xiaobo Zhou, Ph.D.

McWilliams School of Biomedical Informatics

The University of Texas Health Science Center at Houston

7000 Fannin St., Houston, TX 77030

Phone: 713-500-3923

Email: Xiaobo.Zhou@uth.tmc.edu

Abstract

Large language models (LLMs) are a class of artificial intelligence models based on deep learning, which have great performance in various tasks, especially in natural language processing (NLP). Large language models typically consist of artificial neural networks with numerous parameters, trained on large amounts of unlabeled input using self-supervised or semi-supervised learning. However, their potential for solving bioinformatics problems may even exceed their proficiency in modeling human language. In this review, we will provide a comprehensive overview of the essential components of large language models (LLMs) in bioinformatics, spanning genomics, transcriptomics, proteomics, drug discovery, and single-cell analysis. Key aspects covered include tokenization methods for diverse data types, the architecture of transformer models, the core attention mechanism, and the pre-training processes underlying these models. Additionally, we will introduce currently available foundation models and highlight their downstream applications across various bioinformatics domains. Finally, drawing from our experience, we will offer practical guidance for both LLM users and developers, emphasizing strategies to optimize their use and foster further innovation in the field.

1. Introduction

Significant progress has been made in the field of natural language processing with the advent of large language models. Examples of these models include OpenAI's GPT-X [1] and Google's BERT [2] models. These models are transformative because they can understand, generate, and manipulate human language at an unprecedented scale. Vast Large language models are typically trained on datasets that encompass a significant portion of the internet's text, enabling them to learn the complexities of language and context. These models are built upon a neural network architecture called transformers [3]. The transformer architecture revolutionized NLP due to its parallelization, scalability, and ability to capture long-range dependencies in text. Instead of relying on recurrent or convolutional layers, transformers use self-attention mechanisms, as previously described, which allow them to assess the importance of every word in a sentence when understanding context. This innovation is key to their remarkable performance.

The training regimen for large language models comprises two phases: pre-training and fine-tuning. During pre-training, the model is trained on an extensive corpus of text data to acquire proficiency in grammar, factual knowledge, reasoning abilities, and word understanding. Fine-tuning tailors

these models for specific tasks like translation, summarization, or question-answering. The adaptability of large language models is a major advantage; they can excel at various NLP tasks without task-specific architectures. However, they have found applications in diverse fields beyond NLP, including biology, healthcare, education, finance, customer service, and more. In particular, there have been many successful applications of large language models in the field of bioinformatics. In this manuscript, we focus on the applications of large language models to several bioinformatic tasks through five areas: DNA level, RNA level, protein level, drug discovery and single-cell analysis. Applications of LLMs in genomics focus on LLMs using DNA sequence; applications of LLMs focus on transcriptomics using RNA sequence; applications of LLMs in proteomics focus on LLMs using protein sequence; applications of LLMs in drug discovery focus on LLMs using Molecular SMILES (seq) and applications of LLMs in single-cell analysis focus on LLMs using scRNA-seq, scMulti-omics and spatial transcriptomics data (**Figure 1**).

2. Understanding the Building Blocks of Large Language Models in Bioinformatics

Building large language models involves several critical components, including tokenization methods, embedding techniques, attention mechanisms, transformer architectures, and the training processes for large-scale models. Each of these elements plays a vital role in enabling the models to process, understand, and generate complex data.

2.1 Tokenization and input embedding

Tokenization methods are essential for processing raw input data, breaking it down into smaller, manageable units (tokens) that can be analyzed and processed by models. The choice of tokenization method varies depending on the type of data being handled (**Figure 2a, Table 1**).

In DNA and RNA sequence data, tokenization converts raw nucleotide sequences (A, T, C, G for DNA or A, U, C, G for RNA) into a numerical format suitable for computational models. A common method is one-hot encoding, where each nucleotide is represented as a binary vector with a ‘1’ indicating its position (e.g., [1, 0, 0, 0] for A in DNA), as used in RNA-FM [4] and RNA-MSM [5]. Another widely adopted approach is k-mer tokenization, which segments sequences into overlapping substrings of fixed length ‘k’ (e.g., for k=3, “ATGC” becomes “ATG” and “TGC”). This method is employed in models like DNABERT[6], DNAGPT [7], and RNABERT [8].

Additionally, specialized tokens such as ‘[IND]’ can be introduced to mark the start or end of sequences or to handle unknown characters or gaps, as demonstrated in RNAErnie [9].

In protein language models, the input data primarily includes multiple sequence alignments (MSAs), protein sequences, biomedical/biological text, and cDNA. The basic units of MSAs and protein sequences are amino acids, leading most protein language models to use Single Amino Acid Tokenization, where protein sequences are segmented into individual amino acids. This approach is akin to the k-mers method used for DNA and RNA sequences and is employed in models such as ESM-1b [10], ProtTrans [11], and ProGen [12]. For biomedical and biological text, including general descriptions, conditioning tags in generative models, and resources like Gene Ontology (GO), tokenization methods from natural language processing (NLP) are widely used. Methods like WordPiece Tokenization build vocabulary using frequency-based greedy algorithms and segment text into discrete tokens, as demonstrated in ProtST [13]. For cDNA data, tokenization is similar to that of protein sequences but differs in the basic unit. Instead of amino acids, sequences are tokenized into codons, or triplets of nucleotides, as seen in CaLM [14].

In drug discovery, small molecule drugs account for 98% of commonly used medications [1]. LLMs leverage four main tokenization methods to uncover molecular patterns and drug-target interactions. Atom-level tokenization treats molecules as sequences of individual atoms, analogous to character-level text representation, as seen in K-BERT [15]. MolGPT [16] utilizes a SMILES tokenizer that segments molecular structures into units such as atoms, bond types, and ring markers. A Graph-based VQ-VAE approach enhances this by encoding atoms into context-aware discrete values, distinguishing roles like aldehyde versus ester carbons, based on latent codes derived from a graph-based Vector Quantized Variational Autoencoder (VQ-VAE). This method categorizes atoms into chemically meaningful sub-classes, enriching the molecular vocabulary. Fingerprint tokens, another method, represent molecules through binary or numerical vectors summarizing molecular properties or structural patterns, as seen in SMILES-BERT [17].

Tokenization methods for single-cell profiles include four main strategies. Gene ranking/reindexing-based methods rank genes by expression levels and create tokens using ranked gene symbols or unique integer identifiers, as seen in Geneformer [18] and tGPT [19]. Binning-based methods divide gene expression into predefined intervals, assigning tokens based on the corresponding bin, used in models like scBERT [20] and scGPT [21]. Gene set or pathway-based methods group genes into biologically meaningful sets, such as pathways or Gene Ontology terms,

with tokens representing the activation of these sets, exemplified by TOSICA [22]. Patch-based methods segment gene expression vectors into equal-sized sub-vectors, as seen in CIForm [152]. Alternatively, convolutional neural networks (CNNs) can be used to transform the reshaped gene expression matrix into several flattened 2D patches, as demonstrated by scTranSort [23]. Another variation involves reshaping the sub-vectors into a gene expression matrix after segmentation, as employed in scCLIP [24]. In addition to the four methods mentioned above, a more direct approach involves projecting gene expression directly, as seen in models like scFoundation [25], and scMulan [26]. Alternatively, some methods tokenize cells instead of genes, as exemplified by models such as CellPLM [27], ScRAT [28], and mBERT [29], which utilize cell tokens during model training (**Table 1**). These strategies allow models to capture biological structure and variability, tailoring tokenization to single-cell data characteristics.

After tokenization, embedding converts tokens into continuous vector representations, capturing the semantic relationships between them. Positional encoding represents the token order by adding vectors that encode the relative or absolute positions of tokens in the sequence. The final step involves combining the token embeddings with the positional embeddings to create a unified input embedding, which is then fed into the model for further processing (**Figure 2b**).

2.2 Architecture of transformer models

Transformers are the foundational architecture in large language models (LLMs) and consist of two main components: the encoder and the decoder. The encoder takes the input data and processes it in parallel across multiple layers to capture relationships within the sequence. The decoder, on the other hand, generates output sequences based on the encoder's processed information, typically used in tasks like translation or text generation. Each component is built on layers of multi-head attention, add and norm layer, and feed-forward layer (**Figure 2c**).

Attention Mechanism: A key innovation of the transformer is the attention mechanism, particularly self-attention [3], which allows the model to weigh the importance of different tokens in a sequence relative to each other. In self-attention, each token computes a score based on how much attention it should pay to other tokens in the sequence. This is done by calculating three key components: Query (Q), Key (K), and Value (V) vectors for each token (**Figure 2d**). The attention score is computed as the dot product between the Query of one token and the Key of another token, followed by a softmax operation to normalize the scores. These scores are then used to weight the

Value vectors, which are aggregated to form the output representation for each token as following [3]:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (1)$$

Multi-head attention extends this idea by running multiple attention mechanisms (or “heads”) in parallel. Each attention head processes the input tokens in a slightly different way by using different sets of learned weights for the Q, K, and V vectors. The results of all heads are concatenated and linearly transformed, allowing the model to capture different aspects of relationships between tokens simultaneously. This mechanism enables the model to focus on various parts of the input sequence at once, learning different types of interactions between tokens. For example, in single-cell foundation models, self-attention can help identify important gene interactions by determining which genes (tokens) should focus on each other during processing. In this way, multi-head attention allows the model to capture complex relationships between genes in single-cell RNA-seq data, where multiple aspects of gene expression (such as co-expression patterns or functional relationships) need to be captured simultaneously.

Add and norm layer: The add and norm layer performs layer normalization and residual connections, which help stabilize training by ensuring that the output from each layer is added to the input before being normalized. This allows for smoother gradient flow and avoids the vanishing gradient problem.

Feed-forward layer: After the attention mechanism, the feed-forward network is a fully connected neural network (**Figure 2e**), helping the model learn complex mappings and capture more abstract representations of the input data.

2.3 BERT and GPT models

BERT and GPT stand as two exceptional language models. Both BERT and GPT leverage the transformer architecture, employing attention mechanisms to grasp dependencies within input data.

BERT (Bidirectional Encoder Representations from Transformers). BERT is basically an encoder stack of transformer architecture, which was introduced by Google in 2018 [2]. BERT is trained using a bidirectional approach, meaning it considers context from both the left and right of each token during training. This enables BERT to capture richer, more context-aware representations. BERT is typically pre-trained using a masked language model (MLM) task, where random tokens in a sequence are masked, and the model is tasked with predicting them. This bidirectional training allows BERT to better understand the full context of a sequence or biological

sentence, such as a cell (**Figure 2f**). For example, scBERT, a single-cell adaptation of BERT, applies this approach to single-cell RNA-seq data. By masking random gene tokens and predicting them during pretraining, scBERT learns complex dependencies and co-expression patterns between genes. This enables it to capture the full transcriptional context of individual cells, improving downstream tasks like cell type classification.

GPT (Generative Pretrained Transformer). Introduced by OpenAI [1], GPT is based on a decoder stack of transformer architecture. Unlike BERT, GPT uses a unidirectional training approach, processing the input sequence from left to right (**Figure 2f**). It is trained using autoregressive learning, where each token is predicted based on the previous ones, making it particularly suited for generational tasks. GPT excels in zero-shot learning, where they perform tasks without needing task-specific training data. For example, DNAGPT leverages its pretrained knowledge to perform tasks like predicting DNA motifs or identifying regulatory elements without explicit task-specific training. When prompted with a sequence such as “Find the transcription factor binding motif in the following DNA sequence: AGCTTAGGCC...”, DNAGPT can identify or generate plausible motifs based on its understanding of DNA patterns learned during pretraining.

3. Foundation models in bioinformatics

3.1 Key components of biological foundation models

Foundation models are a category of large-scale, pre-trained models designed to be versatile and adaptable to various downstream tasks. They are built upon several fundamental components that enable their widespread applicability and effectiveness across domains. First, foundation models are trained on extensive and diverse datasets to capture broad, generalizable patterns. In single-cell biology, for example, datasets with millions of cells spanning multiple tissues and conditions are often used. Second, the architecture of foundation models is typically designed for flexibility and scalability. Their architecture, often transformer-based (e.g., GPT and BERT), are specifically designed for flexibility and scalability. Third, Self-supervised learning is a core training strategy for foundation models. By creating tasks such as masked prediction, contrastive learning, or next-token prediction, models can learn representations without requiring labeled data. Fourth, foundation models exhibit multi-task transferability, leveraging a two-step process of pre-training and fine-tuning (**Figure 3**). During pre-training, these models are trained on large-scale datasets to develop robust generalization capabilities by capturing broad patterns and knowledge. Fine-

tuning involves adapting the pre-trained model to specific tasks by exposing it to unique data and additional training. This approach enables foundation models to adjust effectively to diverse applications while maintaining their versatility across a wide range of domains. Last but not least, training foundation models requires significant computational resources, often involving GPU or TPU clusters. Foundation models typically feature billions or even trillions of parameters.

3.2 Foundation models in different biological domains

DNA foundation models. Currently, DNA sequence-based foundation models are powerful tools that leverage advanced deep learning architectures to analyze and interpret genomic data [30]. These models are built on frameworks like BERT and GPT, which have been adapted for the specific challenges of genomic sequences (**Table 2**). For example, DNABERT [6] is a BERT-based model trained on the human reference genome, enabling it to capture the contextual relationships between nucleotides and perform tasks such as sequence classification and variant prediction. Expanding beyond a single species, Nucleotide Transformer [31] and Genomic Pre-trained Network (GPN) [32] are transformer-based models that incorporate the multiple species reference genomes, providing a broader understanding of genomic diversity. DNABERT-2 [33] takes this a step further by training on multi-species genomic data from 135 species, allowing for cross-species genomic analysis. Similarly, GROVER [34], another BERT-based model, is focused on the human reference genome and is designed for applications such as understanding gene expression and functional genomics. On the other hand, DNAGPT, based on the GPT architecture, is trained not only on the human reference genome but also on reference genomes from nine other species, facilitating tasks such as sequence generation and evolutionary analysis. Together, these DNA sequence-based foundation models represent a leap forward in computational genomics, enabling more accurate predictions, better understanding of genetic variation, and advancements in personalized medicine.

RNA foundation models. RNA sequence-based language models, particularly BERT-based and Transformer-based models, have gained significant traction in the analysis of RNA sequences due to their ability to understand the complex patterns and structures of RNA. These models are trained using a wide variety of RNA types, including non-coding RNAs (ncRNAs), coding RNA, and untranslated regions (UTRs), across diverse organisms (**Table 2**). For instance, RNABERT [8], RNA-FM [4], RNA-MSM [5], and UNI-RNA [35] focus on all ncRNA types from a broad range of species, enabling insights into RNA function and interactions. Models like SpliceBERT [36]

specialize in coding RNA sequences from 72 vertebrates, while 3UTRBERT [37] is specifically designed for human mRNA transcripts, particularly the 3' untranslated regions. Additionally, UTR-LM [38] focuses on 5' UTR sequences from five species, and RNAErnie [9], a Transformer-based model, covers a wide range of ncRNAs. These models are part of a rapidly growing field aimed at advancing RNA sequence analysis, facilitating the study of RNA biology and its role in various biological processes and diseases. Through the use of these RNA-based language models, researchers can make significant strides in understanding RNA structure, function, and regulatory mechanisms.

Protein foundation models. Foundation models for proteins can be directly utilized to obtain high-quality protein embeddings and support various downstream applications. The foundational protein models listed in **Table 2** not only fulfill these requirements but also exhibit unique characteristics. For example, TAPE [39] made a significant contribution by introducing a comprehensive benchmark for protein bioinformatics tasks. ESM-1b [40] applied the transformer architecture of large language models in a highly standardized manner to protein representation learning. This model has since been widely used to generate protein sequence embeddings, and its variants can also be found via the same link provided in **Table 2**. ProtTrans [11], compared to ESM-1b, significantly expanded the model architecture, the number of parameters, and the size of the training dataset. It has been widely adopted as a frozen encoder for protein sequences. ProtGPT2 [41], as its name suggests, extends GPT-2 into the protein domain (with links providing details on the GPT-2 training framework). Recent foundation models like ProtBert [42] and KeAP [43] integrate biomedical text information alongside protein sequences. Notably, KeAP incorporates a knowledge graph to enhance this integration. Both models demonstrate that multimodal fusion within proteomics often produces more expressive features. CaLM [14], on the other hand, represents proteins using cDNA, embedding cross-omics biological information. From the perspective of algorithmic advancements, the integration of multimodal information within a single omics domain, as well as cross-omics data fusion, represents key strategies for constructing unified large-scale biological models.

Drug discovery foundation models. It has been postulated that the total number of potential drug-like candidates range from 10^{23} to 10^{60} molecules[44]. Foundation models leverage diverse tokenization strategies, embedding techniques, and pre-training mechanisms to enhance molecular representation learning, facilitating the optimization of various downstream tasks (**Table 2**). For

instance, Mol-BERT [45] employs a context-aware tokenizer to encode atoms into chemically meaningful discrete values, although this approach results in an unbalanced atom vocabulary. SMILES-BERT [46], a semi-supervised model incorporating an attention-based Transformer architecture, utilizes datasets such as LogP, PM2, and PCBA-686978 to pre-train the model via a Masked SMILES Recovery (MSR) task. This model demonstrates strong generalization capabilities, enabling its application to diverse molecular property prediction tasks through fine-tuning. Similarly, Mol-GPT [47] facilitates the generation of molecules with specific scaffolds and desired molecular properties by conditioning the generation process on scaffold SMILES strings and property values. Notably, SynerGPT [48] enables a pre-trained GPT model to perform in-context learning of “drug synergy functions”, showcasing potential for future advancements in personalized drug discovery. These foundation models developed based on distinct strategies, effectively learn representations from raw sequence data and molecular descriptors. They provide significant insights into the design of small-molecule drugs, drug-drug interactions, and drug-target interactions.

Single-cell foundation models. Foundation models in single-cell analysis are revolutionizing the field by offering scalable and versatile solutions for a wide range of tasks, leveraging both cell and gene-level representations (**Table 2**). Models like scBERT [20], tGPT [19], scMulan [26], UCE [49] and CancerFoundation [50] focus on learning robust cell representations, effectively supporting applications such as cell clustering, cell type annotation, batch effect correction, trajectory inference and drug response prediction. These models excel at analyzing heterogeneous cellular populations and uncovering cellular dynamics. In contrast, models like scGPT [21], scFoundation [25], Geneformer [18] GeneCompass [51] and scPRINT [52] combine the ability to learn both cell and gene-level representations. They capture inter-gene relationships and regulatory networks, making them highly effective for tasks such as gene expression profiling, gene regulatory network (GRN) inference, gene perturbation prediction, and drug dose-response prediction. Notably, scGPT can also handle single-cell multi-omics data, facilitating tasks like scRNA-seq and scATAC-seq integration. Another notable model is Nichefomer [53], a foundation model specifically designed for spatial transcriptomics. It focuses on learning cell representations while being highly adaptable to various downstream tasks in spatial transcriptomics, such as spatial label prediction (e.g., cell type, niche, and region labels), niche composition analysis, and neighborhood density prediction. Additionally, Nichefomer can

generate joint embeddings of scRNA-seq and spatial transcriptomics data, facilitating the integration of these modalities for a more comprehensive understanding of cellular and spatial interactions.

4. Applications of large language models in bioinformatics

Large language models (LLMs) have seen numerous successful applications in bioinformatics, addressing a wide array of tasks across DNA, RNA, protein, drug discovery, and single-cell analysis (**Figure 4**). These applications highlight the adaptability and potential of LLMs in overcoming bioinformatic challenges, enabling deeper insights into complex biological systems and fostering advancements across multiple domains.

4.1 Applications of large language models in genomics

The DNA language models take DNA sequence as input, use transformer, BERT, GPT models to solve multiple biological tasks, including genome-wide variant effects prediction, DNA cis-regulatory regions prediction, DNA-protein interaction prediction, DNA methylation (6mA,4mC 5hmC) prediction, splice sites prediction from DNA sequence (**Table 3, Supplementary Figure 1**). A detailed list of DNA language models, their downstream tasks, and the datasets used can be found in **Supplementary Table 1**.

Genome-wide variant effects prediction. Genome-wide variant effects prediction is crucial for understanding the role of DNA mutations in species diversity. Genome-wide association studies (GWAS) provide valuable insights but often struggle to identify specific causal variants [30, 54]. The Genome Prediction Network (GPN) [32] addresses this by using unsupervised pre-training on genomic DNA sequences. During this process, GPN predicts nucleotides at masked positions within a 512-bp DNA sequence. This model is particularly effective at predicting rare variant effects, often missed by traditional GWAS methods. Additionally, models like DNABERT, DNABERT-2, and the Nucleotide Transformer also predict variant effects from DNA sequences. These advancements highlight ongoing efforts to better understand how DNA mutations contribute to biological diversity.

Cis-regulatory regions prediction. Cis-regulatory sequences, such as enhancers and promoters, play crucial roles in gene expression regulation, influencing development and physiology [55]. However, identifying these sequences remains a major challenge [56]. Pre-trained models like DNABERT, DNABERT-2, GROVER, and DNAGPT have been developed to predict promoter

regions and their activities with high accuracy. BERT-Promoter [57] utilizes a pre-trained BERT model for feature representation and SHAP analysis to filter data, improving prediction performance and generalization over traditional methods. Enhancers, which bind transcription factors to regulate gene expression [58, 59], are predicted by iEnhancer-BERT [60], which leverages DNABERT and uses a novel transfer learning approach. This model employs output from all transformer encoder layers and classifies features with a Convolutional Neural Network (CNN). These advancements highlight the growing trend of treating biological sequences as a natural language for computational modeling, offering new tools for identifying cis-regulatory regions and understanding their roles in diseases.

DNA-protein interaction prediction. Accurate identification of DNA-protein interactions is crucial for gene expression regulation and understanding evolutionary processes [61]. Several DNA language models, including DNABERT, DNABERT-2, and GROVER, have been developed to predict protein-DNA binding from ChIP-seq data. TFBert [62] is a pre-trained model specifically designed for DNA-protein binding prediction, which treats DNA sequences as natural sentences and k-mer nucleotides as words, allowing effective context extraction. Pre-trained on 690 ChIP-seq datasets, TFBert delivers strong performance with minimal fine-tuning. The MoDNA [63] framework introduces domain knowledge by incorporating common DNA functional motifs. During self-supervised pre-training, MoDNA performs tasks such as k-mer and motif prediction. Pre-training on extensive unlabeled genome data, MoDNA acquires semantic-level genome representations, enhancing predictions for promoter regions and transcription factor binding sites. Essentially, MoDNA functions as a biological language model for DNA-protein binding prediction.

DNA methylation prediction. DNA methylation is a key biological process in epigenetic regulation and is linked to various medical conditions and applications, such as metagenomic binning [64]. DNA methylation types depend on the nucleotide where the methyl group attaches [65]. Several models predict DNA methylation with varying accuracy. BERT6mA [66] is designed for predicting 6-methyadenine (6mA) sites, while iDNA-ABT [67], iDNA-ABF [68], and MuLan-Methyl [69] are versatile models predicting various methylation types (6mA, 5hmC, 4mC). iDNA-ABT, a deep learning model, integrates BERT with transductive information maximization (TIM), though it has yet to fully explore feature representation. iDNA-ABF uses a multi-scale architecture, applying multiple tokenizers for diverse embeddings, and MuLan-Methyl employs four

transformer-based models (DistilBERT [70], ALBERT[71], XLNet [72], and ELECTRA [73]) to predict methylation sites, enhancing performance through joint model utilization.

DNA level splice site identification. Accurate pre-mRNA splicing is essential for proper protein translation, driven by splice site selection. Identifying splice sites is challenging, particularly with prevalent GT-AG sequences [74]. To address this, DNABERT and DNABERT-2 were developed, trained on 10,000 donors, acceptor, and non-splice site sequences from the human reference genome to predict splice sites. DNABERT showed high attention to intronic regions, suggesting the functional role of intronic splicing enhancers and silencers as cis-regulatory elements in splicing regulation. This highlights DNABERT's potential in understanding splicing mechanisms.

4.2 Applications of large language models in transcriptomics

The RNA language models take RNA sequences as input, use transformer, BERT, GPT models to solve multiple biological tasks, including RNA 2D/3D structure prediction, RNA structural alignment,, RNA family clustering, RNA splice sites prediction from RNA sequence, RNA N7-methylguanosine modification prediction, RNA 2'-O-methylation modifications prediction, multiple types of RNA modifications prediction, predicting the association between miRNA, lncRNA and disease, identifying lncRNAs, lncRNAs' coding potential prediction, protein expression and mRNA degradation prediction (**Table 3, Supplementary Figure 1**). A detailed list of RNA language models, their downstream tasks, and the datasets used can be found in **Supplementary Table 1**.

Secondary structure prediction. RNA secondary structure prediction is a major challenge for RNA structural biologists, with models holding potential for RNA-targeting drug development [75]. Several RNA language models, such as RNABERT [8], RNA-MSM [5], RNA-FM [4], and UNI-RNA [35], have been developed to predict RNA structures with varying sophistication. RNABERT uses BERT architecture to predict structural features like base-pairing and stem loops. RNA-MSM integrates sequence and structural information to predict local and long-range folding patterns. RNA-FM focuses on RNA folding, stability, and energetics, including pseudoknots. UNI-RNA combines sequence and structure predictions across various RNA types. These models advance RNA structure prediction by applying deep learning and advanced techniques to improve understanding of RNA folding and function.

RNA splicing prediction. RNA splicing is crucial for gene expression in eukaryotes, and advancements have been made in sequence-based splicing modeling through models like

SpliceBERT [36] and UNI-RNA [35]. SpliceBERT, based on BERT, is trained to predict RNA splicing events by capturing long-range dependencies, identifying splice sites, and predicting alternative splicing events. UNI-RNA, a more generalized model, integrates multiple RNA tasks, including splicing, and combines sequence and structural data to predict splicing regulatory elements and interactions with splicing factors. These models enhance the understanding of RNA splicing, gene regulation, and its role in diseases, providing powerful tools for studying splicing defects and mutations.

lncRNAs identification and lncRNAs' coding potential prediction. Long non-coding RNAs (lncRNAs) play significant regulatory roles in cancer and diseases, and their small Open Reading Frames (sORFs), once thought weak in protein translation, are now known to encode peptides [76]. Identifying lncRNAs with sORFs is crucial for discovering new regulatory factors. LncCat [77] addresses this challenge by using category boosting and ORF-attention features, including BERT for peptide sequence representation, to improve prediction accuracy for both long ORF and sORF datasets. It demonstrates effectiveness across multiple species and Ribo-seq datasets in identifying lncRNAs with sORFs. In predicting translatable sORFs in lncRNAs (lncRNA-sORFs), LSCPP-BERT [78] is a novel method designed for plants, leveraging pre-trained transformer models for reliable coding potential prediction. LSCPP-BERT is poised to impact drug development and agriculture by enhancing understanding of lncRNA coding potential.

RNA–RBP interactions prediction. RNA sequences differ from DNA sequences by a single base (thymine to uracil), maintaining largely congruent syntax and semantics. BERT's versatility extends to Cross-linking and Immunoprecipitation data, particularly in predicting RNA-binding protein (RBP) binding preferences. BERT-RBP [79] is a model pre-trained on a human reference genome, designed to forecast RNA-RBP interactions. It outperforms existing models when tested on eCLIP-seq data from 154 RBPs and can identify transcript regions and RNA secondary structures based on sequence alone. BERT-RBP demonstrates BERT's adaptability in biological contexts and its potential to advance RNA-protein interaction understanding.

RNA–RNA interaction prediction. RNA–RNA interactions occur between various RNA species, including long non-coding RNAs, mRNAs, and small RNAs (e.g., miRNAs and lncRNAs), driven by complementary sequences, secondary structures, and other motifs [80]. Accurate prediction of these interactions provides insights into RNA-mediated regulation, enhancing understanding of biological processes like gene expression, splicing, and translation. RNAErnie, used for this

purpose, employs a TBTH architecture combining RNAErnie with a hybrid network (CNN, Bi-LSTM, and MLP) to predict RNA–RNA interactions. This approach demonstrates RNAErnie's potential in advancing RNA-based regulatory network studies.

RNA modification prediction. Post-transcriptional RNA modifications, such as N7-methylguanosine (m7G) and 2'-O-methylation (Nm), regulate gene expression and are linked to diseases [76, 81]. Identifying modification sites is essential but challenging due to the high cost and time required by experimental methods. Computational tools like BERT-m7G [82] and Bert2Ome [83] address this issue. BERT-m7G uses a stacking ensemble approach to identify m7G sites directly from RNA sequences, offering an efficient, cost-effective alternative. Bert2Ome combines BERT and CNN to predict 2'-O-methylation sites, outperforming existing methods across datasets and species. These tools enhance the accuracy, scalability, and efficiency of RNA modification site identification, advancing research into RNA modifications and their roles in gene regulation and disease.

Protein expression and mRNA degradation prediction. mRNA vaccines are a cost-effective, rapid, and safe alternative to traditional vaccines, showing high potency [84]. These vaccines work by introducing mRNA that encodes a viral protein. CodonBERT [85] is a model specifically designed for mRNA sequences to predict protein expression. It uses a multi-head attention transformer architecture and was pre-trained on 10 million mRNA sequences from various organisms. This pre-training enables CodonBERT to excel in tasks like protein expression and mRNA degradation prediction. Its ability to integrate new biological information makes it a valuable tool for mRNA vaccine development. CodonBERT surpasses existing methods, optimizing mRNA vaccine design and improving efficacy and applicability in immunization. Its strength in predicting protein expression enhances mRNA vaccine development efficiency and effectiveness.

5' UTR-based mean ribosome loading prediction and mRNA subcellular localization prediction. The 5' UTR sequence plays a critical role in regulating translation efficiency. RNA sequence models like 3UTRBERT, UNI-RNA, UTR-LM, RNA-FM, and Nucleotide Transformer have been developed to predict key features of the 5' UTR, focusing on ribosome loading efficiency and mRNA localization. These models use Transformer-based architecture to analyze sequence patterns, motifs, and structural elements. For example, 3UTRBERT [37] and RNA-FM [4] predict ribosome loading efficiency, identifying regions likely to recruit ribosomes for

translation initiation. UTR-LM [38], UNI-RNA [35], and Nucleotide Transformer [31] predict mRNA subcellular localization, determining where mRNA will localize in the cell (cytoplasm, ribosomes, or nucleus), which is crucial for regulating mRNA stability and translation. Together, these models provide valuable insights into gene expression, translation control, and RNA localization, advancing molecular biology research.

4.3 Applications of large language models in proteomics

Protein is an indispensable molecule in life, assuming a pivotal role in the construction and sustenance of vital processes. As the field of protein research advances, there has been a substantial surge in the accumulation of protein data [86]. In this context, the utilization of large language models emerges as a viable approach to extract pertinent and valuable information from these vast reservoirs of data. Several pre-trained protein language models (PPLMs) have been proposed to learn the characteristic representations of proteins data (e.g., protein sequences, gene ontology annotations, property descriptions), then applied to different tasks by fine-tuning, adding or altering downstream networks, such as protein structure, post-translational modifications (PTMs), and biophysical properties, which align with corresponding downstream tasks like secondary structure prediction, major PTMs prediction, and stability prediction [87, 88].

Even though antibodies are classified as proteins, the datasets of antibodies and subsequent tasks differ significantly from those of proteins. Through the establishment and continuous updates of the Observed Antibody Space (OAS) database [89], a substantial amount of antibody sequence data has become available, which can be utilized to facilitate the development of pre-trained antibody large language models (PALMs). PALMs primarily delve into downstream topics encompassing therapeutic antibody binding mechanisms, immune evolution, and antibody discovery, which correspond to tasks like paratope prediction, B cell maturation analysis, and antibody sequence classification (**Table 3, Supplementary Figure 2**).

In this section, some of the popular protein-related large language models of recent years are introduced, as well as corresponding important downstream tasks. It is important to emphasize that both PPLM and PALM are capable of handling not only the downstream tasks introduced in this section. For further details, additional information can be referenced within **Supplementary Table 2**.

Secondary structure and contact prediction. Protein structure is critical to its function and interactions [90]. However, traditional experimental techniques for protein structure analysis are

time-consuming and labor-intensive. With the rise of deep learning, large language models have demonstrated significant advantages in computational efficiency and prediction accuracy for protein structure prediction [91]. MSA Transformer [92] introduces a protein language model that processes MSAs using a unique mechanism of interleaved row and column attention. Trained with a MLM objective across diverse protein families, it outperformed earlier unsupervised approaches and showed greater parameter efficiency than previous models. Drawing on insights from BERT, large parameter models tend to achieve better performance for predicting secondary structures and contacts. Few models have more parameters than the largest models in ProtTrans [11], which includes a series of autoregressive models (Transformer-XL [93], XLNet [72]) and four encoder (BERT [2], Albert [71], Electra [73], T5 [94]) trained on datasets like UniRef [95] and BFD [96], comprising up to 393 billion amino acids. Model sizes vary from millions to billions of parameters. Notably, ProtTrans made a significant breakthrough in per-residue predictions.

Protein sequence generation. Protein sequence generation holds significant potential in drug design and protein engineering [97]. Using machine learning or deep learning, generated sequences aim for good foldability, stable 3D structures, and specific functional properties, such as enzyme activity and antibody binding. The development of large language models, combined with conditional models, has greatly advanced protein generation [98]. ProGen [12] incorporates UniprotKB keywords as conditional tags, covering over 1,100 categories like 'biological process' and 'molecular function'. Proteins generated by ProGen, assessed for sequence similarity, secondary structure, and conformational energy, exhibit desirable structural properties. In 2022, ProtGPT2 [41] inspired by GPT-x was developed. ProtGPT2-generated proteins show amino acid propensities like natural proteins. Prediction of disorder and secondary structure reveals that 88% of these proteins are globular, resembling natural sequences. Employing AlphaFold [99, 100] on ProtGPT2 sequences produces well-folded, non-idealized structures with unique topologies not seen in current databases, suggesting ProtGPT2 has effectively learned "protein language".

Protein function prediction. Proteins are essential in cellular metabolism, signal transduction, and structural support, making their function critical for drug development and disease analysis. However, predicting and annotating protein functions is challenging due to their complexity. PPLMs offer effective solutions to these challenges [101, 102]. ProtST [103] introduced a multimodal framework combining a PPLM for sequences and a biomedical language model (BLM) for protein property descriptions. Through three pre-training tasks, unimodal mask prediction,

multimodal representation alignment, and multimodal mask prediction, the model excels in tasks like protein function annotation, zero-shot classification, and functional protein retrieval from large databases. While most methods focus on increasing model parameters to improve performance, CaLM [14] introduces an alternative representation, the cDNA sequence, akin to an amino acid sequence, as input. The core idea lies in the relationship between synonymous codon usage and protein structure [104], and the information encoded in codons is no less than that of amino acids. Experimental results demonstrate that even with a small parameter language model, using cDNA sequences as input enhances performance in tasks such as protein function prediction, species recognition, prediction of protein and transcript abundance, and melting point estimation.

Major post-translational modification prediction. Post-translational modifications (PTMs) are chemical changes, such as phosphorylation, methylation, and acetylation, that alter protein structure and function after translation. PTMs influence protein stability, localization, interactions, and function, making their study crucial for disease diagnosis and therapeutic strategies [105, 106]. Language models can effectively predict PTMs and related tasks like signal peptide prediction. ProteinBERT [42], with only ~16M parameters, is not large enough but performs well due to its inclusion of Gene Ontology (GO) annotation tasks. By incorporating GO interactions with protein sequences, ProteinBERT achieves strong performance on PTM prediction and other protein property benchmarks, outperforming models with larger parameter sizes.

Evolution and mutation prediction. Protein evolution and mutation drive functional diversity, aiding adaptation to environmental changes and offering insights into protein function origin, which can inform drug development and disease treatment [107, 108]. UniRep [109], built on the LSTM architecture, was trained on the UniRef50 [95] and excelled in tasks like remote homology detection and mutation effect prediction. ESM-1b [40], a deep transformer model trained on 250 million sequences, with 33 layers and 650 million parameters, captures essential protein sequence patterns through self-supervised learning. ESM-1b is also integral to frameworks like PLMSearch [110] and DHR [111], which enable fast, sensitive homology searches. PLMSearch uses supervised training, while DHR relies on unsupervised contrastive learning and enhances structure prediction models like AlphaFold2 [100].

Biophysical properties prediction. Biophysical properties of proteins, such as fluorescence and stability landscapes [112], are crucial for understanding protein folding, stability, and conformational changes, with significant implications for drug design, protein engineering, and

enzyme engineering. Deep learning advancements have enabled more accurate prediction of these properties using PPLMs. TAPE benchmark [39] established standardized tasks for evaluating protein, including fluorescence and stability landscape prediction. In 2022, PromptProtein [113], a prompt-based pre-trained model, incorporated multi-task pre-training and a fine-tuning module to improve task-specific performance. It outperformed existing methods in function and biophysical properties prediction, demonstrating substantial gains in predictive accuracy.

Protein-protein interaction and binding affinity prediction. Protein-protein interactions (PPIs) are crucial for biological functions, and their prediction is also vital for drug discovery and design. PPLMs provide efficient, accurate predictions of PPI types and binding affinities [114, 115]. KeAP model [43], like ProtST, aims to integrate fine-grained information beyond OntoProtein [116]. KeAP uses a triplet format (Protein, Relation, Attribute) as input, processed by encoders and a cascaded decoder based on the Transformer architecture. Using MLM for pre-training, KeAP employs a cross-attention fusion mechanism to capture detailed protein information, achieving superior performance on tasks such as PPI identification and binding affinity estimation.

Antigen-receptor binding and antigen-antibody binding prediction. Antigen proteins are processed into neoantigen peptides that bind to the Major Histocompatibility Complex (MHC), forming pMHC complexes. These complexes are presented to T-cells, stimulating antibody production by B-cells, which triggers an immune response [117]. Predicting peptide binding to MHC molecules is a key focus of language models in this process [118, 119]. MHCRoBERTa [120] uses a pretrained BERT model to predict pMHC-I binding by learning the biological meaning of amino acid sequences. BERTMHC [121], trained on 2,413 MHC-peptide pairs, focuses on pMHC-II binding prediction, filling a gap in this area.

Another goal is predicting the binding specificity of adaptive immune receptors (AIRs), particularly TCRs. TCR-BERT [122] learns TCR CDR3 sequences to predict antigen specificity but lacks the ability to model the interaction between TCR chains. SC-AIR-BERT [123] addresses this by pre-training a model that outperforms others in TCR and BCR binding specificity. Additionally, the Antiformer [124] integrates RNA-seq and BCT-seq data in a graph-based framework to improve antibody development. In antibody modeling, three recent models focus on unique tasks. AbLang [125], built on RoBERTa [126], excels at restoring lost residues during sequencing and outperforms other models in accuracy and efficiency. AntiBERTa [127] understands antibody "language" through tasks like predicting immunogenicity and binding sites.

EATLM [128], with its unique pre-training tasks (Ancestor Germline Prediction and Mutation Position Prediction), contributes a reliable benchmark for antibody language models.

4.4 Applications of large language models in drug discovery

Drug discovery is an expensive and long-term process that exhibits a low success rate. During the early stages of drug discovery, computer-aided drug discovery, employing empirical or expert knowledge algorithms, machine learning algorithms, and deep learning algorithms, serve to accelerate the generation and screening of drug molecules and their lead compounds [129-131]. It speeds up the entire drug discovery process, especially the development of small molecule drugs. Among commonly used medications, small molecule drugs can account for up to 98% of the total [132]. The structure of small molecule drugs exhibits excellent spatial dispersibility, and their chemical properties determine their good drug-like properties and pharmacokinetic properties [133]. With the development of deep learning and the proposal of large language models, it has become easy to apply these methods to discover hidden patterns of molecules and interactions between molecules for drugs (such as small molecules) and targets (such as proteins and RNA) that can be easily represented as sequence data. The Simplified Molecular-Input Line-Entry System (SMILES) string and chemical fingerprint are commonly used to represent molecules. Additionally, through the pooling process of graph neural networks(GNN), small molecules can be transformed into sequential representations [134]. With the protein sequence, large language models can engage in drug discovery through various inputs. Within this section, key tasks within the early drug discovery process that have effectively leveraged large language models will be introduced (**Table 3, Supplementary Figure 3**). A detailed list of drug discovery language models, their downstream tasks, and the datasets used can be found in **Supplementary Table 3**.

Drug-like molecular properties prediction. In drug discovery, significant focus is placed on properties like ADMET and PK to develop more effective, accessible, and safe drugs[135, 136]. Large language models (LLMs) are used for molecular property prediction, including these properties. Since molecular SMILES representations are consistent, models can be easily improved and fine-tuned for specific tasks based on researchers' requirements. SMILES-BERT [17] departed from the usage of knowledge-based molecular fingerprints as input. Instead, it adopted a representation method where molecules were encoded as SMILES sequences and employed as input for both pre-training and fine-tuning within a BERT-based model. This novel approach yielded superior outcomes across various downstream molecular property prediction tasks,

surpassing the performance of previous models reliant on molecular fingerprints. ChemBERTa [137] is a BERT-based model that focuses on the scalability of large language models, exploring the impact of pre-training dataset size, tokenizer, and string representation. Subsequently, ChemBERTa-2[138] improved upon ChemBERTa by using a larger dataset of 77 million compounds from PubChem, enhancing its ability to learn from diverse chemical structures. It also integrates advanced self-supervised learning techniques and fine-tuning strategies, resulting in better generalization performance across various downstream tasks. K-BERT [15] stands out by using three pre-training tasks: atom feature prediction, molecular feature prediction, and contrastive learning. This approach enables the model to understand the essence of SMILES representations, resulting in exceptional performance across 15 drug datasets, highlighting its effectiveness in drug discovery. Given the importance of graph neural networks in the development of molecular pre-training models, Mole-BERT [139] introduces atom-level Masked Atoms Modeling (MAM) task and graph-level Triplet Masked Contrastive Learning (TMCL) task. These tasks enable the network to acquire a comprehensive understanding of the “language” embedded within molecular graphs. By adopting this approach, the network demonstrates exceptional performance across eight downstream tasks, showcasing its adaptability and effectiveness in diverse applications.

Drug-like molecules generation. It is very difficult to chase the full coverage of the enormous drug-like chemical space (estimated at more than 10^{63} compounds), and traditional virtual screening libraries usually contain less than 10^7 compounds and are sometimes not available. In such circumstances, the utilization of deep learning methods to generate molecules exhibiting drug-like properties emerges as a viable approach [140, 141]. Inspired by the generative pre-training model GPT, MolGPT [16] model was introduced. In addition to performing the next token prediction task, MolGPT incorporates an extra training task for conditional prediction, facilitating the capability of conditional generation. Beyond its capacity to generate innovative and efficacious molecules, the model has demonstrated an enhanced ability to capture the statistical characteristics within the dataset.

Drug-target interaction predictions. The investigation of Drug-Target Interaction (DTI) holds paramount significance in the realm of drug development and the optimization of drug therapy. Understanding drug-target interactions aids in pharmaceutical design, accelerates drug development, and reduces time and resource costs in lab experimentation and trial-and-error

methods [142, 143]. During the exploration of DTI, diligent focus is placed on the prediction of drug-target binding affinity. DTI-BERT employs a fine-tuned ProtBERT [144] model to process protein sequences and applies discrete wavelet transform to drug molecular fingerprints.. TransDTI [145] is a multi-class classification and regression workflow. This model not only uses fine-tuned SMILES-BERT to extract drug features, but also expands the selection of fine-tuned large protein models. After acquiring potential representations of drug-target pairs, the authors subject the representations to downstream neural networks for the completion of a multi-classification task. Additionally, The Chemical-Chemical Protein-Protein Transferred DTA (C2P2) [146] method uses pre-trained protein and molecular large language models to capture the interaction information within molecules. Given the relatively limited scale of the DTI dataset, C2P2 leverages protein-protein interaction (PPI) and chemical-chemical interaction (CCI) tasks to acquire knowledge of intermolecular interactions and subsequently transfer this knowledge to affinity prediction tasks [147]. It is worth highlighting that in scenarios involving the docking or when emphasizing the spatial structure of a complex, methodologies incorporating 3D convolution networks, point clouds-based networks, and graph networks are often employed [148-151]. In situations where the molecular structure is unknown, but the sequence is available, the prediction of DTI using large-scale models still holds significant promise.

Drug synergistic effects predictions. Combination therapy is common for complex diseases like cancer, infections, and neurological disorders, often surpassing single-drug treatments. Predicting drug pair synergy, where combining drugs boosts therapeutic effects, is vital in drug development. However, it's challenging due to many drug combinations and complex biology [152, 153]. Various computational methods, including machine learning, help predict drug pair synergy. Carl Edwards et al. introduced SynerGPT [48], which is based on GPT trained to in-context learn drug synergy functions without relying on domain-specific knowledge. Wei Zhang et al. [154] introduced DCE-DForest [154], a model for predicting drug combination synergies. It uses a pretrained drug BERT model to encode the drug SMILES and then predicts synergistic effects based on the embedding vectors of drugs and cell lines using the deep forest method. Mengdie Xua et al. [155] utilized a fine-tuned pre-trained large language model and a dual feature fusion mechanism to predict synergistic drug combinations. Its input includes hashed atom pair molecular fingerprints of drugs, SMILES string encodings, and cell line gene expressions. They conducted ablation analyses on the dual feature fusion network for drug-drug synergy prediction, highlighting

the significant role of fingerprint inputs in ensuring high-quality drug synergy predictions.

4.5 Applications of large language models in single-cell analysis

Large language models have demonstrated significant applications in single-cell analysis, including cell-level tasks such as identifying cell types, determining cell states, and discovering novel cell populations; gene-level tasks like inferring gene regulatory networks; and multi-omics tasks, such as integrating single-cell multi-omics (scMulti-omics) data (**Supplementary Figure 4**). Additionally, this section will explore emerging language models based on spatial transcriptomics (**Table 3**). A detailed list of single-cell large language models, their downstream tasks, and the datasets used can be found in **Supplementary Table 4**.

Cell-level tasks. Cell-level tasks, such as cell clustering, cell type annotation, novel cell type discovery, batch effect removal and trajectory inference, are critical in single-cell analysis. These tasks often rely on cell representations learned during pretraining, which are subsequently fine-tuned for different tasks. Single-cell language models derive cell representations in two primary ways. The first method utilizes a special class token (<cls>) appended to the input sequence; its embedding is updated through the transformer layers, and the final embedding at the <cls> position serves as the cell representation. The second method generates a cell embedding matrix from the model output, where each row represents a specific cell. Both approaches facilitate downstream tasks, as demonstrated by TOSICA [22], which uses the <cls> token to predict cell type probabilities using the whole conjunction neural network cell type classifier to annotate single cells, and iSEEK [156], which generates cell embedding for cell clustering, cell type annotation, and developmental trajectory exploration. Models like scBERT [20] and UCE [49] leverage multi-head attention mechanisms to extract information from diverse representation subspaces, discerning subtle differences between novel and known cell types. Their large receptive fields capture long-range gene-gene interactions, enabling comprehensive characterization of novel cellular states. Addressing batch effects, which arise from variations in species, tissues, operators, and experimental protocols, remains a significant challenge in single-cell analysis. Large language models, pretrained on extensive datasets, utilize attention mechanisms to incorporate prior biological knowledge, enabling batch-insensitive data annotation. Without relying on explicit batch information, models such as CIForm [152] have demonstrated effectiveness in both intra-dataset and inter-dataset scenarios. They handle annotations across diverse species, organs, tissues, and technologies while also supporting the integration of reference and query data from various

sequencing platforms or studies. This capability allows them to address batch effects in single-cell analysis. Drug response or sensitivity prediction is a classification task akin to cell type annotation, where a classifier is appended to the learned cell embeddings to predict whether a cell will respond to or exhibit sensitivity to a specific drug. Models like scFoundation [25] and CellLM [157] effectively utilize this approach, leveraging the robust cell representations learned during pretraining to enhance prediction accuracy.

Gene-level tasks. Gene-level tasks, such as gene expression prediction, gene regulatory network (GRN) inference, gene perturbation prediction, and drug dose-response prediction, are integral to understanding single-cell transcriptomics. Self-attention mechanisms have transformed deep learning by enabling context-aware models that prioritize relevant elements in large input spaces. These models, particularly transformers, are well-suited for modeling the context-dependent dynamics of gene regulatory networks. By focusing on key interactions, transformers can effectively capture the complexities of regulatory relationships, such as the attention matrix in Geneformer [18] and scGPT [21] reflect which genes that gene pays attention to and which genes pay attention to that gene, aiding to infer gene regulation network. Geneformer is pretrained on a vast repository of single-cell transcriptomes to learn gene relationships for diverse downstream applications, including predicting dosage-sensitive disease genes, identifying downstream targets, forecasting chromatin dynamics, and modeling network dynamics. In addition, after pretraining and fine-tuning, single-cell language models output gene embeddings that can be utilized for functional analysis of scRNA-seq data. For instance, scGPT [21] serves as a generalizable feature extractor leveraging zero-shot learning, enabling applications in gene expression prediction and genetic perturbation prediction. Similarly, in scFoundation [25], zero-expressed genes and masked genes are combined with the output from the transformer-based encoder. This combined information is then input into the decoder and projected to gene expression values through a multilayer perceptron (MLP). The gene context expression is employed to formulate a cell-specific gene graph, facilitating the prediction of perturbations using the GEARS [158] model. It is worth noting that genes have a lot of prior knowledge that can be used to enhance many gene-level tasks. For example, GeneCompass [51] incorporates four types of biological prior knowledge, including GRN, promoter information, gene family annotation and gene-co-expressed relationship, making it capable for various gene tasks.

scMulti-omics tasks. Studying single-cell multi-omics data requires integrating diverse information from genomics, transcriptomics, epigenomics, and proteomics at the single-cell level. The adaptability, generalization capabilities, and feature extraction strengths of large language models make them effective in addressing challenges such as feature variance, data sparsity, and cell heterogeneity inherent in single-cell multi-omics datasets. scMulti-omics integration can be viewed as a specialized form of batch effect removal. For example, scGPT [21] treats each modality as a distinct batch and incorporates a special modality token to represent input features (such as genes, regions, or proteins) associated with each modality. This approach helps the transformer model balance attention across modalities, preventing overemphasis on intra-modality features while integrating inter-modality relationships effectively. Another approach involves processing different modalities through separate transformers before projecting their embeddings into a common latent space. Models like scMVP [159] use mask attention-based encoders for scRNA-seq data and transformer-based multi-head self-attention encoders for scATAC-seq. By aligning variations between different omics in this latent space, scMVP captures joint profiling of scRNA-seq and scATAC-seq, achieving paired integration where gene expression and chromatin accessibility are studied within the same cells. Graphs are increasingly recognized as powerful tools for characterizing feature heterogeneity in scMulti-omics integration. For example, DeepMAPS [160] leverages graph transformers to construct cell and gene graphs, learning both local and global features that build cell-cell and gene-gene relationships for data integration, inference of biological networks from scMulti-omics data and cell-cell communication.

Recent advances in sequencing technologies that capture multiple modalities within the same cell have enabled the development of computational tools for cross-modality prediction. One approach involves training large language models on paired datasets to predict one modality from another. For instance, scTranslator [161], pre-trained on paired bulk and single-cell data, fine-tunes to infer protein abundance from scRNA-seq data by minimizing the mean squared error (MSE) between predicted and actual protein levels. Another strategy leverages graph learning with prior knowledge to model feature relationships. For example, scMoFormer [162] can not only translate gene expression to protein abundance, but is also applicable to multi-omics predictions, including protein abundance to gene expression, chromatin accessibility to gene expression, gene expression to chromatin accessibility using graph transformers. Taking protein prediction task as an example, scMoFormer constructs cell-gene graph, gene-gene graph, protein-protein graph, and gene-protein

graph based on gene expression profiles and prior knowledge from STRING database [163]. Each modality has a separate transformer to learn the global information that may not be included in prior knowledge. Message-passing graph neural networks (GNNs) link nodes across various graphs, while transformers are employed to precisely map gene expression to protein abundance.

Spatial transcriptomics tasks. The rapid development of single-cell and spatial transcriptomics has advanced our understanding of cellular heterogeneity and tissue architecture. Spatial transcriptomics retains cells' native spatial context, enabling insights into cellular interactions. Large language models address the challenge of high-dimensional spatial data analysis by integrating spatial and molecular information, enhancing tissue-specific pattern interpretation. For example, Nicheformer [53] is the latest large language model in spatial transcriptomics. It integrates extensive spatial transcriptomics and single-cell transcriptomics data, leveraging metadata across multiple modalities, species, and sequencing technologies. By doing so, Nicheformer is capable of learning joint information from single-cell and spatial transcriptomics, enabling the resolution of various spatial prediction tasks even with limited data. Spaformer [164] is another transformer-based model designed for spatial transcriptomics data. Spaformer is designed to address two key challenges: how to encode spatial information of cells into a transformer model and how to train a transformer to overcome the sparsity of spatial transcriptomics data, enabling data imputation. Spatial transcriptomics, as one of the most popular technologies in recent years, focuses on integrating single-cell resolution gene expression data with tissue spatial information to reveal spatial relationships and functional characteristics among cells. However, large language models (LLMs) specifically designed for spatial transcriptomics are still in their early stages of development. The creation of these models faces unique challenges, such as effectively integrating high-dimensional gene expression data with complex spatial information and addressing the sparsity and irregularity of the data.

In addition to the single-cell large language models discussed above, another category of single-cell prediction models leverages natural language, utilizing textual data such as human-readable descriptions of gene functions and biological features to support various single-cell analyses. For example, GPT-4 [165] leverages its strong contextual understanding for interpreting high-dimensional single-cell analysis for accurate cell type annotation. GenePT [166] utilizes OpenAI's ChatGPT text embedding to classify gene properties and cell types effectively. More and more models demonstrate that natural language pretraining can significantly boost performance on

single-cell downstream tasks, including cell generation [167], cell identity (e.g., cell type, pathway, and disease information) [167-171], and gene enrichment analysis [169]. These models demonstrate significant potential in advancing single-cell analysis by integrating natural language processing techniques. However, the reliance on textual data may constrain performance in less-annotated or novel datasets.

5. Conclusion and Suggestions on large language models in bioinformatics

5.1 Summary of large language models in bioinformatics

Large language models (LLMs) have catalyzed transformative progress across biological disciplines, including genomics, transcriptomics, proteomics, drug discovery, and single-cell analysis. These models, trained on vast datasets, address challenges like the sparsity, high dimensionality, and heterogeneity of biological data while capturing the complexity of sequence relationships. Tokenization methods are pivotal, converting sequences into manageable formats, such as, for genomics and transcriptomics, k-mer encoding is prevalent, segmenting DNA/RNA sequences into overlapping units. In proteomics, amino acid residue-based tokenization captures protein structure and function. These preprocessing strategies enable LLMs to interpret biological language effectively.

Representation learning allows LLMs to uncover contextual and hierarchical relationships within biological data, forming the basis for various downstream applications. These tasks can be grouped into four primary categories: 1) Classification/Prediction Tasks: Examples include identifying functional genomic elements (e.g., promoters, enhancers), predicting protein structures and interactions, and cell type annotation in single-cell data. 2) Generation Tasks: LLMs can create biologically relevant sequences, such as gene expression imputation and synthetic DNA, RNA, or protein sequences, aiding in vaccine development or enzyme engineering. 3) Interaction Tasks: These involve modeling interactions like drug-target binding, cell-cell interaction, protein-protein interactions, or cross-omics relationships (e.g., gene expression to protein abundance). 4) Transfer Learning Tasks: Pretrained LLMs, such as DNABERT and scGPT, are fine-tuned for specific applications, including single-cell data annotation or predicting RNA modifications like N6-methyladenosine sites. Despite their capabilities, challenges persist. Biological data often exhibit sparsity, as seen in single-cell and spatial transcriptomics, and irregularity due to sequencing errors or noise. To address this, LLMs must effectively integrate multi-modal data, balance computational efficiency, and ensure interpretability of their outputs. As foundational models

evolve, their ability to unify diverse biological datasets into a single framework for prediction, generation, interaction, and transfer learning tasks will continue to reshape our understanding and applications of biological systems.

5.2 Guidance on how to use and develop LLMs in practice

Large Language Models offer immense potential in bioinformatics and other fields, but their effective utilization and development require distinct approaches for end-users and developers (**Figure 5**).

For LLM users, the process begins by clearly defining the research domain and task, specifying the relevant omics level (e.g., genomics, transcriptomics, proteomics) and identifying whether the objective involves classification or prediction, generation, interaction, or transfer learning. A well-defined objective streamlines the selection of appropriate models and workflows. Next, users should choose models pretrained on data relevant to their domain, as detailed in **Table 2**, which includes information on foundation models, their training data types, and availability. For instance, DNABERT is ideal for genomics tasks, while scGPT is tailored for single-cell analysis. Additionally, users must assess computational requirements and ensure compatibility with their dataset size and complexity. Proper data preparation is critical, including aligning data with model requirements, addressing missing values, and incorporating metadata like cell types or genomic regions. **Table 1** provides common tokenization methods for reference. To leverage transfer learning, users can fine-tune foundation models listed in **Table 2** for their specific dataset, optimizing performance through hyperparameter tuning, early stopping, and cross-validation. Alternatively, users can utilize predeveloped models listed in **Supplementary Tables 1–4** for similar tasks to obtain results directly. Finally, rigorous evaluation using metrics like accuracy, precision, and recall is essential, complemented by interpretation tools such as attention maps or feature embeddings to extract meaningful biological insights (**Figure 5**).

For LLM developers, it is essential to **first** understand domain-specific challenges to address issues like sparsity, heterogeneity, and high dimensionality. For example, single-cell and spatial transcriptomics datasets often suffer from sparsity and noise, necessitating innovative solutions in model architecture. **Second**, developers should choose or develop tokenization strategies tailored to biological data. For instance, k-mer encoding works well for DNA/RNA sequences, while gene ranking-based tokenization is effective for scRNA-seq data. Exploring hybrid tokenization can enhance cross-modal understanding. **Third**, in model development, developers should employ or

design novel transformer structures. For example, scBERT utilizes Performer to improve scalability. Incorporating knowledge-based information into model training can further enhance performance. For instance, GeneCompass integrates four types of biological prior knowledge including GRNs, promoter information, gene family annotation, and gene co-expression relationships, making it versatile for various gene-related tasks. Similarly, basic protein language models, which are often limited to MSA and protein sequences, can be improved by incorporating additional modalities like 3D structural data. This can be achieved by converting such modalities into sequence formats or integrating large models to collectively capture multi-modal information using fusion techniques. Moreover, combining Graph Neural Networks (GNNs) with transformers has led to significant advancements. For example, scMoFormer constructs cell-gene, gene-gene, protein-protein, and gene-protein graphs for multi-omics predictions, while DeepMAPS uses cell-gene graphs to estimate gene importance. GNNs excel in capturing local interactions, while transformers effectively model long-range dependencies, enabling comprehensive representations of intricate relationships in single-cell data. **Fourth**, novel tasks that can be explored in developing LLMs for bioinformatics include causal inference in multi-omics, such as determining how DNA variations influence mRNA abundance or protein expression. Spatial transcriptomics interpretation can model cell spatial organization within tissues. Epigenetic modulation prediction focuses on regulatory roles of histone modifications, DNA methylation, or chromatin accessibility. Synthetic biology applications can involve generating optimized gene or protein sequences, while cross-species genomics identifies conserved functional genomic elements. These tasks exemplify how LLMs can tackle emerging challenges in biological research. **Fifth**, developers should expand LLMs to accommodate emerging data types, such as CODEX imaging data and long-read sequencing data, which bring unique challenges in terms of data structure, preprocessing, and representation. **Lastly**, validation, application, and interpretability should be prioritized. Developers should not only evaluate models on specific tasks but also ensure that foundational challenges, such as the impact of sparsity in scRNA-seq data on cell type annotation performance, are fully addressed to enhance the robustness and utility of the models (**Figure 5**).

Acknowledgements

We would like to express our gratitude to our colleagues and friends who provided invaluable advice and support throughout the duration of this study.

Funding

This work was partially supported by the National Institutes of Health [R01LM014156, R01GM153822, R01CA241930 to X.Z] and the National Science Foundation [2217515, 2326879 to X.Z]. The funders had no role in study design, data collection and analysis, decision to publish or preparation of the manuscript. Funding for open access charge: Dr & Mrs Carl V. Vartian Chair Professorship Funds to Dr. Zhou from the University of Texas Health Science Center at Houston.

Conflict of interest statement. None declared.

References

1. Radford, A., et al., *Improving language understanding by generative pre-training*. 2018.
2. Devlin, J., et al., *Bert: Pre-training of deep bidirectional transformers for language understanding*. arXiv preprint arXiv:1810.04805, 2018.
3. Vaswani, A., et al., *Attention is all you need*. Advances in neural information processing systems, 2017. **30**.
4. Chen, J., et al., *Interpretable RNA foundation model from unannotated data for highly accurate RNA structure and function predictions*. bioRxiv, 2022: p. 2022.08. 06.503062.
5. Zhang, Y., et al., *Multiple sequence alignment-based RNA language model and its application to structural inference*. Nucleic Acids Research, 2024. **52**(1): p. e3-e3.
6. Ji, Y., et al., *DNABERT: pre-trained Bidirectional Encoder Representations from Transformers model for DNA-language in genome*. Bioinformatics, 2021. **37**(15): p. 2112-2120.
7. Zhang, D., et al., *DNAGPT: A Generalized Pretrained Tool for Multiple DNA Sequence Analysis Tasks*. bioRxiv, 2023: p. 2023.07. 11.548628.
8. Akiyama, M. and Y. Sakakibara, *Informative RNA base embedding for RNA structural alignment and clustering by deep representation learning*. NAR genomics and bioinformatics, 2022. **4**(1): p. lqac012.
9. Wang, N., et al., *Multi-purpose RNA language modelling with motif-aware pretraining and type-guided fine-tuning*. Nature Machine Intelligence, 2024: p. 1-10.
10. Rives, A., et al., *Biological structure and function emerge from scaling unsupervised learning to 250 million protein sequences*. bioRxiv. 2019.
11. Elnaggar, A., et al., *ProtTrans: Towards Cracking the Language of Lifes Code Through Self-Supervised Deep Learning and High Performance Computing*. IEEE Transactions on Pattern Analysis and Machine Intelligence, 2021: p. 1-1.
12. Madani, A., et al., *Progen: Language modeling for protein generation*. arXiv preprint arXiv:2004.03497, 2020.
13. Xu, M., et al., *Protst: Multi-modality learning of protein sequences and biomedical texts*. arXiv preprint arXiv:2301.12040, 2023.
14. Outeiral, C. and C.M. Deane, *Codon language embeddings provide strong signals for use in protein engineering*. Nature Machine Intelligence, 2024. **6**(2): p. 170-179.
15. Wu, Z., et al., *Knowledge-based BERT: a method to extract molecular features like computational chemists*. Briefings in Bioinformatics, 2022. **23**(3): p. bbac131.

16. Bagal, V., et al., *MolGPT: molecular generation using a transformer-decoder model*. Journal of Chemical Information and Modeling, 2021. **62**(9): p. 2064-2076.
17. Wang, S., et al. *Smiles-bert: large scale unsupervised pre-training for molecular property prediction*. in *Proceedings of the 10th ACM international conference on bioinformatics, computational biology and health informatics*. 2019.
18. Theodoris, C.V., et al., *Transfer learning enables predictions in network biology*. Nature, 2023. **618**(7965): p. 616-624.
19. Shen, H., et al., *Generative pretraining from large-scale transcriptomes for single-cell deciphering*. iScience, 2023. **26**(5): p. 106536.
20. Yang, F., et al., *scBERT as a large-scale pretrained deep language model for cell type annotation of single-cell RNA-seq data*. Nature Machine Intelligence, 2022. **4**(10): p. 852-866.
21. Cui, H., et al., *scGPT: toward building a foundation model for single-cell multi-omics using generative AI*. Nat Methods, 2024. **21**(8): p. 1470-1480.
22. Chen, J., et al., *Transformer for one stop interpretable cell type annotation*. Nat Commun, 2023. **14**(1): p. 223.
23. Jiao, L., et al., *scTransSort: Transformers for Intelligent Annotation of Cell Types by Gene Embeddings*. Biomolecules, 2023. **13**(4).
24. Xiong, L., T. Chen, and M. Kellis. *scCLIP: Multi-modal Single-cell Contrastive Learning Integration Pre-training*. in *NeurIPS 2023 AI for Science Workshop*.
25. Hao, M., et al., *Large-scale foundation model on single-cell transcriptomics*. Nat Methods, 2024. **21**(8): p. 1481-1491.
26. Bian, H., et al. *scMulan: a multitask generative pre-trained language model for single-cell analysis*. in *International Conference on Research in Computational Molecular Biology*. 2024. Springer.
27. Wen, H., et al., *CellPLM: pre-training of cell language model beyond single cells*. bioRxiv, 2023: p. 2023.10.03.560734.
28. Mao, Y., et al., *Phenotype prediction from single-cell RNA-seq data using attention-based neural networks*. Bioinformatics, 2024. **40**(2).
29. Querfurth, B.v., et al., *mcBERT: Patient-Level Single-cell Transcriptomics Data Representation*. bioRxiv, 2024: p. 2024.11.04.621897.
30. Sarkar, S., *Decoding " coding": Information and DNA*. BioScience, 1996. **46**(11): p. 857-864.
31. Dalla-Torre, H., et al., *The nucleotide transformer: Building and evaluating robust foundation models for human genomics*. bioRxiv, 2023: p. 2023.01.11.523679.
32. Benegas, G., S.S. Batra, and Y.S. Song, *DNA language models are powerful predictors of genome-wide variant effects*. Proceedings of the National Academy of Sciences, 2023. **120**(44): p. e2311219120.
33. Zhou, Z., et al., *Dnabert-2: Efficient foundation model and benchmark for multi-species genome*. arXiv preprint arXiv:2306.15006, 2023.
34. Sanabria, M., et al., *DNA language model GROVER learns sequence context in the human genome*. Nature Machine Intelligence, 2024. **6**(8): p. 911-923.
35. Wang, X., et al., *UNI-RNA: universal pre-trained models revolutionize RNA research*. bioRxiv, 2023: p. 2023.07.11.548588.
36. Chen, K., et al., *Self-supervised learning on millions of primary RNA sequences from 72 vertebrates improves sequence-based RNA splicing prediction*. Briefings in Bioinformatics, 2024. **25**(3): p. bbae163.
37. Yang, Y., et al., *Deciphering 3'UTR Mediated Gene Regulation Using Interpretable Deep Representation Learning*. Advanced Science, 2024. **11**(39): p. 2407013.
38. Chu, Y., et al., *A 5' UTR language model for decoding untranslated regions of mRNA and function predictions*. Nature Machine Intelligence, 2024. **6**(4): p. 449-460.

39. Rao, R., et al., *Evaluating protein transfer learning with TAPE*. Advances in neural information processing systems, 2019. **32**.
40. Rives, A., et al., *Biological structure and function emerge from scaling unsupervised learning to 250 million protein sequences*. Proceedings of the National Academy of Sciences, 2021. **118**(15): p. e2016239118.
41. Ferruz, N., S. Schmidt, and B. Hcker, *ProtGPT2 is a deep unsupervised language model for protein design*. Nature communications, 2022. **13**(1): p. 4348.
42. Brandes, N., et al., *ProteinBERT: a universal deep-learning model of protein sequence and function*. Bioinformatics, 2022. **38**(8): p. 2102-2110.
43. Zhou, H.-Y., et al., *Protein Representation Learning via Knowledge Enhanced Primary Structure Modeling*. bioRxiv, 2023: p. 2023-01.
44. Polishchuk, P.G., T.I. Madzhidov, and A.J.J.o.c.-a.m.d. Varnek, *Estimation of the size of drug-like chemical space based on GDB-17 data*. 2013. **27**: p. 675-679.
45. MOLE-BERT: RETHINKING PRE-TRAINING GRAPH NEURAL NETWORKS FOR MOLECULES.
46. Wang, S., et al., *Smiles-Bert*, in *Proceedings of the 10th ACM International Conference on Bioinformatics, Computational Biology and Health Informatics*. 2019. p. 429-436.
47. Bagal, V., et al., *MolGPT: Molecular Generation Using a Transformer-Decoder Model*. J Chem Inf Model, 2022. **62**(9): p. 2064-2076.
48. SynerGPT: In-Context Learning for Personalized Drug Synergy Prediction and Drug Design.
49. Rosen, Y., et al., *Universal cell embeddings: A foundation model for cell biology*. bioRxiv, 2023: p. 2023.11. 28.568918.
50. Theus, A., et al., *CancerFoundation: A single-cell RNA sequencing foundation model to decipher drug resistance in cancer*. bioRxiv, 2024: p. 2024.11. 01.621087.
51. Yang, X., et al., *GeneCompass: deciphering universal gene regulatory mechanisms with a knowledge-informed cross-species foundation model*. Cell Res, 2024. **34**(12): p. 830-845.
52. Kalfon, J., et al., *scPRINT: pre-training on 50 million cells allows robust gene network predictions*. bioRxiv, 2024: p. 2024.07. 29.605556.
53. Schaar, A., et al., *Nicheformer: a foundation model for single -cell and spatial omics*. 2024. Preprint at bioRxiv, 2024. **4**: p. 589472.
54. Sinden, R.R. and R.D. Wells, *DNA structure, mutations, and human genetic disease*. Current opinion in biotechnology, 1992. **3**(6): p. 612-622.
55. Wittkopp, P.J. and G. Kalay, *Cis-regulatory elements: molecular mechanisms and evolutionary processes underlying divergence*. Nature Reviews Genetics, 2012. **13**(1): p. 59-69.
56. Yella, V.R., A. Kumar, and M. Bansal, *Identification of putative promoters in 48 eukaryotic genomes on the basis of DNA free energy*. Scientific reports, 2018. **8**(1): p. 4520.
57. Le, N.Q.K., et al., *BERT-Promoter: An improved sequence-based predictor of DNA promoter using BERT pre-trained model and SHAP feature selection*. Computational Biology and Chemistry, 2022. **99**: p. 107732.
58. Claringbould, A. and J.B. Zaugg, *Enhancers in disease: molecular basis and emerging treatment strategies*. Trends in Molecular Medicine, 2021. **27**(11): p. 1060-1073.
59. Nasser, J., et al., *Genome-wide enhancer maps link risk variants to disease genes*. Nature, 2021. **593**(7858): p. 238-243.
60. Luo, H., et al. *iEnhancer-BERT: A novel transfer learning architecture based on DNA-Language model for identifying enhancers and their strength*. in *International Conference on Intelligent Computing*. 2022. Springer.
61. Ferraz, R.A.C., et al., *DNA–protein interaction studies: a historical and comparative analysis*. Plant Methods, 2021. **17**(1): p. 1-21.

62. Luo, H., et al., *Improving language model of human genome for DNA–protein binding prediction based on task-specific pre-training*. Interdisciplinary Sciences: Computational Life Sciences, 2023. **15**(1): p. 32-43.
63. An, W., et al. *MoDNA: motif-oriented pre-training for DNA language model*. in *Proceedings of the 13th ACM International Conference on Bioinformatics, Computational Biology and Health Informatics*. 2022.
64. Moore, L.D., T. Le, and G. Fan, *DNA methylation and its basic function*. Neuropsychopharmacology, 2013. **38**(1): p. 23-38.
65. Zhang, L., et al., *Comprehensive analysis of DNA 5-methylcytosine and N6-adenine methylation by nanopore sequencing in hepatocellular carcinoma*. Frontiers in cell and developmental biology, 2022. **10**: p. 827391.
66. Tsukiyama, S., et al., *BERT6mA: prediction of DNA N6-methyladenine site using deep learning-based approaches*. Briefings in Bioinformatics, 2022. **23**(2): p. bbac053.
67. Yu, Y., et al., *iDNA-ABT: advanced deep learning model for detecting DNA methylation with adaptive features and transductive information maximization*. Bioinformatics, 2021. **37**(24): p. 4603-4610.
68. Jin, J., et al., *iDNA-ABF: multi-scale deep biological language learning model for the interpretable prediction of DNA methylations*. Genome biology, 2022. **23**(1): p. 1-23.
69. Zeng, W., A. Gautam, and D.H. Huson, *MuLan-Methyl-Multiple Transformer-based Language Models for Accurate DNA Methylation Prediction*. bioRxiv, 2023: p. 2023.01. 04.522704.
70. Sanh, V., et al., *DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter*. arXiv preprint arXiv:1910.01108, 2019.
71. Lan, Z., et al., *Albert: A lite bert for self-supervised learning of language representations*. arXiv preprint arXiv:1909.11942, 2019.
72. Yang, Z., et al., *Xlnet: Generalized autoregressive pretraining for language understanding*. Advances in neural information processing systems, 2019. **32**.
73. Clark, K., et al., *Electra: Pre-training text encoders as discriminators rather than generators*. arXiv preprint arXiv:2003.10555, 2020.
74. Wilkinson, M.E., C. Charenton, and K. Nagai, *RNA splicing by the spliceosome*. Annual review of biochemistry, 2020. **89**: p. 359-388.
75. Zhang, J., et al., *Advances and opportunities in RNA structure experimental determination and computational modeling*. Nature Methods, 2022. **19**(10): p. 1193-1207.
76. Malbec, L., et al., *Dynamic methylome of internal mRNA N 7-methylguanosine and its regulatory role in translation*. Cell research, 2019. **29**(11): p. 927-941.
77. Feng, H., et al., *LncCat: An ORF attention model to identify lncRNA based on ensemble learning strategy and fused sequence information*. Computational and Structural Biotechnology Journal, 2023. **21**: p. 1433-1447.
78. Xia, S., et al. *A multi-granularity information-enhanced pre-training method for predicting the coding potential of sORFs in plant lncRNAs*. in *2023 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*. 2023. IEEE.
79. Yamada, K. and M. Hamada, *Prediction of RNA–protein interactions using a nucleotide language model*. Bioinformatics Advances, 2022. **2**(1): p. vbac023.
80. Fang, Y., X. Pan, and H.-B. Shen, *Recent deep learning methodology development for RNA–RNA interaction prediction*. Symmetry, 2022. **14**(7): p. 1302.
81. Gibb, E.A., C.J. Brown, and W.L. Lam, *The functional role of long non-coding RNA in human carcinomas*. Molecular cancer, 2011. **10**(1): p. 1-17.

82. Zhang, L., et al., *BERT-m7G: a transformer architecture based on BERT and stacking ensemble to identify RNA N7-Methylguanosine sites from sequence information*. Computational and Mathematical Methods in Medicine, 2021. **2021**.
83. Soylu, N.N. and E. Sefer, *BERT2OME: Prediction of 2'-O-methylation Modifications from RNA Sequence by Transformer Architecture Based on BERT*. IEEE/ACM Transactions on Computational Biology and Bioinformatics, 2023.
84. Pardi, N., et al., *mRNA vaccines—a new era in vaccinology*. Nature reviews Drug discovery, 2018. **17**(4): p. 261-279.
85. Babjac, A.N., Z. Lu, and S.J. Emrich. *CodonBERT: Using BERT for Sentiment Analysis to Better Predict Genes with Low Expression*. in *Proceedings of the 14th ACM International Conference on Bioinformatics, Computational Biology, and Health Informatics*. 2023.
86. Gong, H., et al., *Integrated mRNA sequence optimization using deep learning*. Brief Bioinform, 2023. **24**(1).
87. Ding, W., K. Nakai, and H. Gong, *Protein design via deep learning*. Briefings in bioinformatics, 2022. **23**(3): p. bbac102.
88. Qiu, Y. and G.-W. Wei, *Artificial intelligence-aided protein engineering: from topological data analysis to deep protein language models*. arXiv preprint arXiv:2307.14587, 2023.
89. Kovaltsuk, A., et al., *Observed antibody space: a resource for data mining next-generation sequencing of antibody repertoires*. The Journal of Immunology, 2018. **201**(8): p. 2502-2509.
90. Schauperl, M. and R.A. Denny, *AI-based protein structure prediction in drug discovery: impacts and challenges*. Journal of Chemical Information and Modeling, 2022. **62**(13): p. 3142-3156.
91. David, A., et al., *The AlphaFold database of protein structures: a biologist's guide*. Journal of molecular biology, 2022. **434**(2): p. 167336.
92. Rao, R.M., et al. *MSA transformer*. in *International Conference on Machine Learning*. 2021.
93. Dai, Z., et al., *Transformer-xl: Attentive language models beyond a fixed-length context*. arXiv preprint arXiv:1901.02860, 2019.
94. Raffel, C., et al., *Exploring the limits of transfer learning with a unified text-to-text transformer*. The Journal of Machine Learning Research, 2020. **21**(1): p. 5485-5551.
95. UniProt: the universal protein knowledgebase in 2021. Nucleic acids research, 2021. **49**(D1): p. D480-D489.
96. Steinegger, M. and J. Sding, *Clustering huge protein sequence sets in linear time*. Nature communications, 2018. **9**(1): p. 2542.
97. Strokach, A. and P.M. Kim, *Deep generative modeling for protein design*. Current opinion in structural biology, 2022. **72**: p. 226-236.
98. Ferruz, N. and B. Hcker, *Controllable protein design with language models*. Nature Machine Intelligence, 2022. **4**(6): p. 521-532.
99. Mirdita, M., et al., *ColabFold: making protein folding accessible to all*. Nature methods, 2022. **19**(6): p. 679-682.
100. Jumper, J., et al., *Highly accurate protein structure prediction with AlphaFold*. Nature, 2021. **596**(7873): p. 583-589.
101. Zhou, X., et al., *I-TASSER-MTD: a deep-learning-based platform for multi-domain protein structure and function prediction*. Nature Protocols, 2022. **17**(10): p. 2326-2353.
102. Ferruz, N., et al., *From sequence to function through structure: Deep learning for protein design*. Computational and Structural Biotechnology Journal, 2023. **21**: p. 238-250.
103. Xu, M., et al. *Protst: Multi-modality learning of protein sequences and biomedical texts*. in *International Conference on Machine Learning*. 2023. PMLR.

104. Rosenberg, A.A., A. Marx, and A.M. Bronstein, *Codon-specific Ramachandran plots show amino acid backbone conformation depends on identity of the translated codon*. Nature communications, 2022. **13**(1): p. 2815.
105. Wang, H., et al., *Protein post-translational modifications in the regulation of cancer hallmarks*. Cancer Gene Therapy, 2023. **30**(4): p. 529-547.
106. de Brevern, A.G. and J. Rebehmed, *Current status of PTMs structural databases: applications, limitations and prospects*. Amino Acids, 2022. **54**(4): p. 575-590.
107. Savino, S., T. Desmet, and J. Franceus, *Insertions and deletions in protein evolution and engineering*. Biotechnology Advances, 2022. **60**: p. 108010.
108. Horne, J. and D. Shukla, *Recent advances in machine learning variant effect prediction tools for protein engineering*. Industrial \& engineering chemistry research, 2022. **61**(19): p. 6235-6245.
109. Alley, E.C., et al., *Unified rational protein engineering with sequence-based deep representation learning*. Nature methods, 2019. **16**(12): p. 1315-1322.
110. Liu, W., et al., *PLMSearch: Protein language model powers accurate and fast sequence search for remote homology*. Nature communications, 2024. **15**(1): p. 2775.
111. Hong, L., et al., *Fast, sensitive detection of protein homologs using deep dense retrieval*. Nature Biotechnology, 2024: p. 1-13.
112. Pucci, F., M. Schwersensky, and M. Rooman, *Artificial intelligence challenges for predicting the impact of mutations on protein stability*. Current opinion in structural biology, 2022. **72**: p. 161-168.
113. Wang, Z., et al. *Multi-level Protein Structure Pre-training via Prompt Learning*. in *The Eleventh International Conference on Learning Representations*. 2022.
114. Tang, T., et al., *Machine learning on protein--protein interaction prediction: models, challenges and trends*. Briefings in Bioinformatics, 2023. **24**(2): p. bbad076.
115. Durham, J., et al., *Recent advances in predicting and modeling protein--protein interactions*. Trends in biochemical sciences, 2023.
116. Zhang, N., et al., *Ontoprotein: Protein pretraining with gene ontology embedding*. arXiv preprint arXiv:2201.11147, 2022.
117. Janeway, C., et al., *Immunobiology: the immune system in health and disease*. Vol. 2. 2001: Garland Pub. New York.
118. Peters, B., M. Nielsen, and A.J.A.R.o.I. Sette, *T cell epitope predictions*. 2020. **38**: p. 123-145.
119. O'Donnell, T.J., A. Rubinsteyn, and U.J.C.s. Laserson, *MHCflurry 2.0: improved pan-allele prediction of MHC class I-presented peptides by incorporating antigen processing*. 2020. **11**(1): p. 42-48. e7.
120. Wang, F., et al., *MHCRoBERTa: pan-specific peptide-MHC class I binding prediction through transfer learning with label-agnostic protein sequences*. Brief Bioinform, 2022. **23**(3).
121. Cheng, J., et al., *BERTMHC: improved MHC-peptide class II interaction prediction with transformer and multiple instance learning*. Bioinformatics, 2021. **37**(22): p. 4172-4179.
122. Wu, K., et al., *TCR-BERT: learning the grammar of T-cell receptors for flexible antigenbinding analyses*. 2021.
123. Zhao, Y., et al., *SC-AIR-BERT: a pre-trained single-cell model for predicting the antigen-binding specificity of the adaptive immune receptor*. Brief Bioinform, 2023. **24**(4).
124. Wang, Q., et al., *AntiFormer: graph enhanced large language model for binding affinity prediction*. Briefings in Bioinformatics, 2024. **25**(5).
125. Olsen, T.H., I.H. Moal, and C.M. Deane, *AbLang: an antibody language model for completing antibody sequences*. Bioinformatics Advances, 2022. **2**(1): p. vbac046.
126. Liu, Y., et al., *Roberta: A robustly optimized bert pretraining approach*. arXiv preprint arXiv:1907.11692, 2019.

127. Leem, J., et al., *Deciphering the language of antibodies using self-supervised learning*. Patterns, 2022. **3**(7).
128. Wang, D., F. Ye, and H. Zhou, *On pre-trained language models for antibody*. bioRxiv, 2023: p. 2023-01.
129. Askr, H., et al., *Deep learning in drug discovery: an integrative review and future challenges*. Artificial Intelligence Review, 2023. **56**(7): p. 5975-6037.
130. Xiaobo, Z. and S.T.C. Wong, *High content cellular imaging for drug development*. IEEE Signal Processing Magazine, 2006. **23**(2): p. 170-174.
131. Sun, X., et al., *Multi-scale agent-based brain cancer modeling and prediction of TKI treatment response: incorporating EGFR signaling pathway and angiogenesis*. BMC Bioinformatics, 2012. **13**: p. 218.
132. Vargason, A.M., A.C. Anselmo, and S.J.N.b.e. Mitragotri, *The evolution of commercial drug delivery technologies*. 2021. **5**(9): p. 951-967.
133. Leeson, P.D. and B.J.N.r.D.d. Springthorpe, *The influence of drug-like concepts on decision-making in medicinal chemistry*. 2007. **6**(11): p. 881-890.
134. Ozcelik, R., et al., *Structure-Based Drug Discovery with Deep Learning*. Chembiochem, 2023. **24**(13): p. e202200776.
135. Li, Z., et al., *Deep learning methods for molecular representation and property prediction*. Drug Discovery Today, 2022: p. 103373.
136. Chen, W., et al., *Artificial intelligence for drug discovery: Resources, methods, and applications*. Molecular Therapy-Nucleic Acids, 2023.
137. Chithrananda, S., G. Grand, and B. Ramsundar, *ChemBERTa: large-scale self-supervised pretraining for molecular property prediction*. arXiv preprint arXiv:2010.09885, 2020.
138. *ChemBERTa-2: Towards Chemical Foundation Models*.
139. Xia, J., et al. *Mole-bert: Rethinking pre-training graph neural networks for molecules*. in *The Eleventh International Conference on Learning Representations*. 2022.
140. Bilodeau, C., et al., *Generative models for molecular discovery: Recent advances and challenges*. Wiley Interdisciplinary Reviews: Computational Molecular Science, 2022. **12**(5): p. e1608.
141. Meyers, J., B. Fabian, and N. Brown, *De novo molecular design and generative models*. Drug Discovery Today, 2021. **26**(11): p. 2707-2715.
142. Abbasi, K., et al., *Deep learning in drug target interaction prediction: current and future perspectives*. Current Medicinal Chemistry, 2021. **28**(11): p. 2100-2113.
143. Zhang, Z., et al., *Graph neural network approaches for drug-target interactions*. Current Opinion in Structural Biology, 2022. **73**: p. 102327.
144. Zheng, J., X. Xiao, and W.-R. Qiu, *DTI-BERT: identifying drug-target interactions in cellular networking based on BERT and deep learning method*. Frontiers in Genetics, 2022. **13**: p. 859188.
145. Kalakoti, Y., S. Yadav, and D. Sundar, *TransDTI: transformer-based language models for estimating DTIs and building a drug recommendation workflow*. ACS omega, 2022. **7**(3): p. 2706-2717.
146. Kang, H., et al., *Fine-tuning of bert model to accurately predict drug--target interactions*. Pharmaceutics, 2022. **14**(8): p. 1710.
147. Nguyen, T.M., T. Nguyen, and T. Tran, *Mitigating cold-start problems in drug-target affinity prediction with interaction knowledge transferring*. Briefings in Bioinformatics, 2022. **23**(4): p. bbac269.
148. Ragoza, M., et al., *Protein--ligand scoring with convolutional neural networks*. Journal of chemical information and modeling, 2017. **57**(4): p. 942-957.
149. Li, S., et al. *Structure-aware interactive graph neural networks for the prediction of protein-ligand binding affinity*. in *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery \& Data Mining*. 2021.

150. Jiang, D., et al., *InteractionGraphNet: a novel and efficient deep graph representation learning framework for accurate protein--ligand interaction predictions*. Journal of medicinal chemistry, 2021. **64**(24): p. 18209-18232.
151. Wang, Y., et al., *A point cloud-based deep learning strategy for protein--ligand binding affinity prediction*. Briefings in Bioinformatics, 2022. **23**(1): p. bbab474.
152. Hecht, J.R., et al., *A randomized phase IIIB trial of chemotherapy, bevacizumab, and panitumumab compared with chemotherapy and bevacizumab alone for metastatic colorectal cancer*. 2009. **27**(5): p. 672-680.
153. Tol, J., et al., *Chemotherapy, bevacizumab, and cetuximab in metastatic colorectal cancer*. 2009. **360**(6): p. 563-572.
154. Zhang, W., et al., *DCE-DForest: a deep forest model for the prediction of anticancer drug combination effects*. Computational and Mathematical Methods in Medicine, 2022. **2022**.
155. Xu, M., et al., *DFFNDDs: prediction of synergistic drug combinations with dual feature fusion networks*. Journal of Cheminformatics, 2023. **15**(1): p. 1-12.
156. Shen, H., et al., *A universal approach for integrating super large-scale single-cell transcriptomes by exploring gene rankings*. Brief Bioinform, 2022. **23**(2).
157. Zhao, S., J. Zhang, and Z. Nie, *Large-scale cell representation learning via divide-and-conquer contrastive learning*. arXiv preprint arXiv:2306.04371, 2023.
158. Roohani, Y., K. Huang, and J. Leskovec, *GEARS: Predicting transcriptional outcomes of novel multi-gene perturbations*. BioRxiv, 2022: p. 2022.07. 12.499735.
159. Li, G., et al., *A deep generative model for multi-view profiling of single-cell RNA-seq and ATAC-seq data*. Genome Biol, 2022. **23**(1): p. 20.
160. Ma, A., et al., *Single-cell biological network inference using a heterogeneous graph transformer*. Nat Commun, 2023. **14**(1): p. 964.
161. Linjing, L., et al., *A pre-trained large language model for translating single-cell transcriptome to proteome*. bioRxiv, 2023: p. 2023.07.04.547619.
162. Tang, W., et al. *Single-cell multimodal prediction via transformers*. in *Proceedings of the 32nd ACM International Conference on Information and Knowledge Management*. 2023.
163. Szklarczyk, D., et al., *The STRING database in 2023: protein-protein association networks and functional enrichment analyses for any sequenced genome of interest*. Nucleic Acids Res, 2023. **51**(D1): p. D638-D646.
164. Wen, H., et al., *Single cells are spatial tokens: Transformers for spatial transcriptomic data imputation*. arXiv preprint arXiv:2302.03038, 2023.
165. Hou, W. and Z. Ji, *Assessing GPT-4 for cell type annotation in single-cell RNA-seq analysis*. Nat Methods, 2024. **21**(8): p. 1462-1465.
166. Chen, Y. and J. Zou, *GenePT: a simple but effective foundation model for genes and cells built from ChatGPT*. bioRxiv, 2024: p. 2023.10. 16.562533.
167. Levine, D., et al., *Cell2Sentence: teaching large language models the language of biology*. BioRxiv, 2023: p. 2023.09. 11.557287.
168. Zhao, S., et al., *Langcell: Language-cell pre-training for cell identity understanding*. arXiv preprint arXiv:2405.06708, 2024.
169. Lu, Y.-C., et al., *scChat: A Large Language Model-Powered Co-Pilot for Contextualized Single-Cell RNA Sequencing Analysis*. bioRxiv, 2024: p. 2024.10. 01.616063.
170. Liu, T., et al., *scelmo: Embeddings from language models are good learners for single-cell data analysis*. bioRxiv, 2023: p. 2023.12. 07.569910.
171. Heimberg, G., et al., *Scalable querying of human cell atlases via a foundational model reveals commonalities across fibrosis-associated macrophages*. BioRxiv, 2023: p. 2023.07. 18.549537.

Main figures

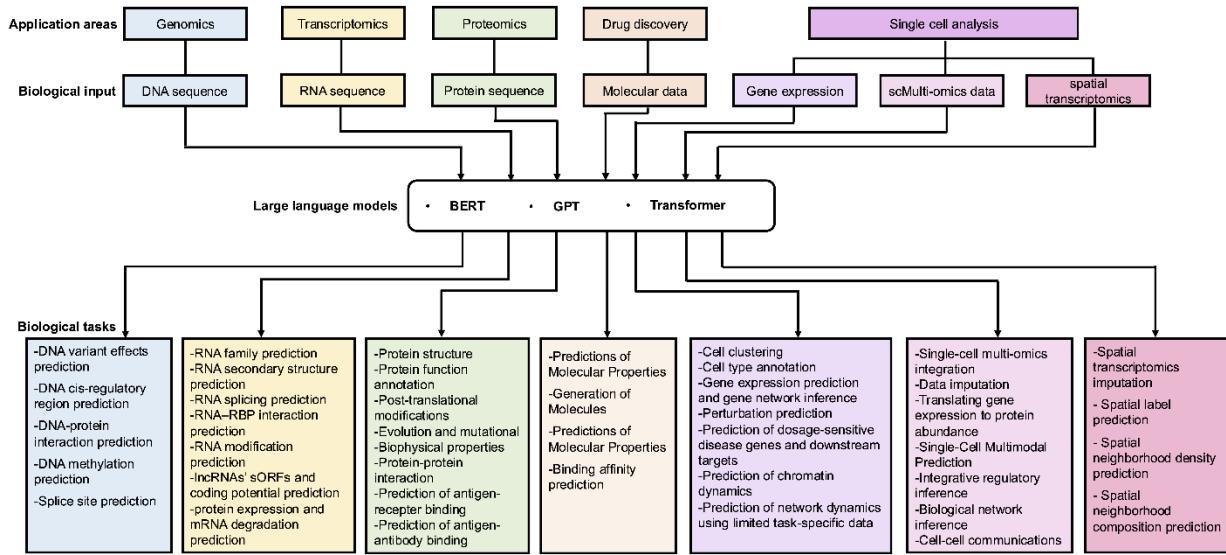
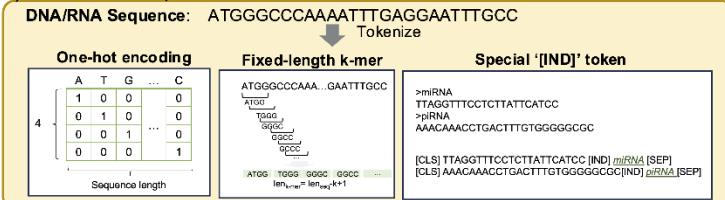


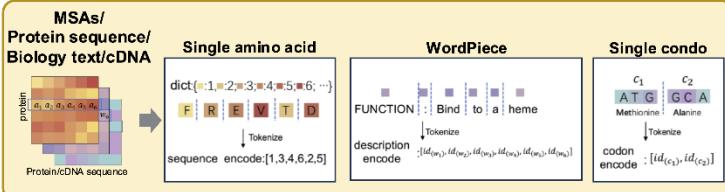
Figure 1. Summary of the application of large language models in bioinformatics in this review. Applications of large language models in bioinformatics include applications in genomics, transcriptomics, proteomics, drug discovery and single-cell analysis. Applications of LLMs in genomics focus on LLMs using DNA sequence; applications of LLMs in transcriptomics focus on using RNA sequence; applications of LLMs in proteomics focus on LLMs using protein sequence; applications of LLMs in drug discovery focus on LLMs using molecular data and applications of LLMs in single-cell analysis focus on LLMs using scRNA-seq, scMulti-omics and spatial transcriptomics data. Each corresponds to a variety of biological downstream tasks.

a Tokenization

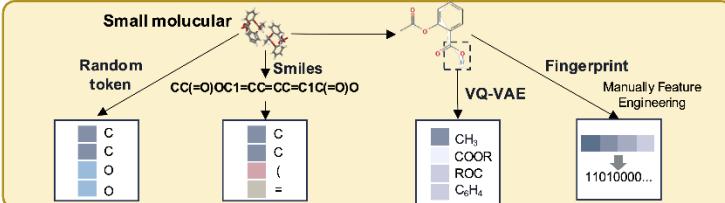
i) DNA/RNA sequence



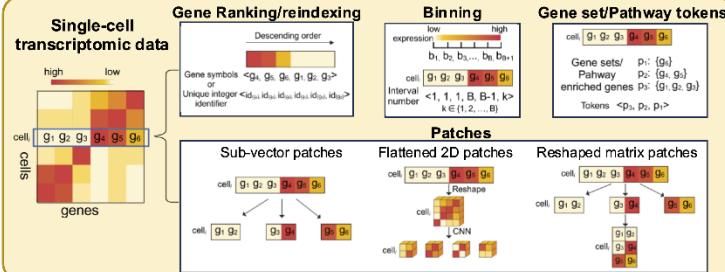
ii) Protein



iii) Small molecular

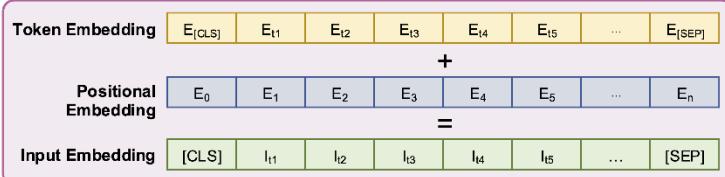


iv) Single-cell profiles

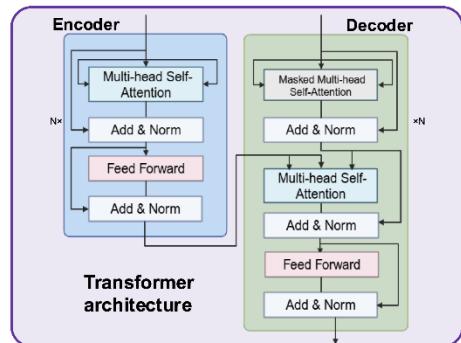


b Input Embedding

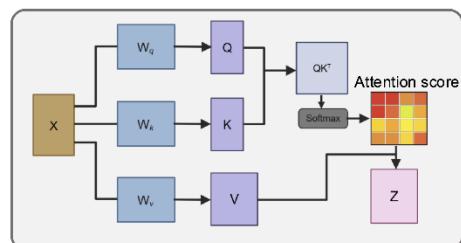
Feed to Embedding layer



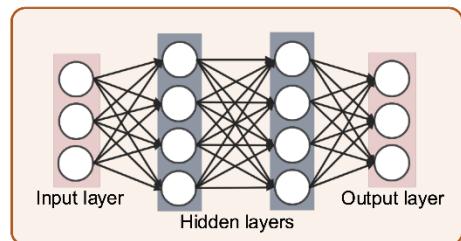
c Transformer architecture



d Attention mechanism



e Feed-forward neural network



f LLM pre-training process

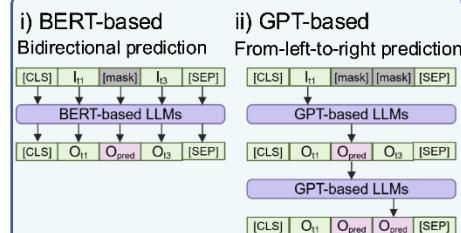


Figure 2. Building blocks of large language models in bioinformatics. **a**, tokenization methods tailored to various data types, including DNA/RNA sequences, proteins, small molecules, and single-cell data. **b**, input embedding strategies used in large language models to encode tokenized data. **c**, schematic representation of the transformer architecture, a foundational structure in LLMs. **d**, the attention mechanism, enabling models to focus on important features in sequences. **e**, the feed-forward network, a critical component of transformers for learning hierarchical

representations. **f**, pre-training processes for BERT and GPT-based models, highlighting BERT's bidirectional prediction approach and GPT's left-to-right prediction strategy.

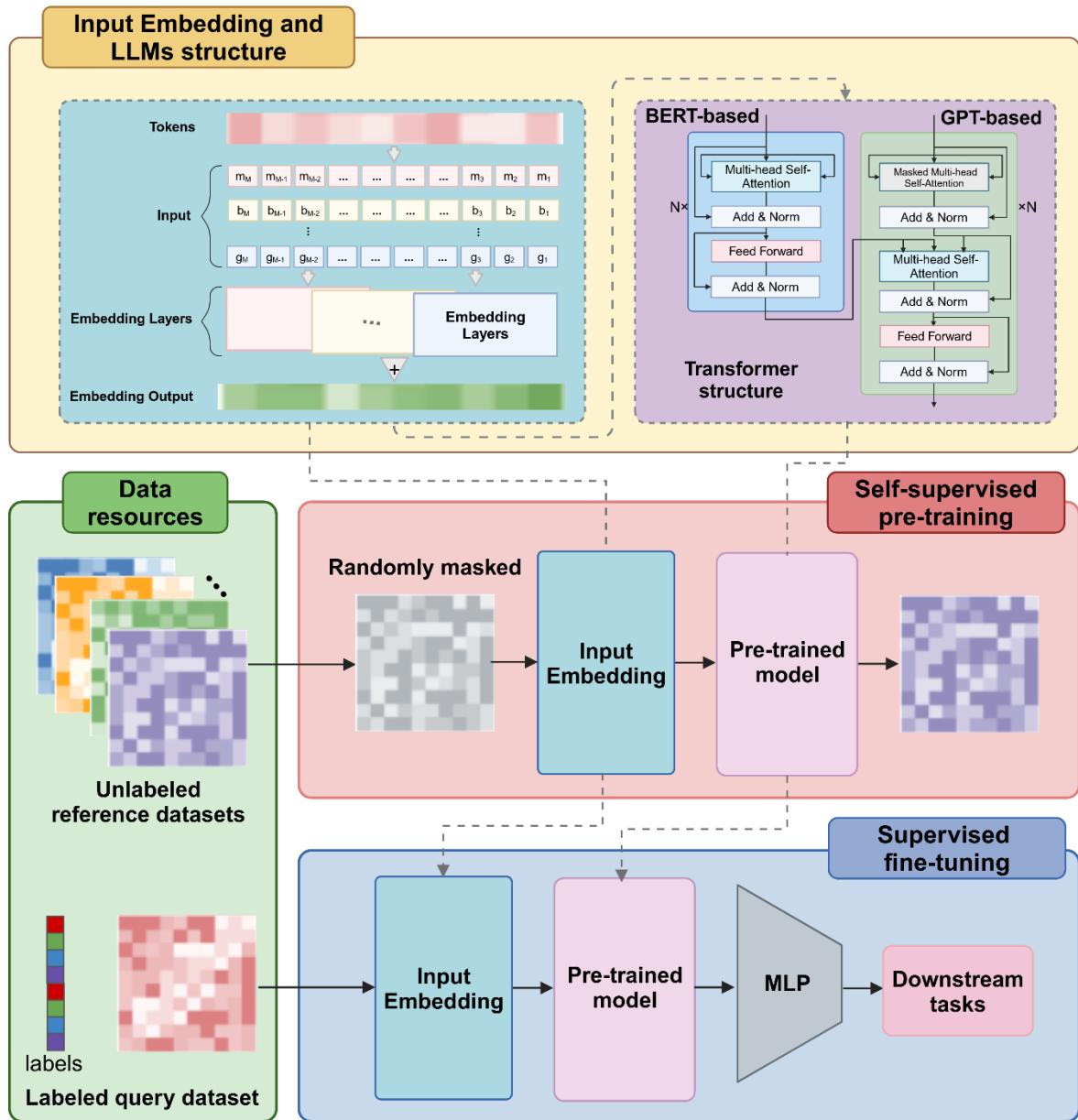
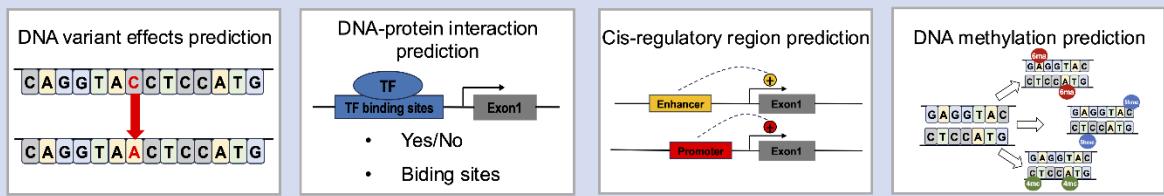


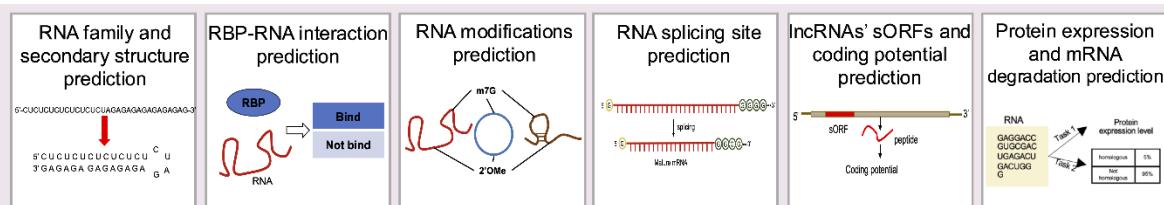
Figure 3. Schematic diagram of the large language model pretraining and fine-tuning process.

The workflow begins with tokenizing the input data, which is then fed into the embedding layer and transformer models. The training process comprises two stages: pretraining and fine-tuning. Pretraining employs self-supervised learning on large-scale, unlabeled reference datasets to develop a general-purpose model with robust generalization capabilities. Fine-tuning builds upon the pretrained model, involving task-specific training to optimize performance for designated applications.

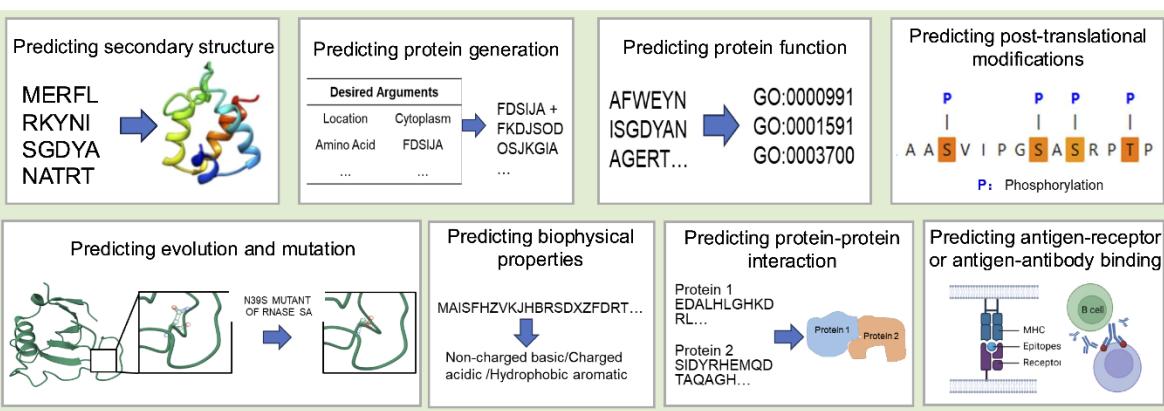
Downstream tasks in genomics



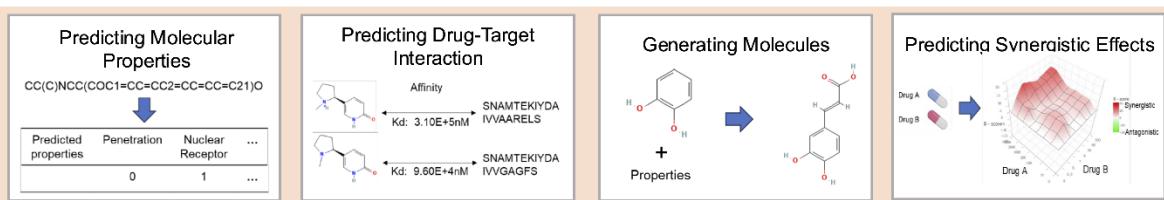
Downstream tasks in transcriptomics



Downstream tasks in proteomics



Downstream tasks in drug discovery



Downstream tasks in single-cell analysis

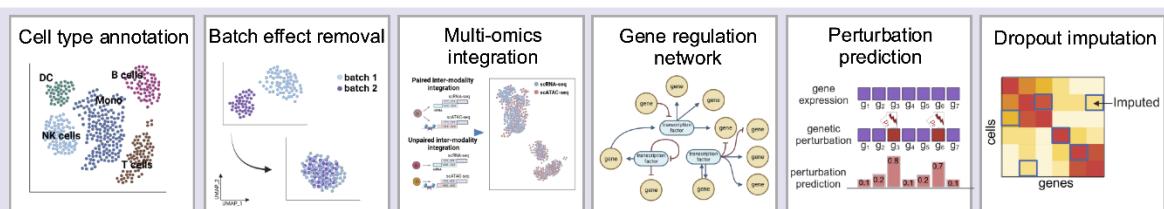


Figure 4. Downstream tasks of large language models in bioinformatics. Large language models (LLMs) have seen numerous successful applications in bioinformatics, addressing a wide array of tasks across DNA, RNA, protein, drug discovery, and single-cell analysis.

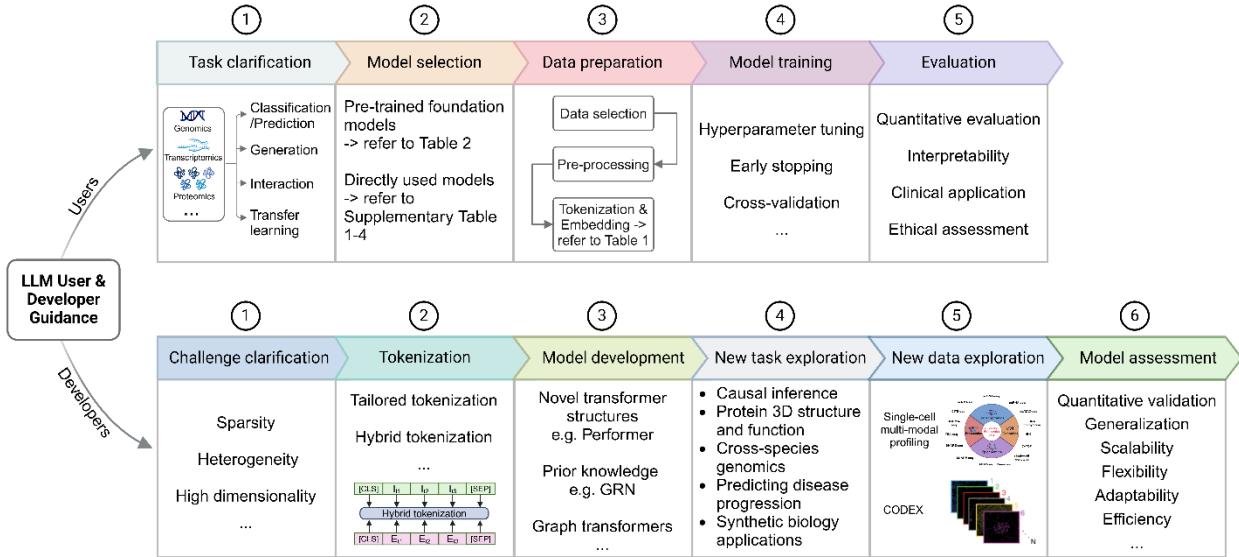


Figure 5. Guidance for LLM users and developers on how to use and develop LLM in practice. Guidance for LLM users includes steps such as clarifying the task, selecting an appropriate model, preparing the dataset, training the model, and evaluating its performance. For LLM developers, the focus involves identifying domain-specific challenges, designing tokenization strategies, advancing model architectures, exploring novel tasks and data types, and assessing model capabilities comprehensively.

Tables

Table 1. Tokenization methods for different types of data

Application area	Data type	Method	Example
Genomics/Transcriptomics	DNA/RNA sequence	One-hot encoding	RNA-FM, RNA-MSM
		Fixed-length k-mers	DNABERT, Nucleotide Transformer, DNABERT-2, DNAGPT, RNABERT
		Special '[IND]' token	RNAErnie
Proteomics	MSAs/Protein sequences	Single Amino Acid Tokenization	MSA Transformer/TAPE, ESM-1b, ProtTrans, Progen
	Biomedical text	WordPiece	ProtST
	cDNA	Single condo Tokenization	CaLM
Drug discovery	Simplified Molecular-Input Line-Entry system (SMILES)	Random token	K-BERT
		SmilesTokenizer	ChemBERTa, ChemBERTa-2, MolGPT
		Graph VQ-VAE	Mole-BERT
		fingerprint	SMILES-BERT
Single-cell analysis	Expression profiles	Gene expression Ranking	Geneformer, tGPT, iSEEEK
		Binning	scBERT, scGPT, scFormer, CellLM, BioFormers, CancerFoundation
		Gene set/Pathway tokens	TOSICA
		Patches	CIForm, scTranSort, scCLIP
		Gene value projection	scTranslator, scFounfation, scMulan, scGREAT
		Cell tokens	CellPLM, ScRAT, mcBERT

Table 2. Foundation models in bioinformatics

Application area	Model	Architecture	Pre-training Data	Code available
Genomics	GPN	Transformer-based	Reference genomes from 8 species	https://github.com/songlab-cal/gpn
	Nucleotide Transformer	Transformer-based	3.2 billion nucleotides in GRCh38/hg38 reference assembly, 20.5 trillion nucleotides including 125 million mutations (111 million SNPs, 14 million indels), and 174 billion nucleotides from 850 species	https://github.com/instadeepai/nucleotide-transformer
	DNABERT	BERT-based	2.75 billion nucleotide based human genome dataset	https://github.com/jerryji1993/DNABERT
	DNABERT-2	BERT-based	2.75 billion nucleotide based human genome dataset and 32.49 billion nucleotide bases from 135 species, spread across 6 categories	https://github.com/MAGICSLAB/DNABERT_2
	MoDNA	BERT-based	Same as Nucleotide Transformer	https://github.com/uta-smile/MoDNA
	GROVER	BERT-based	Homo sapiens (human) genome assembly GRCh37 (hg19)	https://github.com/rowanz/grover
	MuLan-Methyl	BERT-based	3 main types of DNA methylation sites (6mA, 4mC, and 5hmC) across 12 genomes, in total 250,599 positive samples	https://github.com/husonlab/mulan-methyl
	iDNA-ABF	BERT-based	Same as MuLan-Methyl	https://github.com/FakeEnd/iDNA_ABF
	iDNA-ABT	BERT-based	Same as MuLan-Methyl	https://github.com/YUYING07/iDNA_ABТ
Transcriptomics	DNAGPT	GPT-based	Reference genomes from the Ensembl database include 3 billion bps, with a total of 1,594,129,992 bps across 9 species	https://github.com/TencentAILabHealthcare/eDNAGPT
	RNABERT	BERT-based	76,237 human-derived small ncRNAs from RNACentral	https://github.com/mana438/RNABERT
	RNA-FM	BERT-based	About 27 million ncRNA sequences across 47 different databases	https://github.com/ml4bio/RNA-FM
	RNA-MSM	BERT-based	4069 RNA families from rfam	https://github.com/yikunpu/RNA-MSM
	SpliceBERT	BERT-based	2 million sequences and approximately covering 65 billion nucleotides of 72 vertebrates from UCSC genome browser	https://github.com/biomed-AI/SpliceBERT
	UNI-RNA	BERT-based	23 million ncRNA sequences obtained from the RNACentral database	https://github.com/ComDec/unirna-tools
	3UTRBERT	BERT-based	108,573 unique mRNA transcripts from the GENCODE and each contains 3,754 nucleotides (median 3048 nts) on average.	https://github.com/yangyn533/3UTRBERT
	UTR-LM	BERT-based	214,349 unlabeled 5' UTR sequences from Ensembl across 5 species	https://github.com/a96123155/UTR-LM
	RNAErnie	Transformer-based	23 million ncRNA sequences obtained from the RNACentral database	https://github.com/CatIIIILIIII/RAErnie
Proteomics	TAPE	Transformer-based	31 million protein sequences from Pfam	https://github.com/songlab-cal/tape
	ESM-1b	Transformer-based	250 million protein sequences from UniRef50	https://github.com/facebookresearch/esm
	ProtTrans	Transformer-XL, XLNet, BERT, Albert, Electra, T5	About 2.3 billion protein sequences from UniRef and BFD	https://github.com/agemagician/ProtTrans
	ProtGPT2	GPT-based	50 million protein sequences from UniRef50	https://huggingface.co/docs/transformers/main_classes/trainer
	ProteinBERT	BERT-based	106 million protein sequences with GO annotations from UniRef50	https://github.com/nadavbra/protein_bert
	KeAP	BERT-based	5 million Triplet in the format of (Protein, Relation, Attribute) with nearly 600k protein, 50k attribute terms, and 31 relation terms included	https://github.com/RL4M/KeAP

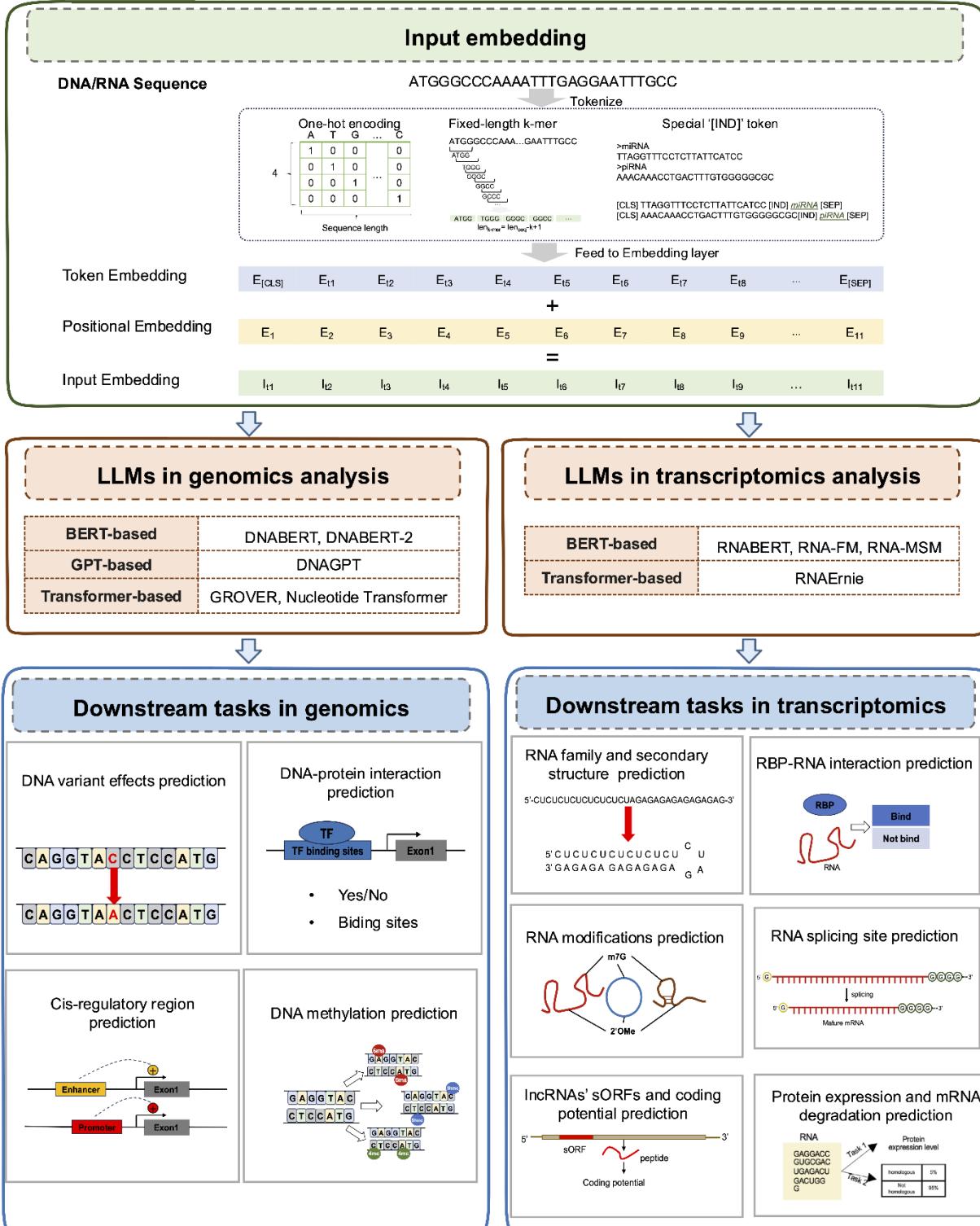
	CaLM	Transformer-based	9,858,385 cDNA sequences of seven model organisms	https://github.com/oxpig/CalM
Drug discovery	SMILES-BERT	BERT-based	Two datasets from NCATS (NIH) and 128 datasets from PubChem	https://github.com/uta-smile/SMILES-BERT
	ChemBERTa	BERT-based	77 million unique SMILES	https://github.com/seyonechithrananda/bert-loves-chemistry
	K-BERT	BERT-based	Book review dataset contains 20,000 positive and 20,000 negative reviews collected from Douban	https://github.com/autoliuweijie/K-BERT
	Mole-BERT	BERT-based	2 million molecules	https://github.com/junxia97/Mole-BERT
	MolGPT	GPT-based	Datasets from MOSES and GuacaMol	https://github.com/devalab/molgpt
	ProtBERT	BERT-based	Datasets from UniRef50, UniRef100 and BFD	https://github.com/agemagician/ProtTrans/
	DeepDDS	BERT-based	Datasets from NCI-ALMANAC	https://github.com/sorachel/DFFNDDSS
	SynerGPT	GPT-based	Datasets from DrugCombDB	Code will be made available upon publication
Single-cell analysis	scBERT	BERT-based	1,126,580 cells from 209 datasets across 74 tissues and 451,513 cells from four sequencing platforms	https://github.com/TencentAILabHealthcare/scBERT
	scGPT	GPT-based	33 million human cells from the CellXGene collection	https://github.com/bowang-lab/scGPT
	Geneformer	BERT-based	29.9 million human single-cell transcriptomes	https://huggingface.co/ctheodoris/Genefomer
	scFoundation	BERT-based	About 50 million human single-cell transcriptomic profiles	https://github.com/biomap-research/scFoundation
	tGPT	GPT-based	22.3 million single-cell transcriptomes	https://github.com/deeplearningplus/tGPT
	GeneCompass	BERT-based	over 120 million single-cell transcriptomes from humans and mice	https://github.com/xCompass-AI/GeneCompass
	scMulan	GPT-based	More than 10 million manually annotated single-cell RNA-seq data	https://github.com/SuperBianC/scMulan
	UCE	BERT-based	300 datasets from the CellXGene corpus includes over 36 million cells, 1,000+ cell types, dozens of tissues, and eight species	https://github.com/snap-stanford/uce
	scPRINT	BERT-based	More than 50M cells from theCellXGene database	https://github.com/cantinilab/scPRINT
	CancerFoundation	BERT-based	50 million cells with roughly a quarter being tumor cells	https://github.com/BoevaLab/CancerFoundation
	Nicheformer	BERT-based	57 million dissociated and 53 million spatially resolved cells across 73 tissues from both human and mouse	https://github.com/theislab/nicheformer

Table 3. Large language models for downstream tasks in bioinformatics

Input data	Biological tasks	Models
DNA sequence	Genome-wide variant effects prediction	DNABERT, DNABERT-2, GPN, Nucleotide Transformer
	DNA cis-regulatory regions prediction	DNABERT, DNABERT-2, BERT-Promoter, iEnhancer-BERT, Nucleotide Transformer
	DNA-protein interaction prediction	DNABERT, DNABERT-2, TFBert, GROVER, and MoDNA
	DNA methylation (6mA,4mC 5hmC) prediction RNA splice sites prediction from DNA sequence	BERT6mA, iDNA-ABF, iDNA-ABT, and MuLan-Methyl DNABERT, DNABERT-2
RNA sequence	RNA 2D/3D structure prediction	RNA-FM, RNA-MSM, and RNA-FM
	RNA structural alignment, RNA family clustering	RNABERT
	RNA splice sites prediction from RNA sequence	SpliceBERT
	RNA N7-Methylguanosine modification prediction	BERT-m7G
	RNA 2'-O-methylation Modifications prediction	Bert2Ome
	Multiple types of RNA modifications prediction	Rm-LR
	Predicting the association between miRNA, lncRNA and disease	BertNDA
	Identifying lncRNAs	LncCat
	Protein expression and mRNA degradation prediction	CodonBERT
Protein sequences	Secondary structure and contact prediction	MSA Transformer, ProtTrans, SPRoBERTa, TAPE, KeAP
	Protein sequence generation	ProGen, ProtGPT2
	Protein function prediction	SPRoBERTa, ProtST, PromptProtein, CaLM
	Major PTMs prediction	ProteinBERT
Gene ontology annotations	Evolution and mutation prediction	SPRoBERTa, UniRep, ESM-1b, TAPE, PLMsearch, DHR
	Biophysical properties prediction	TAPE, PromptProtein
Triplets of protein-relation-attribute	Protein-protein interaction and binding affinity prediction	KeAP
	Antigen-Receptor binding prediction	MHCRoBERTa, BERTMHC, TCR-BERT, SC-AIR-BERT, Antiformer
	Antigen-Antibody binding prediction	AbLang, AntiBERTa, EATLM
Molecular SMILES	Predicting Molecular Properties	SMILES-BERT, ChemBERTa, K-BERT
	Generating Molecules	MolGPT
Molecular graphs	Predicting Molecular Properties	MOLE-BERT
	Predicting Drug-Target Interaction	TransDTI, FG-BERT
Molecular fingerprints and protein sequences		

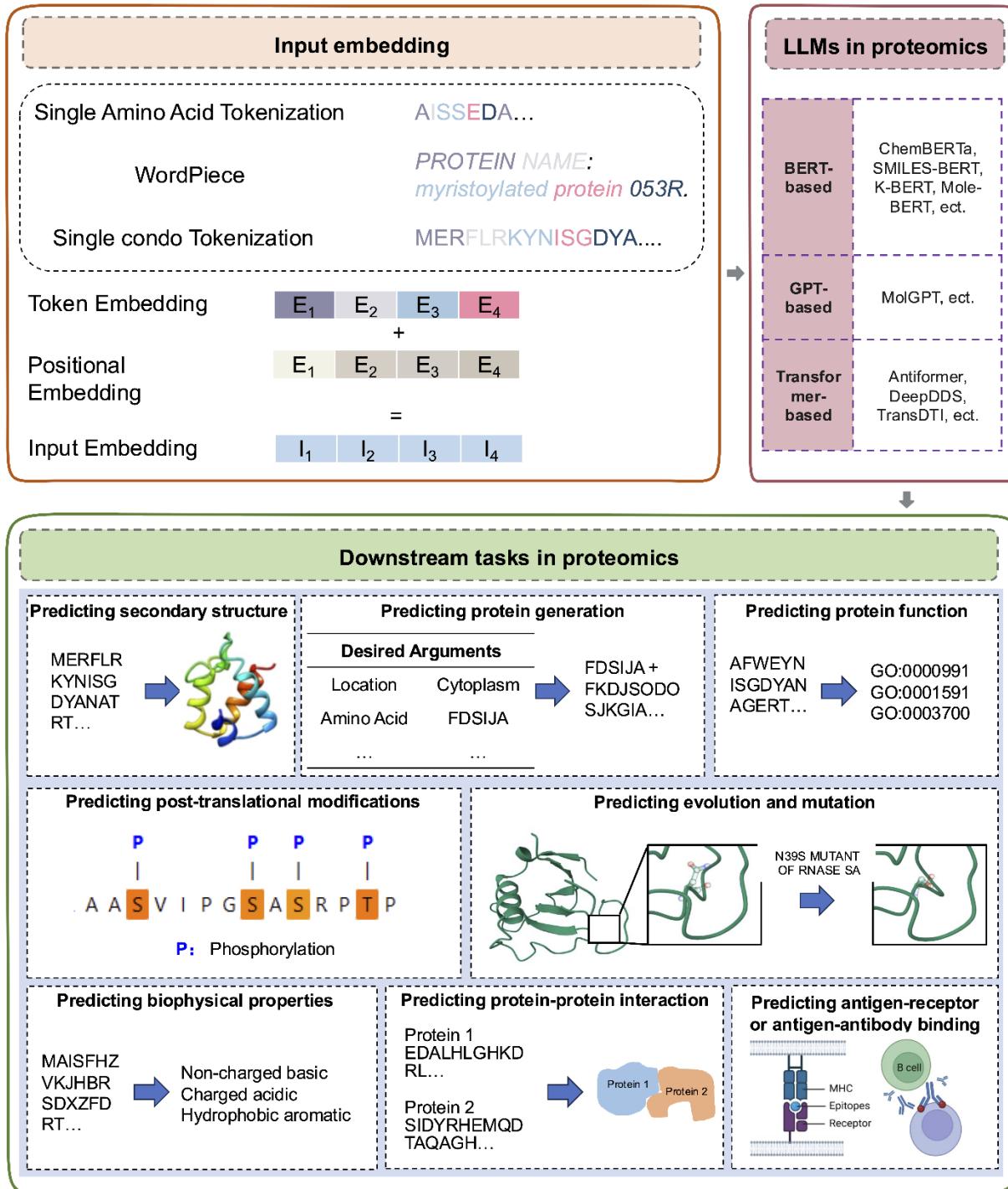
Molecular SMILES and protein sequences	Predicting Synergistic Effects	SynerGPT, C2P2
scRNA-seq data	Cell clustering	tGPT, scFoundation, UCE, iSEEK, CellPLM, BioFormers, mcBERT
	Cell type annotation	scBERT, scGPT, CIFoRM, TOSICA, scTransSort, TransCluster, Geneformer, GeneCompass, scMulan, CellLM, CellPLM, scPRINT
	New cell type identification	scBERT, TOSICA, UCE
	Batch effect removal	scBERT, scGPT, CIFoRM, TOSICA, Geneformer, scMulan, iSEEK, scPRINT, CancerFoundation, mcBERT
	Trajectory inference/Pseudotime analysis	tGPT, scMVP, iSEEK
	Drug response/sensitivity prediction	scFoundation, CellLM, CancerFoundation
	Gene network inference	scGPT, Geneformer, GeneCompass, iSEEK, scGREAT, BioFormers, scPRINT
	Gene perturbation prediction	scGPT, scFoundation, GeneCompass, CellPLM, BioFormers
	Gene expression prediction	scGPT, scMVP, scFoundation, GeneCompass, CellPLM, BioFormers
	cis-regulatory element identification	scMVP
scMuti-omics data	Drug dose-response prediction, Gene dosage sensitivity prediction	GeneCompass
	Single-cell multi-omics integration	scGPT, scMVP, DeepMAPS, scCLIP
	Biological network inference	DeepMAPS
	Cell-cell communications	
	Translating gene expression to protein abundance	scTranslator, scMoFormer
	Single-cell multimodal prediction	scMoFormer
Single-cell spatial transcriptomics data	Integrative regulatory inference	scTranslator
	Spatial transcriptomics imputation	CellPLM, Nicheformer, SpaFormer
	Spatial label prediction	
	Spatial neighborhood density prediction	Nicheformer
	Spatial neighborhood composition prediction	

Supplementary figures



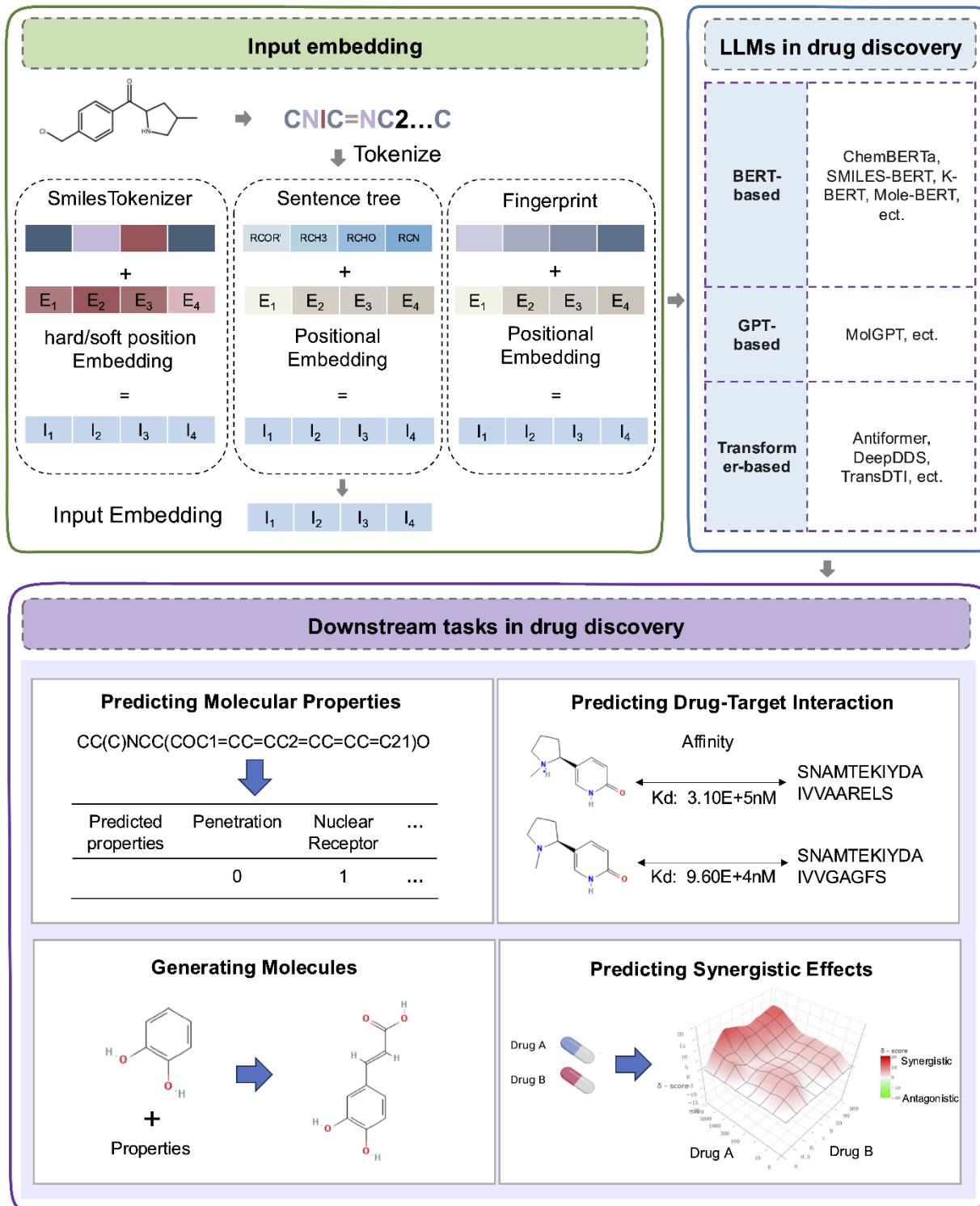
Supplementary figure 1. Applications of large language models in genomics and transcriptomics.

GPT models to solve multiple biological tasks, including genome-wide variant effects prediction, DNA cis-regulatory regions prediction, DNA-protein interaction prediction, DNA methylation (6mA,4mC 5hmC) prediction, splice sites prediction from DNA sequence. The RNA language models take RNA sequences as input, use transformer, BERT, GPT models to solve multiple biological tasks, including RNA 2D/3D structure prediction, RNA structural alignment,, RNA family clustering, RNA splice sites prediction from RNA sequence, RNA N7-methylguanosine modification prediction, RNA 2'-O-methylation modifications prediction, multiple types of RNA modifications prediction, predicting the association between miRNA, lncRNA and disease, identifying lncRNAs, lncRNAs' coding potential prediction, protein expression and mRNA degradation prediction.



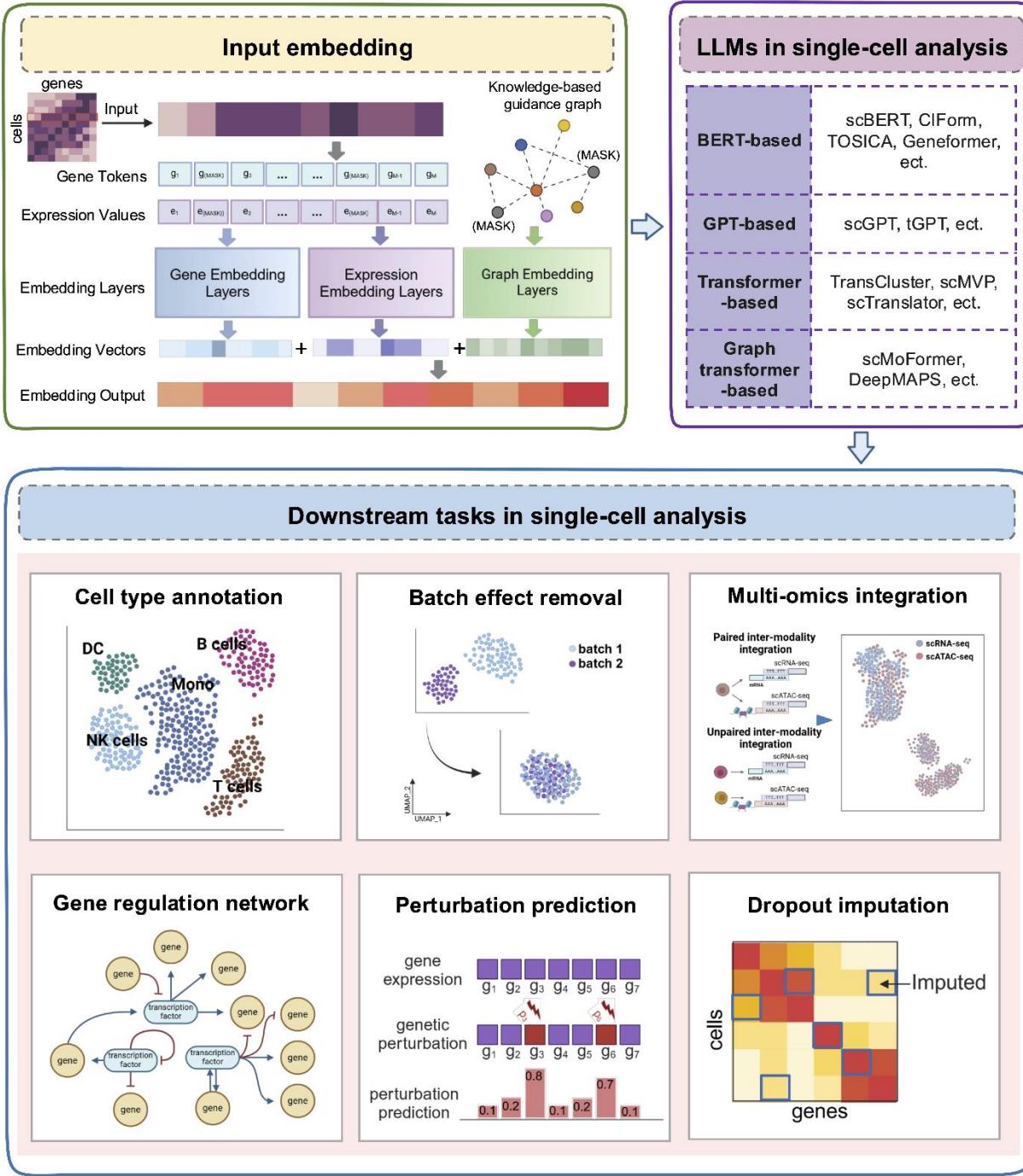
Supplementary figure 2. Applications of large language models in proteomics. The protein language models take multiple sequence alignment, protein sequence, gene ontology and protein-relation-attribute as input, use transformer, BERT, GPT models to solve multiple biological tasks, including predicting secondary structure, predicting protein generation, predicting protein function,

predicting post-translational modifications, predicting evolution and mutation, predicting biophysical properties, predicting protein-protein interaction and predicting antigen-receptor or antigen-antibody binding.



Supplementary figure 3. Applications of large language models in drug discovery. The language models for drug discovery take molecular SMILES, protein sequence, molecular fingerprints and molecular graphs as input, use transformer, BERT, GPT models to solve multiple

biological tasks, including predicting molecular properties, predicting drug-target interaction, generating molecules and predicting synergistic effects.



Supplementary figure 4. Applications of large language models in single-cell analysis. The single-cell language models take gene expression or single-cell multi-omics data as input, use transformer, BERT, GPT models to solve multiple biological tasks, including cell type annotation, batch effect removal, multi-omics integration, gene regulation network inference, perturbation prediction, dropout imputation.

Supplementary Tables

Supplementary Table 1. Detailed information of large language models for genomic and transcriptomic tasks

Application area	Models	Ref	Publication time	Parameters	Architecture	Fine-tuning datasets			Downstream tasks
						Data type	Source	Size	
DNA sequence language model	DNABERT	[1]	Aug 2021	12 transformer layers with 768 hidden units and 12 attention heads in each layer	BERT-based	DNA sequence	Human TATA and non-TATA promoters of 10 000 bp length[2] and ChIP-seq dataset [3]	3,065 human TATA and 26,533 non-TATA promoter-containing sequences and 690 ChIP-seq dataset covers 161 transcription factor binding profiles in 91 human cell lines	Transcription factor binding sites prediction
							DNA sequence	-	Motif analysis
							Assembly GRCh38 FASTA file [4]	10,000 donor, acceptor, and non-splice site sequences	Splice donor and acceptor sites prediction
							dbSNP release 153 [5]	700 million short genetic variants	Identifying effects of genetic variants
	DNABERT-2	[6]	July 2023	batch size is 32, warmup step is 50, and weight decay is 0.01	BERT-based	DNA sequence	TATA and non-TATA promoters downloaded from Eukaryotic Promoter Database (EPDnew) [2]	3,065 human TATA and 26,533 non-TATA promoter-containing sequences	Promoter detection and core promoter detection
							ChIP-seq datasets [7]	161 TF binding profiles in 91 human cell lines(human) and 78 mouse ENCODE ChIP-seq data	Transcription factor binding site prediction
							Ensembl GRCh38 human reference genome [4]	10,000 splice donors, acceptors, and non-splice site sequences.	Splice site prediction
							Histone modification (Yeast)	H3, H3K14ac, H3K36me3, H3K4me1, H3K4me2, H3K4me3, H3K79me3, H3K9ac, H4, H4ac	Epigenetic marks prediction
							SARS_CoV_2 variants [8]	9 types of SARS_CoV_2 variants, including <i>Alpha, Beta, Delta, Eta, Gamma, Iota, Kappa, Lambda and Zeta</i> .	Covid variant prediction
	Nucleotide	[9]	Jan 2023	2 hidden layers	Transfor	DNA	Annotated DNA sequence	90,000 sequences annotated by	Detect known genomic

	Transformer				mer-based	sequence	[10]	Ensembl (“5’ UTR”, “3’ UTR”, “exon”, “intron”, “enhancer”, “promoter”, “CTCF binding site”, “open chromatin”, and “transcription factor binding sites”.)	elements
	DNAGPT	[13]	July 2023	12 layers of transformer blocks based on unidirectional attention, with each layer containing 12 attention heads and a hidden layer size of 768	GPT-based	DNA sequence	DNA sequence with SNP [11]	Independent dataset of genetically diverse human genomes, originating from 7 different meta-populations	Detect human genetic variation
							1000 Genomes Project SNPs [12]	chromosome 22 sequence with 17 variant categories (e.g. stop gained, missense, intergenic)	Predict the impact of mutations
							DNA sequence from DeepGSR [14]	20,933, 18,693, 12,082, and 27,203 true polyadenylation signals data; and 28,244, 25,205, 17,558, and 30,283 true translation initiation sites for human, mouse, bovine, and fruit fly, respectively which are used as ground-truth, non-genomic signals and regions sequences from the genome sequences and combined them with the true cases	Genomic signals and regions prediction
	GROVER	[16]	July 2023	12 transformer layers, 5,000 embeddings	BERT-based	DNA sequence	DNA sequence from Xpresso [15]	18,377 and 21,856 promoters as well as the mRNA half-lives in human and mouse respectively and held out 1000 cases in each specie	mRNA expression level prediction
							CTCF ChIP-seq data [17]	~85,000 binding motifs, only ~32,000 are indeed bound by CTCF	Protein-DNA binding prediction
	GPN	[18]	BioRxiv posted April 2023	25 convolutional blocks with a feed-forward layer, 512embedding sizes of the pre-trained foundation model	BERT-based	DNA sequence	DNA sequence	-	DNA motifs predictions
							1001 Genomes Project [12]	10 million SNPs	Variant effect prediction

	BERT-Promoter	[19]	Aug 2022	BERT model included 12 layers, 768-hidden, 12 heads, and 110,000,000 parameters	BERT-based	DNA sequence	ChIP-chip data, gSELEX peaks, ChIP-exo plus RNA-seq [20, 21]	3382 promoters (1591 strong promoter samples and 1791 weak promoter samples) and 3382 non-promoters	DNA promoter prediction
	TFBert	[22]	Mar 2023	12-layer encoder	BERT-based	DNA sequence	ChIP-seq datasets [7]	690 ChIP-seq dataset contains a training set (80%) and a corresponding test set (20%)	DNA–protein binding sites prediction
MoDNA	[23]	Aug 2022	-	BERT-based	DNA sequence	Same experiment data with DNABERT [2, 3]	3,065 human TATA and 26,533 non-TATA promoter-containing sequences and 690 ChIP-seq dataset covers 161 transcription factor binding profiles in 91 human cell lines	Promoter Prediction	
						CHIP-Seq datasets [3]	690 CHIP-Seq datasets of uniform TFBS contains 161 TFs covering 91 human cell types	Transcription Factor Binding Sites Prediction	
iEnhancer-BERT	[24]	Aug 2022	12-layer transformer architecture with simple fine-tuning	BERT-based	DNA sequence	15 chromatin states of 9 cell types [25]	2968 samples including 1484 non-enhancers, 742 strong enhancers and 742 weak enhancers	Identifying Enhancers and Their Strength	
BERT6mA	[26]	Mar 2022	The hidden size of the LSTM unit is set to 128.DNA sequence embedding with Word2vec.	BERT-based	DNA sequence	Nuclei purification, MNase-seq and ChIP-seq [27-31]	6mA and non6mA data in 11 species including <i>Arabidopsis thaliana</i> (31873 6mAs and non-6mAs), <i>Caenorhabditis elegans</i> (79616 mAs and non-6mAs), <i>Casuarina equisetifolia</i> (6066 6mAs and non-6mAs), <i>Drosophila melanogaster</i> (11 191 6mAs and non-6mAs), <i>Fragaria vesca</i> (3102 6mAs and non-6mAs), <i>H. sapiens</i> (18 335 6mAs and non-6mAs), <i>Rosa chinensis</i> (599 6mAs and non-6mAs), <i>Saccharomyces cerevisiae</i> (37866mAs and non-6mAs), <i>Thermus thermophilus</i> (107 600 6mAs and non-6mAs), <i>Ts. SUP5-1</i> (3379	DNA N6-methyladenine site prediction	

							6mAs and non6mAs) and Xoc. BLS256 (17 215 6mAs and non-6mAs)	
RNA sequence language model	iDNA-ABF	[32]	Oct 2022	12 transformer layers with 768 hidden units and 12 attention heads in each layer	BERT-based	DNA sequence	ChIP-seq data, ATAC-seq data, and histone modifications (HM) data of three human cell lines [33, 34], and DNA methylation dataset from the iDNA-MS [35]	3 main types of DNA methylation sites (6mA, 4mC, and 5hmC) across 12 genomes (1 bacteria and 11 eukaryotes), in total 250,599 positive samples
	iDNA-ABT	[36]	Sep 2021	12 transformer layers with 12 attention heads in each layer.	BERT-based	DNA sequence	ChIP-seq data, ATAC-seq data, and histone modifications (HM) data of three human cell lines [33, 34], and DNA methylation dataset from the iDNA-MS [35]	3 main types of DNA methylation sites (6mA, 4mC, and 5hmC) across 12 genomes (1 bacteria and 11 eukaryotes), in total 250,599 positive samples
	MuLan-Methyl	[37]	July 2023	12 layers in the encoder stack, 768 hidden units for feed-forward networks, and 12 attention heads.	BERT-based	DNA sequence	ChIP-seq data, ATAC-seq data, and histone modifications (HM) data of three human cell lines [33, 34], and DNA methylation dataset from the iDNA-MS [35]	3 main types of DNA methylation sites (6mA, 4mC, and 5hmC) across 12 genomes (1 bacteria and 11 eukaryotes), in total 250,599 positive samples
RNA sequence language model	RNA-MSM	[38]	Nov 2023	12 attention heads with embedding size of 768	BERT-based	RNA sequence	RNA secondary structure and three-dimensional RNA structures [39]	The training, validation, and test sets have 405, 40, and 70 RNAs.
							RNA secondary structure and three-dimensional RNA structures [39]	The training, validation, and test sets have 405, 40, and 70 RNAs.
	RNA-FM	[40]	Arxiv posted Apr 2022	12 transformer-based bidirectional encoder blocks and 640 embedding	BERT-based	RNA sequence	RNA secondary structure [41, 42]	37149 structures from 8 RNA types of RNAStralign and 3975 RNA structures from 10 RNA types of ArchiveII
							RNA secondary structure [41, 42]	37149 structures from 8 RNA types of RNAStralign and 3975 RNA structures from 10 RNA types of ArchiveII

							whole genome of Severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2) [43]	Whole genome	SARS-CoV-2 genome structure and evolution prediction
							In vivo RNA secondary structure profiles for RNA-protein interaction [44]	-	Protein-RNA interaction prediction
							Human 5'UTR library [45]	83,919 5'UTRs of 75 different lengths and their corresponding mean ribosome loadings	mRNA 5' UTR-based mean ribosome loading prediction
RNABERT	[46]	Feb 2022	6 hidden layers of BERT, an embedding layer, one bidirectional-LSTM unit, two dense layers one with ReLU activation and a softmax output layer of LSTM	BERT-based	RNA sequence	RNA (ncRNA) families from RFam database[46]	31 RNA families	classifying RNA families RNA secondary structure prediction	
SpliceBERT	[47]	Mar 2024	6 transformer encoder layers, 512 hidden layer and 16 attention heads	BERT-based	RNA sequence	Reference genomes in fasta [48]	The pre-mRNA sequences from 72 vertebrate genomes for pre-training	Estimating splice sites	
BERT-m7G	[49]	Aug 2021	-	BERT-based	RNA sequence	RNA sequence with N7-methylguanosine sites and RNA sequence without N7-methylguanosine sites AlkAniline-Seq, MeRIP-seq, and miCLIP-seq [50]	741 RNA sequences with N7-methylguanosine sites and 741 RNA sequences without N7-methylguanosine sites	RNA N7-methylguanosine sites prediction	
M6A-BERT-Stacking	[51]	March 2023	12 transformer layers with 12 attention heads in each layer.	BERT-based	RNA sequence	RNA sequence with m6A sites and RNA sequence without m6A sites identified from MeRIP, m6A-seq, PAm6A-seq, and miCLIP [52]	11 datasets including 3000~16000 RNA sequences for each dataset	RNA m6A sites prediction	
Bert2Ome	[53]	May 2023	16 heads, 12 layers, and 1024 hidden units	BERT-based	RNA sequence	2-O-methylation modification sites from RMBase database [54]	215 positive, 215 negative instances for the training part and 46 positive, 114 negative instances for the testing part.	RNA 2-O-methylation prediction	
Rm-LR	[55]	Sep 2023	6 transformer	BERT-	RNA	Transcriptomic-wide	20 different epi-transcriptome	Multiple types of RNA	

				encoder layers, with a hidden layer size of 512 and 16 attention heads	based	sequence	profiling data derived from the MultiRM, GEO, RMBase, RADAR [54, 56-58]	profiles based on various base resolution techniques	modifications prediction
BertNDA	[59]	Nov 2023	8 layers	BERT-based	RNA sequence	miRNA-disease associations, lncRNA-disease associations, and miRNA-lncRNA associations [60-63]	1000 positive pairs and 1000 negative pairs	ncRNA-Disease Association Prediction	
LncCat	[64]	Feb 2023	-	BERT-based	RNA sequence	lncRNAs and protein-coding transcripts of five species [10, 63, 65]	22960 coding transcripts and 21081 lncRNAs of human. 20707 coding transcripts and 10707 lncRNAs of mouse. 15891 coding transcripts and 4382 lncRNAs of zebrafish. 4693 coding transcripts and 5377 lncRNAs of wheat 20584 coding transcripts and 3897 lncRNAs of chicken	Identify lncRNA	
LSCPP-BERT	[66]	Dec 2023	4 identical layers and each layer is divided into two sublayers	BERT-based	RNA sequence	lncRNAs sequences from multispecies [67]	593251 plant lncRNAs sequences	lncRNA-sORFs coding potential prediction	
CodonBERT	[68]	Oct 2023	12 layers of bidirectional transformer encoders, Each transformer layer with 12 self-attention heads	BERT-based	RNA sequence	mRFP Expression dataset; [69] Fungal expression dataset; E. coli proteins dataset; mRNA stability dataset; Te-Riboswitches dataset [70]	Experimental data for protein expression (2308 low expression proteins, 2067 medium expression proteins, and 1973 high expression proteins, respectively);	mRNA properties prediction	
						SARS-CoV-2 Vaccine degradation dataset	-	Vaccine expression prediction	
RNA-TorsionBERT	[71]	Jun 2024	18 layer	BERT-based	RNA sequence	PDB structure and removed the structures from the nonredundant Training	4,267 structures with sequences from 11 to 508 nucleotides	RNA 3D structure prediction	
UNI-RNA	[72]	bioRxiv posted Jul 2023	-	BERT-based	RNA sequence	non-coding RNA sequences from RNACentral, nucleic acid data from NCBI's database, and genomic data from repositories such as Genome Warehouse[73-75]	37,149 RNA structures 7,600 held-out real human 5'UTRs 3 million distinct UTR sequences	RNA secondary structure prediction RNA distance map prediction mRNA 5'-UTR mean ribosome load prediction Alternative	

									polyadenylation isoform prediction RNA splice site prediction ncRNA classification RNA modification prediction
UTR-LM	[76]	Apr 2024	Six-layer transformer with 16 multi-head self-attention.	Transformer-based	RNA sequence	Unlabeled 5' UTR sequences from three sources: the Ensembl database	214,349 unlabeled 5' UTR sequences	Mean ribosome loading prediction. mRNA expression level and translation efficiency prediction Internal ribosome entry site identification Attention-based motif detection	
3UTRBERT	[77]	Aug 2024	12 identical Transformer components, Each layer contained a multi-head self-attention module and a position-wise fully connected feed-forward layer	BERT-based	RNA sequence	3'UTR of human mRNA transcript eCLIP datasets RNA localization data	108 573 unique mRNA transcripts 1000 samples RNA binding samples with 1000 sequences 17 023 mRNAs	mRNA subcellular localization prediction	
RNAErnie	[78]	May 2024	12-layer transformer and a hidden state dimension of 768	Transformer-based	RNA sequence	non-coding RNA sequences from RNacentral		RNA grouping RNA sequence classification RNA–RNA interaction prediction RNA secondary structure prediction	

References

1. Ji, Y., et al., *DNABERT: pre-trained Bidirectional Encoder Representations from Transformers model for DNA-language in genome*. Bioinformatics, 2021. **37**(15): p. 2112-2120.
2. Dreos, R., et al., *EPD and EPDnew, high-quality promoter resources in the next-generation sequencing era*. Nucleic acids research, 2013. **41**(D1): p. D157-D164.
3. Consortium, E.P., *An integrated encyclopedia of DNA elements in the human genome*. Nature, 2012. **489**(7414): p. 57.
4. Cunningham, F., et al., *Ensembl 2019*. Nucleic acids research, 2019. **47**(D1): p. D745-D751.
5. Sherry, S.T., et al., *dbSNP: the NCBI database of genetic variation*. Nucleic acids research, 2001. **29**(1): p. 308-311.
6. Zhou, Z., et al., *Dnabert-2: Efficient foundation model and benchmark for multi-species genome*. arXiv preprint arXiv:2306.15006, 2023.
7. Zeng, H., et al., *Convolutional neural network architectures for predicting DNA–protein binding*. Bioinformatics, 2016. **32**(12): p. i121-i127.
8. Chen, K., H. Zhao, and Y. Yang, *Capturing large genomic contexts for accurately predicting enhancer-promoter interactions*. Briefings in Bioinformatics, 2022. **23**(2): p. bbab577.
9. Dalla-Torre, H., et al., *The nucleotide transformer: Building and evaluating robust foundation models for human genomics*. bioRxiv, 2023: p. 2023.01. 11.523679.
10. Howe, K.L., et al., *Ensembl 2021*. Nucleic acids research, 2021. **49**(D1): p. D884-D891.
11. Bergström, A., et al., *Insights into human genetic variation and population history from 929 diverse genomes*. Science, 2020. **367**(6484): p. eaay5012.
12. Alonso-Blanco, C., et al., *1,135 genomes reveal the global pattern of polymorphism in *Arabidopsis thaliana**. Cell, 2016. **166**(2): p. 481-491.
13. Zhang, D., et al., *DNAGPT: A Generalized Pretrained Tool for Multiple DNA Sequence Analysis Tasks*. bioRxiv, 2023: p. 2023.07. 11.548628.
14. Kalkatawi, M., et al., *DeepGSR: an optimized deep-learning structure for the recognition of genomic signals and regions*. Bioinformatics, 2019. **35**(7): p. 1125-1132.
15. Agarwal, V. and J. Shendure, *Predicting mRNA abundance directly from genomic sequence using deep convolutional neural networks*. Cell

- reports, 2020. **31**(7).
- 16. Sanabria, M., et al., *DNA language model GROVER learns sequence context in the human genome*. Nature Machine Intelligence, 2024. **6**(8): p. 911-923.
 - 17. de Souza, N., *The ENCODE project*. Nature methods, 2012. **9**(11): p. 1046-1046.
 - 18. Benegas, G., S.S. Batra, and Y.S. Song, *DNA language models are powerful zero-shot predictors of genome-wide variant effects*. bioRxiv, 2022: p. 2022.08. 22.504706.
 - 19. Le, N.Q.K., et al., *BERT-Promoter: An improved sequence-based predictor of DNA promoter using BERT pre-trained model and SHAP feature selection*. Computational Biology and Chemistry, 2022. **99**: p. 107732.
 - 20. Gama-Castro, S., et al., *RegulonDB version 9.0: high-level integration of gene regulation, coexpression, motif clustering and beyond*. Nucleic acids research, 2016. **44**(D1): p. D133-D143.
 - 21. Xiao, X., et al., *iPSW (2L)-PseKNC: A two-layer predictor for identifying promoters and their strength by hybrid features via pseudo K-tuple nucleotide composition*. Genomics, 2019. **111**(6): p. 1785-1793.
 - 22. Luo, H., et al., *Improving language model of human genome for DNA–protein binding prediction based on task-specific pre-training*. Interdisciplinary Sciences: Computational Life Sciences, 2023. **15**(1): p. 32-43.
 - 23. An, W., et al. *MoDNA: motif-oriented pre-training for DNA language model*. in *Proceedings of the 13th ACM International Conference on Bioinformatics, Computational Biology and Health Informatics*. 2022.
 - 24. Luo, H., et al. *iEnhancer-BERT: A novel transfer learning architecture based on DNA-Language model for identifying enhancers and their strength*. in *International Conference on Intelligent Computing*. 2022. Springer.
 - 25. Ernst, J., et al., *Mapping and analysis of chromatin state dynamics in nine human cell types*. Nature, 2011. **473**(7345): p. 43-49.
 - 26. Tsukiyama, S., et al., *BERT6mA: prediction of DNA N6-methyladenine site using deep learning-based approaches*. Briefings in Bioinformatics, 2022. **23**(2): p. bbac053.
 - 27. Xiao, C.-L., et al., *N6-methyladenine DNA modification in the human genome*. Molecular cell, 2018. **71**(2): p. 306-318. e7.
 - 28. Ye, G., et al., *De novo genome assembly of the stress tolerant forest species Casuarina equisetifolia provides insight into secondary growth*. The Plant Journal, 2019. **97**(4): p. 779-794.
 - 29. Ye, P., et al., *MethSMRT: an integrative database for DNA N6-methyladenine and N4-methylcytosine generated by single-molecular real-*

- time sequencing*. Nucleic acids research, 2016: p. gkw950.
- 30. Liu, Z.-Y., et al., *MDR: an integrative DNA N6-methyladenine and N4-methylcytosine modification database for Rosaceae*. Horticulture research, 2019. **6**.
 - 31. Wang, Y., et al., *N6-adenine DNA methylation is associated with the linker DNA of H2A. Z-containing well-positioned nucleosomes in Pol II-transcribed genes in Tetrahymena*. Nucleic acids research, 2017. **45**(20): p. 11594-11606.
 - 32. Jin, J., et al., *iDNA-ABF: multi-scale deep biological language learning model for the interpretable prediction of DNA methylations*. Genome biology, 2022. **23**(1): p. 1-23.
 - 33. Luo, Y., et al., *New developments on the Encyclopedia of DNA Elements (ENCODE) data portal*. Nucleic acids research, 2020. **48**(D1): p. D882-D889.
 - 34. Zhang, J., et al., *An integrative ENCODE resource for cancer genomics*. Nature communications, 2020. **11**(1): p. 3696.
 - 35. Lv, H., et al., *iDNA-MS: an integrated computational tool for detecting DNA modification sites in multiple genomes*. Iscience, 2020. **23**(4).
 - 36. Yu, Y., et al., *iDNA-ABT: advanced deep learning model for detecting DNA methylation with adaptive features and transductive information maximization*. Bioinformatics, 2021. **37**(24): p. 4603-4610.
 - 37. Zeng, W., A. Gautam, and D.H. Huson, *MuLan-Methyl-Multiple Transformer-based Language Models for Accurate DNA Methylation Prediction*. bioRxiv, 2023: p. 2023.01.04.522704.
 - 38. Zhang, Y., et al., *Multiple sequence alignment-based RNA language model and its application to structural inference*. Nucleic Acids Research, 2024. **52**(1): p. e3-e3.
 - 39. Singh, J., et al., *RNA secondary structure prediction using an ensemble of two-dimensional deep neural networks and transfer learning*. Nature communications, 2019. **10**(1): p. 5407.
 - 40. Chen, J., et al., *Interpretable RNA foundation model from unannotated data for highly accurate RNA structure and function predictions*. bioRxiv, 2022: p. 2022.08.06.503062.
 - 41. Tan, Z., et al., *TurboFold II: RNA structural alignment and secondary structure prediction informed by multiple homologs*. Nucleic acids research, 2017. **45**(20): p. 11570-11581.
 - 42. Sloma, M.F. and D.H. Mathews, *Exact calculation of loop formation probability identifies folding motifs in RNA secondary structures*. RNA, 2016. **22**(12): p. 1808-1818.

43. Wu, F., et al., *A new coronavirus associated with human respiratory disease in China*. Nature, 2020. **579**(7798): p. 265-269.
44. Sun, L., et al., *Predicting dynamic cellular protein–RNA interactions by deep learning using in vivo RNA structures*. Cell research, 2021. **31**(5): p. 495-516.
45. Sample, P.J., et al., *Human 5' UTR design and variant effect prediction from a massively parallel translation assay*. Nature biotechnology, 2019. **37**(7): p. 803-809.
46. Akiyama, M. and Y. Sakakibara, *Informative RNA base embedding for RNA structural alignment and clustering by deep representation learning*. NAR genomics and bioinformatics, 2022. **4**(1): p. lqac012.
47. Chen, K., et al., *Self-supervised learning on millions of primary RNA sequences from 72 vertebrates improves sequence-based RNA splicing prediction*. Briefings in Bioinformatics, 2024. **25**(3): p. bbae163.
48. Haeussler, M., et al., *The UCSC genome browser database: 2019 update*. Nucleic acids research, 2019. **47**(D1): p. D853-D858.
49. Zhang, L., et al., *BERT-m7G: a transformer architecture based on BERT and stacking ensemble to identify RNA N7-Methylguanosine sites from sequence information*. Computational and Mathematical Methods in Medicine, 2021. **2021**.
50. Dai, C., et al., *Iterative feature representation algorithm to improve the predictive performance of N7-methylguanosine sites*. Briefings in Bioinformatics, 2021. **22**(4): p. bbaa278.
51. Li, Q., et al., *M6A-BERT-Stacking: A Tissue-Specific Predictor for Identifying RNA N6-Methyladenosine Sites Based on BERT and Stacking Strategy*. Symmetry, 2023. **15**(3): p. 731.
52. Dao, F.-Y., et al., *Computational identification of N6-methyladenosine sites in multiple tissues of mammals*. Computational and structural biotechnology journal, 2020. **18**: p. 1084-1091.
53. Soylu, N.N. and E. Sefer, *BERT2OME: Prediction of 2'-O-methylation Modifications from RNA Sequence by Transformer Architecture Based on BERT*. IEEE/ACM Transactions on Computational Biology and Bioinformatics, 2023.
54. Xuan, J.-J., et al., *RMBase v2. 0: deciphering the map of RNA modifications from epitranscriptome sequencing data*. Nucleic acids research, 2018. **46**(D1): p. D327-D334.
55. Liang, S., et al., *Rm-LR: A long-range-based deep learning model for predicting multiple types of RNA modifications*. Computers in Biology and Medicine, 2023. **164**: p. 107238.
56. Song, Z., et al., *Attention-based multi-label neural networks for integrated prediction and interpretation of twelve widely occurring RNA*

- modifications.* Nature communications, 2021. **12**(1): p. 4011.
- 57. Barrett, T., et al., *NCBI GEO: archive for functional genomics data sets—update*. Nucleic acids research, 2012. **41**(D1): p. D991-D995.
 - 58. Ramaswami, G. and J.B. Li, *RADAR: a rigorously annotated database of A-to-I RNA editing*. Nucleic acids research, 2014. **42**(D1): p. D109-D113.
 - 59. Ning, Z., et al., *BertNDA: a Model Based on Graph-Bert and Multi-scale Information Fusion for ncRNA-disease Association Prediction*. bioRxiv, 2023: p. 2023.05. 18.541387.
 - 60. Li, Y., et al., *HMDD v2. 0: a database for experimentally supported human microRNA and disease associations*. Nucleic acids research, 2014. **42**(D1): p. D1070-D1074.
 - 61. Jiang, Q., et al., *miR2Disease: a manually curated database for microRNA deregulation in human disease*. Nucleic acids research, 2009. **37**(suppl_1): p. D98-D104.
 - 62. Bao, Z., et al., *LncRNADisease 2.0: an updated database of long non-coding RNA-associated diseases*. Nucleic acids research, 2019. **47**(D1): p. D1034-D1037.
 - 63. Gao, Y., et al., *Lnc2Cancer 3.0: an updated resource for experimentally supported lncRNA/circRNA cancer associations and web tools based on RNA-seq and scRNA-seq data*. Nucleic acids research, 2021. **49**(D1): p. D1251-D1258.
 - 64. Feng, H., et al., *LncCat: An ORF attention model to identify LncRNA based on ensemble learning strategy and fused sequence information*. Computational and Structural Biotechnology Journal, 2023. **21**: p. 1433-1447.
 - 65. O'Leary, N.A., et al., *Reference sequence (RefSeq) database at NCBI: current status, taxonomic expansion, and functional annotation*. Nucleic acids research, 2016. **44**(D1): p. D733-D745.
 - 66. Xia, S., et al. *A multi-granularity information-enhanced pre-training method for predicting the coding potential of sORFs in plant lncRNAs*. in *2023 IEEE International Conference on Bioinformatics and Biomedicine (BIBM)*. 2023. IEEE.
 - 67. Di Marsico, M., et al., *GreeNC 2.0: a comprehensive database of plant long non-coding RNAs*. Nucleic Acids Research, 2022. **50**(D1): p. D1442-D1447.
 - 68. Babjac, A.N., Z. Lu, and S.J. Emrich. *CodonBERT: Using BERT for Sentiment Analysis to Better Predict Genes with Low Expression*. in *Proceedings of the 14th ACM International Conference on Bioinformatics, Computational Biology, and Health Informatics*. 2023.
 - 69. Nieuwkoop, T., et al., *Revealing determinants of translation efficiency via whole-gene codon randomization and machine learning*. Nucleic

- acids research, 2023. **51**(5): p. 2363-2376.
70. Byrska-Bishop, M., et al., *High-coverage whole-genome sequencing of the expanded 1000 Genomes Project cohort including 602 trios*. Cell, 2022. **185**(18): p. 3426-3440. e19.
71. Bernard, C., et al., *RNA-TorsionBERT: leveraging language models for RNA 3D torsion angles prediction*. bioRxiv, 2024: p. 2024.06. 06.597803.
72. Wang, X., et al., *UNI-RNA: universal pre-trained models revolutionize RNA research*. bioRxiv, 2023: p. 2023.07. 11.548588.
73. *RNACentral: a hub of information for non-coding RNA sequences*. Nucleic Acids Research, 2019. **47**(D1): p. D221-D229.
74. Sayers, E.W., et al., *Database resources of the national center for biotechnology information*. Nucleic acids research, 2022. **50**(D1): p. D20-D26.
75. Chen, M., et al., *Genome Warehouse: a public repository housing genome-scale data*. Genomics, Proteomics and Bioinformatics, 2021. **19**(4): p. 584-589.
76. Chu, Y., et al., *A 5' UTR language model for decoding untranslated regions of mRNA and function predictions*. Nature Machine Intelligence, 2024. **6**(4): p. 449-460.
77. Yang, Y., et al., *Deciphering 3'UTR Mediated Gene Regulation Using Interpretable Deep Representation Learning*. Advanced Science, 2024. **11**(39): p. 2407013.
78. Wang, N., et al., *Multi-purpose RNA language modelling with motif-aware pretraining and type-guided fine-tuning*. Nature Machine Intelligence, 2024: p. 1-10.

Supplementary Table 2. Detailed information of large language models for proteomic tasks

Application area	Models	Ref	Publication time	Parameters	Architecture	Datasets			Downstream tasks
						Data type	Source	Size	
Protein Large Language Models	MSA Transformer	[1]	Jul 2021	100M parameters model with 12 layers, 768 embedding size, and 12 attention heads	Transformer-based	MSAs	CAMEO [2]	131 domains (129 evaluated)	Unsupervised contact prediction, supervised contact prediction, secondary structure prediction
							CASP13-FM [3]	31 free modeling domains (from 25 targets)	
							trRosetta training set [4]	15,051 MSAs and structures (14,842 used)	
						Protein sequences	CB513 [5]	513 protein sequences	
	UniRep	[8]	Dec 2019	18.2M parameters (a 1,900-hidden unit mLSTM with amino-acid character embeddings)	LSTM-based		Netsurf dataset [6]	12,185 crystal structures obtained from the PDB [7]	Predicting stability of naturally occurring and de novo designed proteins, Prediction of functional effects of single mutations in diverse proteins
					Protein sequences	UniRef50 [9]	24M		
						Mini protein dataset from [10]	1,432 out of 5,570 test set and 1,416 out of 5,571 validation set		
						DMS dataset [11]	~65,420 variants across the 8 proteins		
						avGFP dataset [12]	32,400 variants derived from 27 homologs of avGFP (used)		
	TAPE	[13]	Dec 2019	38M parameters	(ResNet, Transformer, LSTM)-based	Protein sequences	Netsurf dataset [6]	12,185 protein sequences	Secondary structure (SS) prediction (structure prediction task), contact prediction (structure prediction task), remote homology detection (evolutionary understanding task), fluorescence landscape prediction (protein engineering task), stability landscape prediction (protein engineering task)
							ProteinNet dataset [14]	~332,283,871 protein sequences	
							DeepSF dataset [15]	Training set includes 16,712 proteins spanning 1,195 folds [16]; test datasets include 2,533 protein domains across 550 folds from SCOP 2.06 [16], a subset of SCOP	

							1.75 and the CASP dataset [17, 18]	
						avGFP dataset [12]	~51,715 protein sequences	
						Dataset from [10]	~ 46,800 protein sequences	
ESM-1b	[19]	Dec 2020	A model with ~650M parameters (33 layers)	Transformer-based	Protein sequences	SCOPe [20]	15,297 protein sequences	Remote homology detection, prediction of secondary structure, long-range residue-residue contacts, mutational effect prediction, etc.
						CB513 [5]	513 protein sequences;	
						CASP13 [21]	431 domains	
						Envision (DMS dataset) [11] and DeepSequence [22]	Over 700,000 variant effect measurements from over 100 large-scale experimental mutagenesis datasets	
						CB513 [5]	513 protein sequences	Per-residue secondary structure prediction, per-protein localization & membrane prediction
						TS115 [5]	115 protein sequences	
						CASP12 [24]	~102 protein sequences	
ProtTrans	[23]	Aug 2021	From millions to billions parameters (224M-11B)	(Transformer-XL, XLNet, BERT, Albert, Electra, T5)-based	Protein sequences	NEW364 [23]	364 protein sequences	Per-residue secondary structure prediction, per-protein localization & membrane prediction
						DeepLoc [25]	~19,817 protein sequences	
						SCOPe 2.07 [26]	14,323 protein sequences (non-redundant at PIDE < 40%)	
						Netsurf dataset [6]	12,185 protein sequences	
						CB513 [5]	513 protein sequences	
						CASP12 [24]	~102 protein sequences	
						DeepSF dataset [15]	Consistent with the size used TAPE [13]	
SProBERTa	[27]	Sep 2022	12 Transformer encoder layers, the embedding size is 768, the feed-forward hidden units are 3072 and the attention heads are 12	BERT-based	Protein sequences			Secondary structure prediction, contact prediction, remote homology prediction, protein function prediction or Gene Ontology (GO) term prediction

						DeepFRI [28]	~284,832 protein sequences with GO annotations	
PromptProtein	[29]	Sep 2022	650M parameters with 33 layers and 20 attention heads. The embedding size is 1280	Transformer-based	Protein sequences	EC dataset from [28]	19,199 protein sequences	Enzyme commission and Gene Ontology prediction, stability landscape prediction, fluorescence landscape prediction, thermostability landscape prediction, adeno-associated virus (AAV) landscape prediction, GB1 landscape prediction, antibody-antigen affinity prediction
						GO dataset from [28]	36,641 protein sequences	
						Protein engineering dataset from TAPE and FLIP [13, 30]	68,965 protein sequences used in stability prediction, 54,025 used in fluorescence landscape prediction, 28,131 used in thermostability landscape prediction, 82,583 in AAV prediction, 8,733 in GB1 prediction	
						Uniparc [32], UniprotKB [33], SWISS-PROT [34], TrEMBL [35], Pfam [36], and NCBI taxonomic information [37]	281M	
ProGen	[31]	Mar 2020	1.2B parameters. Sequence length is 512. The model has dimension $d = 1028$, inner dimension $f = 512$, 36 layers, and 8 heads per layer. Dropout with probability 0.1 follows the residual connections in each layer	Transformer-based	Protein sequences, Conditioning tags	Protein Gym (a DMS dataset) [38]	~1.8M protein sequences	Controllable protein generation and two case study: completing VEGFR2 kinase domain, zero-shot fitness selection for protein GB1
Tranception	[38]	Jun 2022	700M parameters	Transformer-based	Protein sequences	UniRef50 [40]	10,000 protein sequences	Fitness prediction
ProtGPT2	[39]	Jul 2022	738M parameters. The model consists of 36 layers with a model dimensionality of 1280. The architecture matches that of the previously released GPT2-large	GPT-based	Protein sequences	ProtGPT2 dataset [39]	Generated 10,000 protein sequences	Sequence dataset generation, homology detection, disorder prediction
ProteinBERT	[41]	Feb 2022	16M parameters. The model architecture consists of two almost parallel paths: one for local representations with	BERT-based	Protein sequences, Gene Ontology (GO) annotations	Secondary structure dataset from	8,678 sequences (train) [13, 42]	(Secondary structure, disorder, Remote homology, fold classes, signal peptide, major PTMs, neuropeptide cleavage,
						Disorder dataset	8,678 sequences (train)	

			d=128 and the other for global representations with d=512			from	[42]	fluorescence, stability) prediction
						Remote homology dataset from	12,312 sequences (train) [13, 43, 44]	
						Fold classes dataset from	15,680 sequences (train) [43, 44]	
						Signal peptide dataset from	16,606 sequences (train) [45]	
						Major PTMs dataset from	43,356 sequences (train) [46]	
						Neuropeptide cleavage dataset from	2,727 sequences (train) [47, 48]	
						Fluorescence dataset from	21,446 sequences (train) [12, 13]	
						Stability dataset from	53,679 sequences (train) [10]	
ProtST	[49]	Jan 2023	Depends on the parameters of the chosen language model	Multi-models-based (Protein Language Models and Biomedical Language Models)	Protein sequences, property descriptions	ProtDescribe [50, 51]	553,052 aligned pairs of protein sequence and property description	Protein localization prediction, fitness landscape prediction, protein function annotation (totally 11 downstream tasks)
KeAP	[52]	Jan 2023	Depends on the parameters of the chosen language model	Multi-cascade Bert-like network	Triplet in the format of (Protein, Relation, Attribute)	ProteinKG25 [53]	5M with nearly 600k protein, 50k attribute terms, and 31 relation terms included	Amino acid contact prediction, protein homology detection, protein stability prediction, protein-protein interaction identification, protein-protein binding affinity prediction, and semantic similarity inference
CaLM	[54]	Jan 2024	86M parameters. 12 transformer layers contain 12 attention heads, with dimension 768. Similar to architectures of ESM family	Transformer-based	Protein-coding DNA (cDNA)	cDNA dataset obtained from the European Nucleotide Archive with a timestamp of April 2022	9,858,385 cDNA sequences of seven model organisms	Melting point prediction, solubility prediction, subcellular localization prediction and function prediction

						Melting temperature dataset [30, 55]	-	
						Subcellular localization dataset [30, 56]	-	
						Solubility dataset [57]	-	
						Gene ontology dataset [58]	-	
						Transcriptomics dataset [59]	-	
						Proteomics dataset [60]	-	
PLMSearch	[61]	Mar 2024	ESM-1b (650M parameters) [19] and ProtT5-XL-UniRef50 (3B parameters) [23]	Transformer-based	Protein sequences	SCOPe40 [20, 62], New protein search test, Swiss-Prot [34], CATHS40 [63]	~489,764 sequences for training	Homologous protein search
DHR	[64]	Jul 2024	2 ESM-1b (650M parameters per encoder) [19]	Transformer-based	Protein sequences	UniRef [9], SCOPe [20, 62]	~2 M	Protein homolog detection
Antibody Large Language Models	MHCRoBERTa	[65]	Dec 2021	Model with 12 multi-heads and 5 self-attention layers.	ROBERTa-based	Protein sequences	UniProtKB [66] 565,254 protein sequences	Predicting the binding of peptide and major histocompatibility complex (MHC)

						Immune Epitope Database (IEDB) [67]	MHC class I transmembrane proteins containing HLA-A(1,777 sequences), HLA-B (2,100 sequences) and HLA-C(1,931 sequences)	
BERTMHC	[68]	Jun 2021	The model has 12 layers with 12 self-attention heads in each layer	BERT-based	Protein sequences	The data from Kamilla Kjærgaard Jensen [69]	2,413 additional MHC-peptide pairs covering 47 MHC class II alleles.	Predicting precisely the binding and presentation of peptides to major histocompatibility complex (MHC) alleles
						Immune Epitope Database (IEDB) [67]	95,638 peptides	
TCR-BERT	[70]	BioRxiv posted Nov 2021	12 stacked transformer blocks with 8 attention heads, utilizing a hidden representation dimensionality of 768 and featuring 12 transformer layers.	BERT-based	Protein sequences	Pan immune repertoire database (PIRD) [71]	47,040 TRB sequences and 4,607 TRA sequences.	Antigen specificity classification
						VDJdb [72]	58,795 human TCRs and 3,353 mouse TCRs.	
						TCRdb [73]	139,00,913 TRB sequences of unknown antigen binding affinity.	
Antiformer	[74]	July 2024	he transformer encoder involves 12 stacking layers of transformer with the multi-head self-attention and feed forward network.	BERT-based	Protein sequences Gene expression	The OAS database [75]	55 BCR-seq datasets containing 600 million sequences	Binding specificity prediction
SC-AIR-BERT	[76]	May 2023	Six standard transformer layers and each layer has four attention heads, the hidden representation dimensionality is 512 and the intermediate representation dimensionality is 2048	BERT-based	Protein sequences	VDJdb [72]	23,358 unique paired TCRs	Binding specificity prediction
						Immune Epitope Database (IEDB) [67]	18,662 paired TCRs and 589 paired BCRs	
						huARdb [77]	612,077 high-confidence paired full-length α/β chains of	

							TCR sequences	
						CoV-AbDab [78]	1,105,906 paired antibody heavy/light chains	
AbLang	[79]	Jun 2022	Consists of three modules (Each of AbRep's 12 transformer blocks has 12 attenuated heads, an inner hidden size of 3072 and a hidden size of 768. From AbRep, the rescodings (768 values for each residue) are obtained. AbHead follows the design of RoBERTa's[80] head model, with a hidden size of 768)	BERT-based	Human antibody sequences	Observed Antibody Space (OAS) database [81]	Training sets of 14,126 724 heavy and 187,068 light sequences, and two evaluation sets of 100,000 heavy and 50,000 light sequences	Sequence specific predictions, residue specific predictions, amino acid predictions
AntiBERTa		Jul 2022	86M parameters. A 12-layer transformer model. Attention heads is 12, embedding dimension is 768, feedforward layer dimension is 3072	BERT-based	Human antibody sequences	Data from [82]	10,000 naive and 10,000 memory B-cell sequences	Trace the B cell origin of the antibody, quantify immunogenicity, predict the antibody's binding site
EATLM	[85]	Jan 2023	86M parameters (12 layers, 12 heads, and 768 hidden states)	Transformer-based	Human antibody sequences	SAbDab [83]	Training/validation/test split of 720/90/90	Accurate antigen-binding prediction, paratope prediction, B cell analysis, antibody discovery
						BCR repertoire dataset [82]	-	
						TheraSAbDab [84]	191 non-redundant therapeutic antibodies	
						Dataset from [86]	Training/validation/test split of 15,128/3,242/3,242	
						Paratope data from [87]	1,662 CDR segments on 277 antibodies	
						Data from [88]	88,094 sequences with 6 maturation stages	
						A subset of the OAS database [81]	Antibody sequences from 133 SARS-CoV-2 patients and 87 health persons	

References

1. Rao, R.M., et al. *MSA transformer*. in *International Conference on Machine Learning*. 2021.
2. Haas, J., et al., *Continuous Automated Model EvaluatiOn (CAMEO) complementing the critical assessment of structure prediction in CASP12*. Proteins: Structure, Function, and Bioinformatics, 2018. **86**: p. 387-398.
3. Shrestha, R., et al., *Assessing the accuracy of contact predictions in CASP13*. Proteins: Structure, Function, and Bioinformatics, 2019. **87**(12): p. 1058-1068.
4. Yang, J., et al., *Improved protein structure prediction using predicted interresidue orientations*. Proceedings of the National Academy of Sciences, 2020. **117**(3): p. 1496-1503.
5. Cuff, J.A. and G.J. Barton, *Evaluation and improvement of multiple sequence methods for protein secondary structure prediction*. Proteins: Structure, Function, and Bioinformatics, 1999. **34**(4): p. 508-519.
6. Klausen, M.S., et al., *NetSurfP-2.0: Improved prediction of protein structural features by integrated deep learning*. Proteins: Structure, Function, and Bioinformatics, 2019. **87**(6): p. 520-527.
7. Berman, H.M., et al., *The protein data bank*. Nucleic acids research, 2000. **28**(1): p. 235-242.
8. Alley, E.C., et al., *Unified rational protein engineering with sequence-based deep representation learning*. Nature methods, 2019. **16**(12): p. 1315-1322.
9. Suzek, B.E., et al., *UniRef clusters: a comprehensive and scalable alternative for improving sequence similarity searches*. Bioinformatics, 2015. **31**(6): p. 926-932.
10. Rocklin, G.J., et al., *Global analysis of protein folding using massively parallel design, synthesis, and testing*. Science, 2017. **357**(6347): p. 168-175.
11. Gray, V.E., et al., *Quantitative missense variant effect prediction using large-scale mutagenesis data*. Cell systems, 2018. **6**(1): p. 116-124.
12. Sarkisyan, K.S., et al., *Local fitness landscape of the green fluorescent protein*. Nature, 2016. **533**(7603): p. 397-401.
13. Rao, R., et al., *Evaluating protein transfer learning with TAPE*. Advances in neural information processing systems, 2019. **32**.
14. AlQuraishi, M., *ProteinNet: a standardized data set for machine learning of protein structure*. BMC bioinformatics, 2019. **20**: p. 1-10.
15. Hou, J., B. Adhikari, and J. Cheng, *DeepSF: deep convolutional neural network for mapping protein sequences to folds*. Bioinformatics, 2018. **34**(8): p. 1295-1303.

16. Murzin, A.G., et al., *SCOP: a structural classification of proteins database for the investigation of sequences and structures*. Journal of molecular biology, 1995. **247**(4): p. 536-540.
17. Kinch, L.N., et al., *CASP9 target classification*. PROTEINS: structure, function, and bioinformatics, 2011. **79**(S10): p. 21-36.
18. Kinch, L.N., et al., *CASP 11 target classification*. Proteins: Structure, Function, and Bioinformatics, 2016. **84**: p. 20-33.
19. Rives, A., et al., *Biological structure and function emerge from scaling unsupervised learning to 250 million protein sequences*. Proceedings of the National Academy of Sciences, 2021. **118**(15): p. e2016239118.
20. Fox, N.K., S.E. Brenner, and J.-M. Chandonia, *SCOPe: Structural Classification of Proteins—extended, integrating SCOP and ASTRAL data and classification of new structures*. Nucleic acids research, 2014. **42**(D1): p. D304-D309.
21. Moult, J., et al., *Critical assessment of methods of protein structure prediction: Progress and new directions in round XI*. Proteins: Structure, Function, and Bioinformatics, 2016. **84**: p. 4-14.
22. Riesselman, A.J., J.B. Ingraham, and D.S. Marks, *Deep generative models of genetic variation capture the effects of mutations*. Nature methods, 2018. **15**(10): p. 816-822.
23. Elnaggar, A., et al., *ProtTrans: Towards Cracking the Language of Lifes Code Through Self-Supervised Deep Learning and High Performance Computing*. IEEE Transactions on Pattern Analysis and Machine Intelligence, 2021: p. 1-1.
24. Abriata, L.A., et al., *Assessment of hard target modeling in CASP12 reveals an emerging role of alignment-based contact prediction methods*. Proteins: Structure, Function, and Bioinformatics, 2018. **86**: p. 97-112.
25. Almagro Armenteros, J.J., et al., *DeepLoc: prediction of protein subcellular localization using deep learning*. Bioinformatics, 2017. **33**(21): p. 3387-3395.
26. Chandonia, J.-M., N.K. Fox, and S.E. Brenner, *SCOPe: classification of large macromolecular structures in the structural classification of proteins—extended database*. Nucleic acids research, 2019. **47**(D1): p. D475-D481.
27. Wu, L., et al., *SProBERTa: protein embedding learning with local fragment modeling*. Briefings in Bioinformatics, 2022. **23**(6): p. bbac401.
28. Gligorijevi, V., et al., *Structure-based protein function prediction using graph convolutional networks*. Nature communications, 2021. **12**(1): p. 3168.
29. Wang, Z., et al. *Multi-level Protein Structure Pre-training via Prompt Learning*. in *The Eleventh International Conference on Learning Representations*. 2022.

30. Dallago, C., et al., *FLIP: Benchmark tasks in fitness landscape inference for proteins*. bioRxiv, 2021: p. 2021-11.
31. Madani, A., et al., *Progen: Language modeling for protein generation*. arXiv preprint arXiv:2004.03497, 2020.
32. Leinonen, R., et al., *UniProt archive*. Bioinformatics, 2004. **20**(17): p. 3236-3237.
33. Bairoch, A., et al., *The universal protein resource (UniProt)*. Nucleic acids research, 2005. **33**(suppl_1): p. D154-D159.
34. Bairoch, A., et al., *Swiss-Prot: juggling between evolution and stability*. Briefings in bioinformatics, 2004. **5**(1): p. 39-55.
35. Boeckmann, B., et al., *The SWISS-PROT protein knowledgebase and its supplement TrEMBL in 2003*. Nucleic acids research, 2003. **31**(1): p. 365-370.
36. Bateman, A., et al., *The Pfam protein families database*. Nucleic acids research, 2004. **32**(suppl_1): p. D138-D141.
37. Federhen, S., *The NCBI taxonomy database*. Nucleic acids research, 2012. **40**(D1): p. D136-D143.
38. Notin, P., et al. *Tranception: protein fitness prediction with autoregressive transformers and inference-time retrieval*. in *International Conference on Machine Learning*. 2022.
39. Ferruz, N., S. Schmidt, and B. Hcker, *ProtGPT2 is a deep unsupervised language model for protein design*. Nature communications, 2022. **13**(1): p. 4348.
40. *UniProt: the universal protein knowledgebase in 2021*. Nucleic acids research, 2021. **49**(D1): p. D480-D489.
41. Brandes, N., et al., *ProteinBERT: a universal deep-learning model of protein sequence and function*. Bioinformatics, 2022. **38**(8): p. 2102-2110.
42. Moult, J., et al., *Critical assessment of methods of protein structure prediction (CASP)—Round XII*. Proteins: Structure, Function, and Bioinformatics, 2018. **86**: p. 7-15.
43. Andreeva, A., et al., *SCOP2 prototype: a new approach to protein structure mining*. Nucleic acids research, 2014. **42**(D1): p. D310-D314.
44. Andreeva, A., et al., *The SCOP database in 2020: expanded classification of representative family and superfamily domains of known protein structures*. Nucleic acids research, 2020. **48**(D1): p. D376-D382.
45. Armenteros, J.J.A., et al., *SignalP 5.0 improves signal peptide predictions using deep neural networks*. Nature biotechnology, 2019. **37**: p. 420-423.
46. Hornbeck, P.V., et al., *PhosphoSitePlus, 2014: mutations, PTMs and recalibrations*. Nucleic acids research, 2015. **43**(D1): p. D512-D520.
47. Ofer, D. and M. Linial, *ProFET: Feature engineering captures high-level protein functions*. Bioinformatics, 2015. **31**(21): p. 3429-3436.

48. Brandes, N., D. Ofer, and M. Linial, *ASAP: a machine learning framework for local protein properties*. Database, 2016. **2016**: p. baw133.
49. Xu, M., et al. *Protst: Multi-modality learning of protein sequences and biomedical texts*. in *International Conference on Machine Learning*. 2023. PMLR.
50. Bairoch, A. and R. Apweiler, *The SWISS-PROT protein sequence database and its supplement TrEMBL in 2000*. Nucleic acids research, 2000. **28**(1): p. 45-48.
51. Xu, M., et al., *Protst: Multi-modality learning of protein sequences and biomedical texts*. arXiv preprint arXiv:2301.12040, 2023.
52. Zhou, H.-Y., et al., *Protein Representation Learning via Knowledge Enhanced Primary Structure Modeling*. bioRxiv, 2023: p. 2023-01.
53. Zhang, N., et al., *Ontoprotein: Protein pretraining with gene ontology embedding*. arXiv preprint arXiv:2201.11147, 2022.
54. Outeiral, C. and C.M. Deane, *Codon language embeddings provide strong signals for use in protein engineering*. Nature Machine Intelligence, 2024. **6**(2): p. 170-179.
55. Jarzab, A., et al., *Meltome atlas—thermal proteome stability across the tree of life*. Nature methods, 2020. **17**(5): p. 495-503.
56. Thumuluri, V., et al., *DeepLoc 2.0: multi-label subcellular localization prediction using protein language models*. Nucleic acids research, 2022. **50**(W1): p. W228-W234.
57. Sridharan, S., et al., *Proteome-wide solubility and thermal stability profiling reveals distinct regulatory roles for ATP*. Nature communications, 2019. **10**(1): p. 1155.
58. Unsal, S., et al., *Learning functional properties of proteins with language models*. Nature Machine Intelligence, 2022. **4**(3): p. 227-245.
59. Uhln, M., et al., *Tissue-based map of the human proteome*. Science, 2015. **347**(6220): p. 1260419.
60. Wang, M., et al., *PaxDb, a database of protein abundance averages across all three domains of life*. Molecular \& cellular proteomics, 2012. **11**(8): p. 492-500.
61. Liu, W., et al., *PLMSearch: Protein language model powers accurate and fast sequence search for remote homology*. Nature communications, 2024. **15**(1): p. 2775.
62. Chandonia, J.-M., et al., *SCOPe: improvements to the structural classification of proteins--extended database to facilitate variant interpretation and machine learning*. Nucleic acids research, 2022. **50**(D1): p. D553-D559.
63. Sillitoe, I., et al., *CATH: increased structural coverage of functional space*. Nucleic acids research, 2021. **49**(D1): p. D266-D273.
64. Hong, L., et al., *Fast, sensitive detection of protein homologs using deep dense retrieval*. Nature Biotechnology, 2024: p. 1-13.

65. Wang, F., et al., *MHCRoBERTa: pan-specific peptide-MHC class I binding prediction through transfer learning with label-agnostic protein sequences*. Brief Bioinform, 2022. **23**(3).
66. Boutet, E., et al., *UniProtKB/Swiss-Prot: the manually annotated section of the UniProt KnowledgeBase*, in *Plant bioinformatics: methods and protocols*. 2007, Springer. p. 89-112.
67. Vita, R., et al., *The Immune Epitope Database (IEDB): 2018 update*. Nucleic Acids Res, 2019. **47**(D1): p. D339-D343.
68. Cheng, J., et al., *BERTMHC: improved MHC-peptide class II interaction prediction with transformer and multiple instance learning*. Bioinformatics, 2021. **37**(22): p. 4172-4179.
69. *Improved methods for predicting peptide binding affinity to MHC class II molecules*. 2017.
70. Wu, K., et al., *TCR-BERT: learning the grammar of T-cell receptors for flexible antigenbinding analyses*. 2021.
71. Zhang, W., et al., *PIRD: Pan Immune Repertoire Database*. Bioinformatics, 2020. **36**(3): p. 897-903.
72. Bagaev, D.V., et al., *VDJdb in 2019: database extension, new analysis infrastructure and a T-cell receptor motif compendium*. Nucleic Acids Research, 2020. **48**(D1): p. D1057-D1062.
73. Chen, S.Y., et al., *TCRdb: a comprehensive database for T-cell receptor sequences with powerful search function*. Nucleic Acids Res, 2021. **49**(D1): p. D468-D474.
74. Wang, Q., et al., *AntiFormer: graph enhanced large language model for binding affinity prediction*. Briefings in Bioinformatics, 2024. **25**(5).
75. Olsen, T.H., F. Boyles, and C.M.J.P.S. Deane, *Observed Antibody Space: A diverse database of cleaned, annotated, and translated unpaired and paired antibody sequences*. 2022. **31**(1): p. 141-146.
76. Zhao, Y., et al., *SC-AIR-BERT: a pre-trained single-cell model for predicting the antigen-binding specificity of the adaptive immune receptor*. Brief Bioinform, 2023. **24**(4).
77. Wu, L., et al., *huARdb: human Antigen Receptor database for interactive clonotype-transcriptome analysis at the single-cell level*. Nucleic Acids Res, 2022. **50**(D1): p. D1244-D1254.
78. Raybould, M.I.J., et al., *CoV-AbDab: the coronavirus antibody database*. Bioinformatics, 2021. **37**(5): p. 734-735.
79. Olsen, T.H., I.H. Moal, and C.M. Deane, *AbLang: an antibody language model for completing antibody sequences*. Bioinformatics Advances, 2022. **2**(1): p. vbac046.

80. Liu, Y., et al., *Roberta: A robustly optimized bert pretraining approach*. arXiv preprint arXiv:1907.11692, 2019.
81. Kovaltsuk, A., et al., *Observed antibody space: a resource for data mining next-generation sequencing of antibody repertoires*. The Journal of Immunology, 2018. **201**(8): p. 2502-2509.
82. Ghraichy, M., et al., *Different B cell subpopulations show distinct patterns in their IgH repertoire metrics*. Elife, 2021. **10**: p. e73111.
83. Dunbar, J., et al., *SAbDab: the structural antibody database*. Nucleic acids research, 2014. **42**(D1): p. D1140-D1146.
84. Marks, C., et al., *Humanization of antibodies using a machine learning approach on large-scale repertoire data*. Bioinformatics, 2021. **37**(22): p. 4041-4047.
85. Wang, D., F. Ye, and H. Zhou, *On pre-trained language models for antibody*. bioRxiv, 2023: p. 2023-01.
86. Mason, D.M., et al., *Optimization of therapeutic antibodies by predicting antigen specificity from antibody sequence via deep learning*. Nature Biomedical Engineering, 2021. **5**(6): p. 600-612.
87. Liberis, E., et al., *Parapred: antibody paratope prediction using convolutional and recurrent neural networks*. Bioinformatics, 2018. **34**(17): p. 2944-2950.
88. Mroczek, E.S., et al., *Differences in the composition of the human antibody repertoire by B cell subsets in the blood*. Frontiers in immunology, 2014. **5**: p. 96.

Supplementary Table 3. Detailed information of large language models for drug-discovery tasks

Application area	Models	Ref	Publication time	Parameters	Architecture	Datasets			Downstream tasks
						Data type	Source	Size	
Predictions of Molecular Properties	SMILES-BERT	[1]	Sep 2019	6 Transformer encoder layers, the feed-forward hidden units are 1024 and the attention heads are 4	BERT-based	SMILES	ZINC [2]	18.69 M used, totally more than 35 M	LogP prediction, PM2 prediction, and molecular property prediction
	ChemBERTa	[3]	Oct 2020	Implementation of RoBERTa uses 12 attention heads and 6 layers, resulting in 72 distinct attention mechanisms.	BERT-based	SMILES	PubChem [4]	Curated a dataset of 77 M unique SMILES from PubChem, divided this dataset into subsets of 100 K, 250 K, 1 M, and 10 M.	Binary classification prediction of barrier permeability properties, binary classification of clinical trial toxicity, whether the compound inhibits HIV replication for binary classification between active and inactive
	ChemBERTa-2	[5]	Sep 2022	Implementation of RoBERTa uses 12 attention heads and 6 layers, resulting in 72 distinct attention mechanisms.	BERT-based	SMILES	PubChem	Over a large corpus of 77 million SMILES strings	brain penetrability, toxicity, solubility, and on-target inhibition
	K-BERT	[6]	Apr 2022	The hidden size of the transformer encoder is 768, and the number of the attention heads is 12. Six transformer encoders were used in KBERT.	BERT-based	SMILES	ChEMBL [7]	1.8 M used	Related tasks on 15 drug-discovery-related datasets including carcinogenicity, respiratory toxicity and drug induced liver injury.

	MOLE-BERT	[8]	Apr 2023	A 5-layer Graph Isomorphism Networks (GINs) whose hidden dimension is 300	BERT-style, graph-based	Molecular graphs	ZINC15 [9]	2 million molecules sampled from the dataset	Related tasks on 8 drug-discovery-related datasets including barrier permeability properties and clinical trial toxicity
Generation of Molecules	MolGPT	[10]	Oct 2021	6 M parameters. Each self-attention layer returns a vector of size 256 that is taken as input by the fully connected network. The hidden layer of the neural network outputs a vector of size 1024 and passes it through GELU activation layer. The final layer of the fully connected neural network returns a vector of size 256, that is then used as input for the next decoder block. MolGPT consists of eight such decoder blocks	GPT-based	SMILES	MOSES [11] and GuacaMol [12]	1.9 M, 1.6 M	Generating molecules
Drug-Target Interaction	DTI-BERT	[13]	Jun 2022	The proteins can be represented via 1024-D vectors (dimensionality of the features extracted by the ProtBert model). Drug molecular fingerprints are represented by 128-D vectors through semi decomposition process discrete wavelet transform (DWT). Secondly, the 1152-D vectors (a concatenation of protein sequence feature and drug feature) are fed into the feature extraction model to generate interaction information	BERT-based	Molecular fingerprints and protein sequence pairs	DrugBanks [14], BRENDA, SuperTarget, and KEGG BRITE [15]	4,803 drug-target pairs in positive subsets 9,606 synthesized negative pairs	A pair belongs to an interactive drug-target pair or non-interactive drug-target pair
	TransDTI	[16]	Jan 2022	Consist of SMILES-BERT and fine-tuned large protein models.	BERT-based	Molecular SMILES and protein sequence pairs	KIBA [17], gold-standard external data sets from DTI-MLCD [18]	30,474 compounds, 961 targets and 61,624 interactions	A three-classification based on binding affinity

	C2P2	[19]	Jul 2022	Consist of ChemBERTa and ESM model	BERT-based	Molecular SMILES and protein sequence pairs	STRING [20], STITCH [21], Davis [22], and PDBBind v2019 [23, 24]	Over 67.6 million proteins with over 20 billion protein–protein pairs, over 0.5 million chemicals with over 1.6 billion interactions, 30,056 interactions, 14,011 interactions	Binding affinity prediction
	Hyeunseok Kang et al.	[25]	Aug 2022	Consist of ChemBERTa and ProtBERT model	BERT-based	Molecular SMILES and protein sequence pairs	BIOSNAP [26], DAVIS [22] and BindingDB [27]	27,482 interactions, 11,103 interactions, 32,601 interactions	Binding affinity prediction

References

1. Wang, S., et al. *Smiles-bert: large scale unsupervised pre-training for molecular property prediction.* in *Proceedings of the 10th ACM international conference on bioinformatics, computational biology and health informatics.* 2019.
2. Irwin, J.J., et al., *ZINC: a free tool to discover chemistry for biology.* Journal of chemical information and modeling, 2012. **52**(7): p. 1757-1768.
3. Chithrananda, S., G. Grand, and B. Ramsundar, *ChemBERTa: large-scale self-supervised pretraining for molecular property prediction.* arXiv preprint arXiv:2010.09885, 2020.
4. Kim, S., et al., *PubChem 2019 update: improved access to chemical data.* Nucleic acids research, 2019. **47**(D1): p. D1102-D1109.
5. *ChemBERTa-2: Towards Chemical Foundation Models.*
6. Wu, Z., et al., *Knowledge-based BERT: a method to extract molecular features like computational chemists.* Briefings in Bioinformatics, 2022. **23**(3): p. bbac131.
7. Mendez, D., et al., *ChEMBL: towards direct deposition of bioassay data.* Nucleic acids research, 2019. **47**(D1): p. D930-D940.
8. Xia, J., et al. *Mole-bert: Rethinking pre-training graph neural networks for molecules.* in *The Eleventh International Conference on Learning Representations.* 2022.
9. Sterling, T. and J.J. Irwin, *ZINC 15--ligand discovery for everyone.* Journal of chemical information and modeling, 2015. **55**(11): p. 2324-2337.
10. Bagal, V., et al., *MolGPT: molecular generation using a transformer-decoder model.* Journal of Chemical Information and Modeling, 2021. **62**(9): p. 2064-2076.
11. Polykovskiy, D., et al., *Molecular sets (MOSES): a benchmarking platform for molecular generation models.* Frontiers in pharmacology, 2020. **11**: p. 565644.
12. Brown, N., et al., *GuacaMol: benchmarking models for de novo molecular design.* Journal of chemical information and modeling, 2019. **59**(3): p. 1096-1108.
13. Zheng, J., X. Xiao, and W.-R. Qiu, *DTI-BERT: identifying drug-target interactions in cellular networking based on BERT and deep learning method.* Frontiers in Genetics, 2022. **13**: p. 859188.
14. Wishart, D.S., et al., *DrugBank 5.0: a major update to the DrugBank database for 2018.* Nucleic acids research, 2018. **46**(D1): p. D1074-D1082.
15. Hu, J., et al., *GPCR--drug interactions prediction using random forest with drug-association-matrix-based post-processing procedure.* Computational biology and chemistry, 2016. **60**: p. 59-71.
16. Kalakoti, Y., S. Yadav, and D. Sundar, *TransDTI: transformer-based language models for estimating DTIs and building a drug recommendation workflow.* ACS omega, 2022. **7**(3): p. 2706-2717.
17. Tang, J., et al., *Making sense of large-scale kinase inhibitor bioactivity data sets: a comparative and integrative analysis.* Journal of Chemical Information and Modeling, 2014. **54**(3): p. 735-743.
18. Chu, Y., et al., *DTI-MLCD: predicting drug-target interactions using multi-label learning with community detection method.* Briefings in bioinformatics, 2021. **22**(3):

- p. bbba205.
- 19. Nguyen, T.M., T. Nguyen, and T. Tran, *Mitigating cold-start problems in drug-target affinity prediction with interaction knowledge transferring*. *Briefings in Bioinformatics*, 2022. **23**(4): p. bbac269.
 - 20. Szklarczyk, D., et al., *The STRING database in 2021: customizable protein--protein networks, and functional characterization of user-uploaded gene/measurement sets*. *Nucleic acids research*, 2021. **49**(D1): p. D605-D612.
 - 21. Kuhn, M., et al., *STITCH: interaction networks of chemicals and proteins*. *Nucleic acids research*, 2007. **36**(suppl_1): p. D684-D688.
 - 22. Davis, M.I., et al., *Comprehensive analysis of kinase inhibitor selectivity*. *Nature Biotechnology*, 2011. **29**(11): p. 1046-1051.
 - 23. Wang, R., et al., *The PDBbind database: Collection of binding affinities for protein- ligand complexes with known three-dimensional structures*. *Journal of medicinal chemistry*, 2004. **47**(12): p. 2977-2980.
 - 24. Wang, R., et al., *The PDBbind database: methodologies and updates*. *Journal of medicinal chemistry*, 2005. **48**(12): p. 4111-4119.
 - 25. Kang, H., et al., *Fine-tuning of bert model to accurately predict drug--target interactions*. *Pharmaceutics*, 2022. **14**(8): p. 1710.
 - 26. Zitnik, M., R. Sosic, and J. Leskovec, *BioSNAP Datasets: Stanford biomedical network dataset collection*. Note: <http://snap.stanford.edu/biodata> Cited by, 2018. **5**(1).
 - 27. Liu, T., et al., *BindingDB: a web-accessible database of experimentally determined protein--ligand binding affinities*. *Nucleic acids research*, 2007. **35**(suppl_1): p. D198-D201.

Supplementary Table 4. Detailed information of large language models for single-cell tasks

Application area	Models	Ref	Publication time	Parameters	Architecture	Fine-tuning datasets			Downstream tasks
						Data type	Source	Size	
Single-cell large language models	scBERT	[1]	Sep 2022	six Performer encoder layers and ten heads for each layer, 200 dimensions of gene embedding using gene2vec	BERT-based	scRNA-seq	The Panglao dataset [2]	209 human single-cell datasets comprising 74 tissues with 1,126,580 cells	Cell type annotation, novel cell type discovery, robustness to batch effects and model interpretability
							Zheng68k dataset [3]	68,450 cells	
							Pancreas datasets [4-7]	-	
							MacParland dataset [8]	8,444 cells	
							Heart datasets [9, 10]	451,513 cells for pretraining and the 287,269 cells for benchmarking	
							Lung dataset [11]	39,778 cells	
							Human Cell Atlas dataset [12]	84,363 cells from 27 cell types among 15 major organs	
	scGPT	[13]	Feb 2024	12 stacked transformer blocks with 8 attention heads, 512 embedding sizes of the pre-trained foundation model, 512 hidden sizes of the fully connected layer	GPT-based	scRNA-seq	CELLxGENE scRNA-seq human PBMC Collection [14]	33 million human PBMC scRNA-seq samples	Gene network inference
							PBMC 10K dataset [15]	Two scRNA-seq data of 7,982 cells and 4,008 cells.	Multi-batch integration
							Immune Human dataset [16]	33,506 cells	
							hPancreas dataset [4-7, 17, 18]	10,600 cells in the reference set and 4,218 cells in the query set	Cell type annotation
							Adamson dataset [19]	87 unique one-gene perturbations, each replicated in around 100 cells	Genetic perturbation prediction
							Norman dataset [20]	131 two-gene perturbations and 105 one-gene perturbations. Each perturbation is replicated in around 300-700 cells.	
	CIForm	[23]	July 2023	64 attention heads	BERT-based	scRNA-seq	10X Multiome PBMC [21]	9,631 cells	Multi-omic integration
							ASAP PBMC [22]	Four datasets each contain 5,023, 3,666, 3,517, and 4,849 cells respectively	
							Pancreas datasets [4-7]	-	Cell type annotation
							Immune datasets [24-26]	-	
							Brain datasets [27-29]	-	
							Tabula Muris dataset [30]	Nearly 100000 cells from 20 organs and tissues	

						Zheng68k dataset [3]	68,450 cells	
						ZhangTdataset [31]	8,530 cells from 20 subtypes	
						Allen mouse brain dataset [32]	12832cells	
TOSICA	[17]	Jan 2023	-	BERT-based	scRNA-seq	human pancreas (hPancreas) [4-7]	10,600 cells for training and 4,218 for testing	Cell type annotation, new cell type discovery and batch correction
						human bone (hBone, GSE152805) [33]	14,615 cells for training and 11,525 for testing	
						human artery (hArtery, GSE159677) [34]	10,960 cells for training and 35,399 for testing	
						mouse brain (mBrain) [28-30, 35]	48,801 cells for training and 7,394 for testing	
						mouse pancreas (mPancreas, GSE132188) [36]	25,465 cells for training and 10,886 for testing	
						mouse atlas (mAtlas, GSE132042) [37]	78,672 cells for training and 277,541 for testing	
						human cell atlas dataset	295,805 cells from 35 tissues	
scTransSort	[38]	Mar 2023	12 layers of transformer	BERT-based	scRNA-seq	mouse cell atlas dataset	105,148 cells from 26 tissues and 103,148 cells from 26 tissues	Cell type annotation
						data processed by Shao X et al. [39]	-	
						Human training data [41]	-	
TransCluster	[40]	Oct 2022	5 attention heads	improved Transformer model	scRNA-seq	The Shao dataset [39]	-	Cell type Identification
						The Baron dataset [4]	-	
Geneformer	[42]	May 2023	six transformer encoder units, input size of 2,048, 256 embedding dimensions, four attention heads per layer and feed forward size of 512	BERT-based	scRNA-seq	iPSC differentiation data	Assayed in parallel on the Dropseq (single cell) or DroNc-seq (single nucleus) platform	Batch effect removal, cell type annotation
						Huggingface Dataset [43]	a largescale pretraining corpus, Genecorpus-30M, comprising 29.9 million human single-cell transcriptomes	discover key network regulators and candidate therapeutic targets
tGPT	[44]	May 2023	8 transformer decoder blocks with 1024 hidden units and 16 attention heads	GPT-based	scRNA-seq	Human Cell Atlas Census of Immune Cells (HCA) [45]	282,588 Bone marrow cells from 64 healthy donors in Human Cell Atlas (HCA) project	single-cell clustering, inference of developmental lineage
						Human cell Landscape (HCL) [41]	586,135 human cells	
						Tabula Mursi dataset [30]	54,865 cells	
						Macaque Retina dataset [46]	124,965 cells	

	DeepMAPS	[49]	Feb 2023	-	BERT-based, graph transformer	Bulk RNA-seq	The Cancer Genome Atlas (TCGA) [47]	9,318 bulk samples	interrogation of feature representation of bulk tissues in relation to genomic alterations, prognosis and treatment response of immunotherapy
							Genotype-Tissue Expression Project (GTEx) [48]	11,688 bulk samples	
							Multiple scRNA-seq data [4, 16]	Three scRNA-seq datasets with 20,125 cells, 14,878 cells and 16,382 cells	
						Single-cell multi-omics	CITE-seq data [50]	Three CITE-seq datasets with 25,171 cells, 32,029 cells and 16,750 cells	
	scMVP	[51]	Jan 2022	8 self-attention heads and each head takes 16-dimension feature	Transformer-based		scRNA-ATAC-seq data	Four scRNA-ATAC-seq datasets with 3,009 cells, 11,898 cells, 3,233 cells and 10,970 cells	Cell clustering, infer cell-type-specific biological networks from scMulti-omics data
					Paired scRNA-seq and scATAC-seq data	sci-CAR cell line dataset [52]	293T cell line, 3T3 cell line, 293T/3T3 cell mixture, and A549 cell line treated with dexamethasone (DEX) for 0 h, 1 h, and 3 h		
						Paired-seq cell line dataset [53]	derived from HEK293, HepG2, and their cell line mixture		
						SNARE-seq cell line dataset [54]	5,081 cells		
						SHARE-seq [55]	67,418 cells		
	scTranslator	[56]	BioRxiv posted July 2023	8-headed attention mechanism in each sub-layer, 117 million parameters.	Transformer-based	Bulk datasets	The Cancer Genome Atlas (TCGA) data [57-60]	31 cancer types and 18,227 samples in total	translate single-cell transcriptome to proteome, predict protein abundance
							The data from Clinical Proteomic Tumor Analysis Consortium (CPTAC) [61-67]		
							The dataset from Broad Institute [68, 69]		
							The dataset from Memorial Sloan Kettering Cancer Center (MSKCC) [70]		
							Single-cell	The Seurat v4 PBMCs dataset	
								161,764 human peripheral blood	

						datasets	[71]	mononuclear cells	
							The REAP-seq PBMCs dataset [72]	4,330 PBMCs with simultaneous measurements of 44 proteins and 21,005 transcriptome genes	
							The CITE-seq CBMCs dataset [73]	simultaneous measurements of 13 cellular proteins and 16,508 transcriptome genes. This dataset includes 8,005 cord blood mononuclear cells (CBMCs)	
scFoundation	[75]	June 2024	100 million parameters	BERT-based	Single-cell datasets	The single-cell pan-cancer dataset [74]	65,698 myeloid cells with only single-cell transcriptome data, involving 15,844 genes	gene expression enhancement, tissue drug response prediction, cell clustering, single-cell drug response classification, and single-cell perturbation prediction	
							Baron dataset [4]	-	
						Zheng68K dataset [3]	68,450 cells		
						Cancer drug response dataset [76]	-		
						Single cell drug response classification dataset	-		
						Perturbation dataset [77]	-		
scMoFormer	[78]	Oct 2023	Graph transformer	Transformer-based	Single-cell multi-omics	Joint measurements of gene expression and surface protein levels datasets from the NeurIPS multimodal single-cell integration competition of the year 2021 [25] and 2022	-	use gene expression (RNA) to predict surface protein level, protein levels to gene expression, gene expression to chromatin accessibility and chromatin accessibility to gene expression	
GeneCompass	[79]	Sep 2024	12-layer transformer framework, 100 million parameters	BERT-based	single-cell transcriptomes	CHIP-Atlas related to PBMC cells on GSE43036 [80]	-	Gene embedding analysis, Gene expression profiling prediction	
						human multiple sclerosis	-	cell type annotation	

							(hMS), lung (hLung) and liver (hLiver) datasets, and mouse brain (mBrain), lung (mLung) and pancreas (mPancreas) datasets.		
							Immune Human [16]	33,506 cells	GRN inference
							dataset provided by Srivatsan et al. [81]	-	Drug dose response prediction
							predefined dosage-sensitive and nonsensitive gene datasets [42]	-	Gene dosage sensitivity predictions.
							Norman dataset [20]	131 two-gene perturbations and 105 one-gene perturbations. Each perturbation is replicated in around 300-700 cells.	In silico perturbation
							mouse embryonic stem cells (ESCs) [82]	-	In silico quantitative perturbation
scMulan	[83]	BioRxiv posted Jan 2024	24-layer transformer, 368 million parameters	GPT-based	single-cell transcriptomic data	AHCA_BoneMarrow [12]	3,000 bone marrow cells from a single adult donor	zero-shot cell type annotation	
						Simonson2023 [84]	60,345 cells from 8 human left ventricle samples		
						Suo2022 [85]	140,000 liver cells from 14 fetal donors at various developmental stages		
						Intestine_HCL_55k dataset [41]	55,214 intestinal cells across 24 cell types		
						Immune cell dataset [86]	274,346 cells spanning 18 batches	batch integration	
						Lung dataset [16]	32,472 cells from 16 donors		
						six organs within the hECA-10M dataset [83]	3,000 conditions	conditional cell generation	
UCE	[87]	BioRxiv posted Nov 2023	a 33-layer model consisting of over 650 million parameters	BERT-based	single-cell gene expression dataset	Tabula Sapiens v2 dataset	human data from 581,430 cells, 27 tissues, batches and 162 unique cell types	zero-shot embedding of new datasets	
						a dataset of green monkey lymph node and lung cells [88]	17 cell types	Cell type embedding for new species	
						naked mole rat spleen and circulating immune cells [89]	24 cell types		
						two distinct chicken datasets (chick retina [90] and	15 cell types in chicken heart dataset		

							developing chick heart [91])		
							mouse renal cells [92]		decode the function of newly discovered cell types
CellLM	[93]	arXiv posted June 2023	Performer model, 10 layers, 16 attention heads, over 50 million parameters	BERT-based	scRNA-seq data	Zheng68k dataset [3]	68,450 human peripheral blood mononuclear cells (PBMCs) with 11 highly related cell types	Cell type annotation	
						The pancreas Baron dataset [4]	8,562 cells categorized into 13 different cell types		
						Human lung cancer cells (GSE149383) [94]	2,739 cells	Single-cell drug sensitivity prediction	
						human oral squamous cancer cells (GSE117872) [95-97]	1,302 cells		
						Single-omics cell line data	cell lines integrating from CCLE [98] and GDSC [99]	Single-omics cell line drug sensitivity prediction	
					Paired scRNA and scATAC-seq data	Fetal Atlas [101, 102]	377, 134 cells	integration of multi-modal single-cell sequencing data	
scCLIP	[100]	Oct 2023	Vanilla transformers	transformer-based encoders		Brain [103]	-		
iSEEK	[104]	Jan 2022	8 transformer layers each with 576 hidden units and 8 attention heads	BERT-based	single-cell expression data	PBMCs [105]	43,073 cells	Cell clustering	
						Human Cell Atlas Census of Immune Cells (HCA) [45]	282,588 Bone marrow cells from 64 healthy donors in Human Cell Atlas (HCA) project		
						Tabula Muris dataset [30]	Nearly 100000 cells from 20 organs and tissues		
						Zheng68k dataset [3]	68,579 cells		
						the dataset of FACS-sorted CD4+/8+ T cells [31, 106, 107]	12,670 CD4+ and 9,012 CD8+ T cells	Identify gene–gene interaction networks	
					scRNA-seq and spatially-resolved transcriptomic (SRT) data	dataset from Li et al. [109]	48, 082 cells	zero-shot clustering	
CellPLM	[108]	BioRxiv posted Oct 2023	over 80 million parameters	BERT-based		PBMC 5K and Jurkat from 10x Genomics	33,538 cells in PBMC 5K and 32,738 cells in Jurkat	scRNA-seq denoising	
						two spatial transcriptomic datasets at single-cell resolution, i.e., Lung2 and Liver2	836,739 cells in Lung2 and 598,141 cells in Liver2	spatial transcriptomic imputation	
						hPancreas [17] and Multiple Sclerosis (MS) [110]	-	cell type annotation	
						the Adamson Perturb-Seq	87 one-gene perturbations in the	perturbation	

							dataset [19] and the Norman Perturb-Seq dataset [20]	Adamson Perturb-Seq dataset, and 131 two-gene perturbations and 105 one-gene perturbations in the Norman Perturb-Seq dataset	prediction
scGREAT	[111]	Feb 2024	-	Transformer-based	single-cell transcriptomics	human embryonic stem cells (hESC) (GSE75748) [112] human mature hepatocytes (hHEP) (GSE81252) [113, 114] mouse dendritic cells (mDC) (GSE48968) [115] mouse embryonic stem cells (mESC) (GSE98664) [116] mouse hematopoietic stem cells with erythroid-lineage (mHSC-E) mouse hematopoietic stem cells with granulocyte-monocyte-lineage (mHSC-GM) mouse hematopoietic stem cells with lymphoid-lineage (mHSC-L) (GSE81682) [117]	- - - - -	gene regulatory network inference	
BioFormers	[118]	bioRxiv posted Dec 2023	An 8-layer transformer encoder model with 8 self-attention heads per layer and a hidden state dimension of 512	BERT-based	scRNA-seq	PBMC [15] PBMC 4k and 8k [119] Perturb-seq dataset [19]	7,982 cells and 3,346 genes 11,990 cells and 2,000 HVGs 87 single-gene perturbations, with ~100 cells per perturbation and a control set of at least 7,000 unperturbed cells	Cell clustering and identification Gene expression prediction genetic perturbation prediction and gene network inference	
scPRINT	[120]	bioRxiv posted July 2024	2M to 100M parameters, 4-layer transformer encoder model with 2 self-attention heads per layer	BERT-based	scRNA-seq	three test datasets of kidney, retina, and colon tissues [121-123] perturb-seq [124] and ChIP-seq [125] 3 test datasets of ciliary body, colon, and retina [122, 126, 127] premalignant neoplasms from	comprising 26 cell types - - -	gene network inference, denoising, batch effect correction, and cell label prediction	

						human prostate tissues [128]			
ScRAT	[129]	Feb 2024	one-layer transformer encoder model with 8 self-attention heads	Transformer-based	scRNA-seq	COMBAT [130] and Haniffa datasets [131]	835,937 and 528,438 cells	disease diagnosis	
						SC4 with COVID samples [132]	501,943 cells for severity prediction 1,289,496 cells for stage prediction	predict severity and stage	
CancerFoundation	[133]	bioRxiv posted Nov 2024	Six transformer layers, 10.8 million parameters	BERT-based	scRNA-seq	glioblastoma dataset [134]	four distinct malignant cell states	batch integration	
						CCLE [98] and GDSC [99]	-	Drug response prediction	
					bulk RNA-seq data	TCGA data [135, 136]	21 cancer types	Survival prediction	
mcBERT	[137]	bioRxiv posted Nov 2024	12 blocks, each with 12 attention heads	BERT-based	scRNA-seq	Heart [9, 84, 138-141]	Refer to Table 1 in [137]	patient-level representation, disease clustering, phenotypical interpretation/batch effect removal	
						Kidney [142-146]			
Nicheformer	[150]	bioRxiv posted Oct 2024	The transformer block leverages 12 transformer encoder units with 16 attention heads per layer and a feed-forward network size of 1,024 to generate a 512-dimensional embedding of the pretraining dataset, resulting in altogether 49.3M parameters.	BERT-based		PBMC [130, 147, 148]			
						Lung [149]			
						MERFISH mouse brain [151]	4.3 million cells across 59 tissue sections	Label prediction (cell type, niche and region labels)	
						CosMx human liver [152]	332,877 healthy cells and 460,441 cancer cells	niche label prediction (healthy data only), niche composition prediction	
						CosMx human lung [152]	five different donors (301,611, 89,975, 227,110, 71,304 and 81,236 cells, respectively)	niche composition prediction	
SpaFormer	[153]	arXiv posted Feb 2023	2 layers and 8 heads	Transformer-based	Spatial transcriptomics data	Xenium human lung from 10X Genomics	295,883 healthy cells and 531,165 cancer cells	neighborhood density prediction	
						Xenium human colon from 10X Genomics	275,822 healthy cells and 587,115 cancer cells		
						Lung 5 data generated by the CosMX platform [152]	99,656 cells		
						Kidney 1139 data generated by the CosMX platform [152]	61,283 cells	spatial transcriptomic imputation	
						Liver normal generated by the CosMX platform [152]	305,730 cells		

References

1. Yang, F., et al., *scBERT as a large-scale pretrained deep language model for cell type annotation of single-cell RNA-seq data*. Nature Machine Intelligence, 2022. **4**(10): p. 852-866.
2. Franzen, O., L.M. Gan, and J.L.M. Bjorkegren, *PanglaoDB: a web server for exploration of mouse and human single-cell RNA sequencing data*. Database (Oxford), 2019. **2019**.
3. Zheng, G.X., et al., *Massively parallel digital transcriptional profiling of single cells*. Nat Commun, 2017. **8**: p. 14049.
4. Baron, M., et al., *A Single-Cell Transcriptomic Map of the Human and Mouse Pancreas Reveals Inter- and Intra-cell Population Structure*. Cell Syst, 2016. **3**(4): p. 346-360 e4.
5. Muraro, M.J., et al., *A Single-Cell Transcriptome Atlas of the Human Pancreas*. Cell Syst, 2016. **3**(4): p. 385-394 e3.
6. Segerstolpe, A., et al., *Single-Cell Transcriptome Profiling of Human Pancreatic Islets in Health and Type 2 Diabetes*. Cell Metab, 2016. **24**(4): p. 593-607.
7. Xin, Y., et al., *RNA Sequencing of Single Human Islet Cells Reveals Type 2 Diabetes Genes*. Cell Metab, 2016. **24**(4): p. 608-615.
8. MacParland, S.A., et al., *Single cell RNA sequencing of human liver reveals distinct intrahepatic macrophage populations*. Nat Commun, 2018. **9**(1): p. 4383.
9. Litvinukova, M., et al., *Cells of the adult human heart*. Nature, 2020. **588**(7838): p. 466-472.
10. Tucker, N.R., et al., *Transcriptional and Cellular Diversity of the Human Heart*. Circulation, 2020. **142**(5): p. 466-482.
11. Lukassen, S., et al., *SARS-CoV-2 receptor ACE2 and TMPRSS2 are primarily expressed in bronchial transient secretory cells*. EMBO J, 2020. **39**(10): p. e105114.
12. He, S., et al., *Single-cell transcriptome profiling of an adult human cell atlas of 15 major organs*. Genome Biol, 2020. **21**(1): p. 294.
13. Cui, H., et al., *scGPT: toward building a foundation model for single-cell multi-omics using generative AI*. Nat Methods, 2024. **21**(8): p. 1470-1480.
14. (n.d.), C.Z.I., *CZ CELLxGENE Discover*. 2022, <https://cellxgene.cziscience.com/>.
15. Gayoso, A., et al., *A Python library for probabilistic analysis of single-cell omics data*. Nat Biotechnol, 2022. **40**(2): p. 163-166.
16. Luecken, M.D., et al., *Benchmarking atlas-level data integration in single-cell genomics*. Nat Methods, 2022. **19**(1): p. 41-50.
17. Chen, J., et al., *Transformer for one stop interpretable cell type annotation*. Nat Commun, 2023. **14**(1): p. 223.

18. Lawlor, N., et al., *Single-cell transcriptomes identify human islet cell signatures and reveal cell-type-specific expression changes in type 2 diabetes*. *Genome Res*, 2017. **27**(2): p. 208-222.
19. Adamson, B., et al., *A Multiplexed Single-Cell CRISPR Screening Platform Enables Systematic Dissection of the Unfolded Protein Response*. *Cell*, 2016. **167**(7): p. 1867-1882 e21.
20. Norman, T.M., et al., *Exploring genetic interaction manifolds constructed from rich single-cell phenotypes*. *Science*, 2019. **365**(6455): p. 786-793.
21. Cusanovich, D.A., et al., *Multiplex single cell profiling of chromatin accessibility by combinatorial cellular indexing*. *Science*, 2015. **348**(6237): p. 910-4.
22. Mimitou, E.P., et al., *Scalable, multimodal profiling of chromatin accessibility, gene expression and protein levels in single cells*. *Nat Biotechnol*, 2021. **39**(10): p. 1246-1258.
23. Xu, J., et al., *CIForm as a Transformer-based model for cell-type annotation of large-scale single-cell RNA-seq data*. *Brief Bioinform*, 2023. **24**(4).
24. Sun, Z., et al., *A Bayesian mixture model for clustering droplet-based single-cell transcriptomic data from population studies*. *Nat Commun*, 2019. **10**(1): p. 1649.
25. Oetjen, K.A., et al., *Human bone marrow assessment by single-cell RNA sequencing, mass cytometry, and flow cytometry*. *JCI Insight*, 2018. **3**(23).
26. Dahlin, J.S., et al., *A single-cell hematopoietic landscape resolves 8 lineage trajectories and defects in Kit mutant mice*. *Blood*, 2018. **131**(21): p. e1-e11.
27. Zeisel, A., et al., *Molecular Architecture of the Mouse Nervous System*. *Cell*, 2018. **174**(4): p. 999-1014 e22.
28. Saunders, A., et al., *Molecular Diversity and Specializations among the Cells of the Adult Mouse Brain*. *Cell*, 2018. **174**(4): p. 1015-1030 e16.
29. Rosenberg, A.B., et al., *Single-cell profiling of the developing mouse brain and spinal cord with split-pool barcoding*. *Science*, 2018. **360**(6385): p. 176-182.
30. Tabula Muris, C., et al., *Single-cell transcriptomics of 20 mouse organs creates a Tabula Muris*. *Nature*, 2018. **562**(7727): p. 367-372.
31. Zhang, L., et al., *Lineage tracking reveals dynamic relationships of T cells in colorectal cancer*. *Nature*, 2018. **564**(7735): p. 268-272.

32. Tasic, B., et al., *Shared and distinct transcriptomic cell types across neocortical areas*. Nature, 2018. **563**(7729): p. 72-78.
33. Chou, C.H., et al., *Synovial cell cross-talk with cartilage plays a major role in the pathogenesis of osteoarthritis*. Sci Rep, 2020. **10**(1): p. 10868.
34. Alsaigh, T., et al., *Decoding the transcriptome of calcified atherosclerotic plaque at single-cell resolution*. Commun Biol, 2022. **5**(1): p. 1084.
35. Zeisel, A., et al., *Brain structure. Cell types in the mouse cortex and hippocampus revealed by single-cell RNA-seq*. Science, 2015. **347**(6226): p. 1138-42.
36. Bastidas-Ponce, A., et al., *Comprehensive single cell mRNA profiling reveals a detailed roadmap for pancreatic endocrinogenesis*. Development, 2019. **146**(12).
37. Tabula Muris, C., *A single-cell transcriptomic atlas characterizes ageing tissues in the mouse*. Nature, 2020. **583**(7817): p. 590-595.
38. Jiao, L., et al., *scTransSort: Transformers for Intelligent Annotation of Cell Types by Gene Embeddings*. Biomolecules, 2023. **13**(4).
39. Shao, X., et al., *scDeepSort: a pre-trained cell-type annotation method for single-cell transcriptomics using deep learning with a weighted graph neural network*. Nucleic Acids Res, 2021. **49**(21): p. e122.
40. Song, T., et al., *TransCluster: A Cell-Type Identification Method for single-cell RNA-Seq data using deep learning based on transformer*. Front Genet, 2022. **13**: p. 1038919.
41. Han, X., et al., *Construction of a human cell landscape at single-cell level*. Nature, 2020. **581**(7808): p. 303-309.
42. Theodoris, C.V., et al., *Transfer learning enables predictions in network biology*. Nature, 2023. **618**(7965): p. 616-624.
43. Lhoest, Q., et al., *Datasets: A community library for natural language processing*. arXiv preprint arXiv:2109.02846, 2021.
44. Shen, H., et al., *Generative pretraining from large-scale transcriptomes for single-cell deciphering*. iScience, 2023. **26**(5): p. 106536.
45. Regev, A., et al., *The human cell atlas white paper*. arXiv preprint arXiv:1810.05192, 2018.
46. Peng, Y.R., et al., *Molecular Classification and Comparative Taxonomies of Foveal and Peripheral Cells in Primate Retina*. Cell, 2019. **176**(5): p. 1222-1237 e22.
47. Johnson, W.E., C. Li, and A. Rabinovic, *Adjusting batch effects in microarray expression data using empirical Bayes methods*. Biostatistics, 2007. **8**(1): p. 118-27.
48. Huang, T.X. and L. Fu, *The immune landscape of esophageal cancer*. Cancer Commun (Lond), 2019. **39**(1): p. 79.

49. Ma, A., et al., *Single-cell biological network inference using a heterogeneous graph transformer*. Nat Commun, 2023. **14**(1): p. 964.
50. Luecken, M.D., et al. *A sandbox for prediction and integration of DNA, RNA, and proteins in single cells*. in *Thirty-fifth conference on neural information processing systems datasets and benchmarks track (Round 2)*. 2021.
51. Li, G., et al., *A deep generative model for multi-view profiling of single-cell RNA-seq and ATAC-seq data*. Genome Biol, 2022. **23**(1): p. 20.
52. Cao, J., et al., *Joint profiling of chromatin accessibility and gene expression in thousands of single cells*. Science, 2018. **361**(6409): p. 1380-1385.
53. Zhu, C., et al., *An ultra high-throughput method for single-cell joint analysis of open chromatin and transcriptome*. Nat Struct Mol Biol, 2019. **26**(11): p. 1063-1070.
54. Chen, S., B.B. Lake, and K. Zhang, *High-throughput sequencing of the transcriptome and chromatin accessibility in the same cell*. Nat Biotechnol, 2019. **37**(12): p. 1452-1457.
55. Ma, S., et al., *Chromatin Potential Identified by Shared Single-Cell Profiling of RNA and Chromatin*. Cell, 2020. **183**(4): p. 1103-1116 e20.
56. Linjing, L., et al., *A pre-trained large language model for translating single-cell transcriptome to proteome*. bioRxiv, 2023: p. 2023.07.04.547619.
57. Cancer Genome Atlas Research, N., *Comprehensive molecular characterization of clear cell renal cell carcinoma*. Nature, 2013. **499**(7456): p. 43-9.
58. Ciriello, G., et al., *Comprehensive Molecular Portraits of Invasive Lobular Breast Cancer*. Cell, 2015. **163**(2): p. 506-19.
59. Fishbein, L., et al., *Comprehensive Molecular Characterization of Pheochromocytoma and Paraganglioma*. Cancer Cell, 2017. **31**(2): p. 181-193.
60. Kahles, A., et al., *Comprehensive Analysis of Alternative Splicing Across Tumors from 8,705 Patients*. Cancer Cell, 2018. **34**(2): p. 211-224 e6.
61. Cao, L., et al., *Proteogenomic characterization of pancreatic ductal adenocarcinoma*. Cell, 2021. **184**(19): p. 5031-5052 e26.
62. Dou, Y., et al., *Proteogenomic Characterization of Endometrial Carcinoma*. Cell, 2020. **180**(4): p. 729-748 e26.
63. Gillette, M.A., et al., *Proteogenomic Characterization Reveals Therapeutic Vulnerabilities in Lung Adenocarcinoma*. Cell, 2020. **182**(1):

- p. 200-225 e35.
- 64. Krug, K., et al., *Proteogenomic Landscape of Breast Cancer Tumorigenesis and Targeted Therapy*. Cell, 2020. **183**(5): p. 1436-1456 e31.
 - 65. Petralia, F., et al., *Integrated Proteogenomic Characterization across Major Histological Types of Pediatric Brain Cancer*. Cell, 2020. **183**(7): p. 1962-1985 e31.
 - 66. Satpathy, S., et al., *A proteogenomic portrait of lung squamous cell carcinoma*. Cell, 2021. **184**(16): p. 4348-4371 e40.
 - 67. Wang, L.B., et al., *Proteogenomic and metabolomic characterization of human glioblastoma*. Cancer Cell, 2021. **39**(4): p. 509-528 e20.
 - 68. Encyclopedia, T.C.C.L., et al., *Consistency of drug profiles and predictors in large-scale cancer cell line data*. Nature, 2015. **528**(7580): p. 84.
 - 69. Nusinow, D.P., et al., *Quantitative Proteomics of the Cancer Cell Line Encyclopedia*. Cell, 2020. **180**(2): p. 387-402 e16.
 - 70. Pietzak, E.J., et al., *Genomic Differences Between "Primary" and "Secondary" Muscle-invasive Bladder Cancer as a Basis for Disparate Outcomes to Cisplatin-based Neoadjuvant Chemotherapy*. Eur Urol, 2019. **75**(2): p. 231-239.
 - 71. Hao, Y., et al., *Integrated analysis of multimodal single-cell data*. Cell, 2021. **184**(13): p. 3573-3587 e29.
 - 72. Peterson, V.M., et al., *Multiplexed quantification of proteins and transcripts in single cells*. Nat Biotechnol, 2017. **35**(10): p. 936-939.
 - 73. Stoeckius, M., et al., *Simultaneous epitope and transcriptome measurement in single cells*. Nat Methods, 2017. **14**(9): p. 865-868.
 - 74. Cheng, S., et al., *A pan-cancer single-cell transcriptional atlas of tumor infiltrating myeloid cells*. Cell, 2021. **184**(3): p. 792-809 e23.
 - 75. Hao, M., et al., *Large-scale foundation model on single-cell transcriptomics*. Nat Methods, 2024. **21**(8): p. 1481-1491.
 - 76. Liu, Q., et al., *DeepCDR: a hybrid graph convolutional network for predicting cancer drug response*. Bioinformatics, 2020. **36**(Suppl_2): p. i911-i918.
 - 77. Roohani, Y., K. Huang, and J. Leskovec, *Predicting transcriptional outcomes of novel multigene perturbations with GEARS*. Nat Biotechnol, 2023.
 - 78. Tang, W., et al. *Single-cell multimodal prediction via transformers*. in *Proceedings of the 32nd ACM International Conference on Information and Knowledge Management*. 2023.
 - 79. Yang, X., et al., *GeneCompass: deciphering universal gene regulatory mechanisms with a knowledge-informed cross-species foundation model*. Cell Res, 2024. **34**(12): p. 830-845.
 - 80. Qiao, Y., et al., *Synergistic activation of inflammatory cytokine genes by interferon-gamma-induced chromatin remodeling and toll-like*

- receptor signaling.* Immunity, 2013. **39**(3): p. 454-69.
- 81. Srivatsan, S.R., et al., *Massively multiplex chemical transcriptomics at single-cell resolution.* Science, 2020. **367**(6473): p. 45-51.
 - 82. Garipler, G., et al., *The BTB transcription factors ZBTB11 and ZFP131 maintain pluripotency by repressing pro-differentiation genes.* Cell Rep, 2022. **38**(11): p. 110524.
 - 83. Bian, H., et al. *scMulan: a multitask generative pre-trained language model for single-cell analysis.* in *International Conference on Research in Computational Molecular Biology*. 2024. Springer.
 - 84. Simonson, B., et al., *Single-nucleus RNA sequencing in ischemic cardiomyopathy reveals common transcriptional profile underlying end-stage heart failure.* Cell Rep, 2023. **42**(2): p. 112086.
 - 85. Suo, C., et al., *Mapping the developing human immune system across organs.* Science, 2022. **376**(6597): p. eabo0510.
 - 86. Lotfollahi, M., et al., *Mapping single-cell data to reference atlases by transfer learning.* Nat Biotechnol, 2022. **40**(1): p. 121-130.
 - 87. Rosen, Y., et al., *Universal cell embeddings: A foundation model for cell biology.* bioRxiv, 2023: p. 2023.11.28.568918.
 - 88. Dominguez Conde, C., et al., *Cross-tissue immune cell analysis reveals tissue-specific features in humans.* Science, 2022. **376**(6594): p. eabl5197.
 - 89. Speranza, E., et al., *Single-cell RNA sequencing reveals SARS-CoV-2 infection dynamics in lungs of African green monkeys.* Sci Transl Med, 2021. **13**(578).
 - 90. Hilton, H.G., et al., *Single-cell transcriptomics of the naked mole-rat reveals unexpected features of mammalian immunity.* PLoS Biol, 2019. **17**(11): p. e3000528.
 - 91. Yamagata, M., W. Yan, and J.R. Sanes, *A cell atlas of the chick retina based on single-cell transcriptomics.* Elife, 2021. **10**.
 - 92. Orozco, L.D., et al., *Integration of eQTL and a Single-Cell Atlas in the Human Eye Identifies Causal Genes for Age-Related Macular Degeneration.* Cell Rep, 2020. **30**(4): p. 1246-1259 e6.
 - 93. Zhao, S., J. Zhang, and Z. Nie, *Large-scale cell representation learning via divide-and-conquer contrastive learning.* arXiv preprint arXiv:2306.04371, 2023.
 - 94. Aissa, A.F., et al., *Single-cell transcriptional changes associated with drug tolerance and response to combination therapies in cancer.* Nat Commun, 2021. **12**(1): p. 1628.
 - 95. Sharma, A., et al., *Longitudinal single-cell RNA sequencing of patient-derived primary cells reveals drug-induced infidelity in stem cell*

- hierarchy.* Nat Commun, 2018. **9**(1): p. 4931.
- 96. Ravasio, A., et al., *Single-cell analysis of EphA clustering phenotypes to probe cancer cell heterogeneity.* Commun Biol, 2020. **3**(1): p. 429.
 - 97. Suphavilai, C., et al., *Predicting heterogeneity in clone-specific therapeutic vulnerabilities using single-cell transcriptomic signatures.* Genome Med, 2021. **13**(1): p. 189.
 - 98. Barretina, J., et al., *The Cancer Cell Line Encyclopedia enables predictive modelling of anticancer drug sensitivity.* Nature, 2012. **483**(7391): p. 603-7.
 - 99. Iorio, F., et al., *A Landscape of Pharmacogenomic Interactions in Cancer.* Cell, 2016. **166**(3): p. 740-754.
 - 100. Xiong, L., T. Chen, and M. Kellis. *scCLIP: Multi-modal Single-cell Contrastive Learning Integration Pre-training.* in NeurIPS 2023 AI for Science Workshop.
 - 101. Cao, J., et al., *A human cell atlas of fetal gene expression.* Science, 2020. **370**(6518).
 - 102. Domcke, S., et al., *A human cell atlas of fetal chromatin accessibility.* Science, 2020. **370**(6518).
 - 103. Anderson, A.G., et al., *Single nucleus multiomics identifies ZEB1 and MAFB as candidate regulators of Alzheimer's disease-specific cis-regulatory elements.* Cell Genom, 2023. **3**(3): p. 100263.
 - 104. Shen, H., et al., *A universal approach for integrating super large-scale single-cell transcriptomes by exploring gene rankings.* Brief Bioinform, 2022. **23**(2).
 - 105. Kang, H.M., et al., *Multiplexed droplet single-cell RNA-sequencing using natural genetic variation.* Nat Biotechnol, 2018. **36**(1): p. 89-94.
 - 106. Guo, X., et al., *Global characterization of T cells in non-small-cell lung cancer by single-cell sequencing.* Nat Med, 2018. **24**(7): p. 978-985.
 - 107. Zheng, C., et al., *Landscape of Infiltrating T Cells in Liver Cancer Revealed by Single-Cell Sequencing.* Cell, 2017. **169**(7): p. 1342-1356 e16.
 - 108. Wen, H., et al., *CellPLM: pre-training of cell language model beyond single cells.* bioRxiv, 2023: p. 2023.10.03.560734.
 - 109. Li, Y., et al., *Single-Cell Transcriptome Analysis Reveals Dynamic Cell Populations and Differential Gene Expression Patterns in Control and Aneurysmal Human Aortic Tissue.* Circulation, 2020. **142**(14): p. 1374-1388.

110. Schirmer, L., et al., *Neuronal vulnerability and multilineage diversity in multiple sclerosis*. Nature, 2019. **573**(7772): p. 75-82.
111. Wang, Y., et al., *scGREAT: Transformer-based deep-language model for gene regulatory network inference from single-cell transcriptomics*. iScience, 2024. **27**(4): p. 109352.
112. Chu, L.F., et al., *Single-cell RNA-seq reveals novel regulators of human embryonic stem cell differentiation to definitive endoderm*. Genome Biol, 2016. **17**(1): p. 173.
113. Mora-Bermudez, F., et al., *Differences and similarities between human and chimpanzee neural progenitors during cerebral cortex development*. Elife, 2016. **5**.
114. Camp, J.G., et al., *Multilineage communication regulates human liver bud development from pluripotency*. Nature, 2017. **546**(7659): p. 533-538.
115. Shalek, A.K., et al., *Single-cell RNA-seq reveals dynamic paracrine control of cellular variation*. Nature, 2014. **510**(7505): p. 363-9.
116. Hayashi, T., et al., *Single-cell full-length total RNA sequencing uncovers dynamics of recursive splicing and enhancer RNAs*. Nat Commun, 2018. **9**(1): p. 619.
117. Nestorowa, S., et al., *A single-cell resolution map of mouse hematopoietic stem and progenitor cell differentiation*. Blood, 2016. **128**(8): p. e20-31.
118. Amara-Belgadi, S., et al., *BIOFORMERS: A SCALABLE FRAMEWORK FOR EXPLORING BIOSTATES USING TRANSFORMERS*. bioRxiv, 2023: p. 2023.11. 29.569320.
119. Zheng, G.X., et al., *Massively parallel digital transcriptional profiling of single cells*. Nature communications, 2017. **8**(1): p. 14049.
120. Kalfon, J., et al., *scPRINT: pre-training on 50 million cells allows robust gene network predictions*. bioRxiv, 2024: p. 2024.07. 29.605556.
121. Kong, L., et al., *The landscape of immune dysregulation in Crohn's disease revealed through single-cell transcriptomic profiling in the ileum and colon*. Immunity, 2023. **56**(2): p. 444-458 e5.
122. Wang, S.K., et al., *Single-cell multiome of the human retina and deep learning nominate causal variants in complex eye diseases*. Cell Genom, 2022. **2**(8).
123. Marshall, J.L., et al., *High-resolution Slide-seqV2 spatial transcriptomics enables discovery of disease-specific cell neighborhoods and pathways*. iScience, 2022. **25**(4): p. 104097.

124. Dixit, A., et al., *Perturb-Seq: Dissecting Molecular Circuits with Scalable Single-Cell RNA Profiling of Pooled Genetic Screens*. Cell, 2016. **167**(7): p. 1853-1866 e17.
125. Park, P.J., *ChIP-seq: advantages and challenges of a maturing technology*. Nat Rev Genet, 2009. **10**(10): p. 669-80.
126. Burclaff, J., et al., *A Proximal-to-Distal Survey of Healthy Adult Human Small Intestine and Colon Epithelium by Single-Cell Transcriptomics*. Cell Mol Gastroenterol Hepatol, 2022. **13**(5): p. 1554-1589.
127. van Zyl, T., et al., *Cell atlas of the human ocular anterior segment: Tissue-specific and shared cell types*. Proc Natl Acad Sci U S A, 2022. **119**(29): p. e2200914119.
128. Joseph, D.B., et al., *Single-cell analysis of mouse and human prostate reveals novel fibroblasts with specialized distribution and microenvironment interactions*. J Pathol, 2021. **255**(2): p. 141-154.
129. Mao, Y., et al., *Phenotype prediction from single-cell RNA-seq data using attention-based neural networks*. Bioinformatics, 2024. **40**(2).
130. julian.knight@well.ox.ac.uk, C.O.-M.-o.B.A.C.E.a. and C.O.-M.-o.B.A. Consortium, *A blood atlas of COVID-19 defines hallmarks of disease severity and specificity*. Cell, 2022. **185**(5): p. 916-938 e58.
131. Stephenson, E., et al., *Single-cell multi-omics analysis of the immune response in COVID-19*. Nat Med, 2021. **27**(5): p. 904-916.
132. Ren, X., et al., *COVID-19 immune features revealed by a large-scale single-cell transcriptome atlas*. Cell, 2021. **184**(7): p. 1895-1913 e19.
133. Theus, A., et al., *CancerFoundation: A single-cell RNA sequencing foundation model to decipher drug resistance in cancer*. bioRxiv, 2024: p. 2024.11. 01.621087.
134. Neftel, C., et al., *An Integrative Model of Cellular States, Plasticity, and Genetics for Glioblastoma*. Cell, 2019. **178**(4): p. 835-849 e21.
135. Wissel, D., et al., *Survboard: standardised benchmarking for multi-omics cancer survival models*. bioRxiv, 2022: p. 2022.11. 18.517043.
136. Cancer Genome Atlas Research, N., et al., *The Cancer Genome Atlas Pan-Cancer analysis project*. Nat Genet, 2013. **45**(10): p. 1113-20.
137. Querfurth, B.v., et al., *mBERT: Patient-Level Single-cell Transcriptomics Data Representation*. bioRxiv, 2024: p. 2024.11. 04.621897.
138. Chaffin, M., et al., *Single-nucleus profiling of human dilated and hypertrophic cardiomyopathy*. Nature, 2022. **608**(7921): p. 174-180.
139. Koenig, A.L., et al., *Single-cell transcriptomics reveals cell-type-specific diversification in human heart failure*. Nat Cardiovasc Res, 2022. **1**(3): p. 263-280.
140. Reichart, D., et al., *Pathogenic variants damage cell composition and single cell transcription in cardiomyopathies*. Science, 2022.

377(6606): p. eab01984.

141. Kuppe, C., et al., *Spatial multi-omic map of human myocardial infarction*. Nature, 2022. **608**(7924): p. 766-777.
142. Lake, B.B., et al., *An atlas of healthy and injured cell states and niches in the human kidney*. Nature, 2023. **619**(7970): p. 585-594.
143. Kuppe, C., et al., *Decoding myofibroblast origins in human kidney fibrosis*. Nature, 2021. **589**(7841): p. 281-286.
144. Muto, Y., et al., *Defining cellular complexity in human autosomal dominant polycystic kidney disease by multimodal single cell analysis*. Nat Commun, 2022. **13**(1): p. 6497.
145. Wilson, P.C., et al., *Multimodal single cell sequencing implicates chromatin accessibility and genetic background in diabetic kidney disease progression*. Nat Commun, 2022. **13**(1): p. 5253.
146. Muto, Y., et al., *Single cell transcriptional and chromatin accessibility profiling redefine cellular heterogeneity in the adult human kidney*. Nat Commun, 2021. **12**(1): p. 2190.
147. Perez, R.K., et al., *Single-cell RNA-seq reveals cell type-specific molecular and genetic associations to lupus*. Science, 2022. **376**(6589): p. eabf1970.
148. Yoshida, M., et al., *Local and systemic responses to SARS-CoV-2 infection in children and adults*. Nature, 2022. **602**(7896): p. 321-327.
149. Sikkema, L., et al., *An integrated cell atlas of the lung in health and disease*. Nat Med, 2023. **29**(6): p. 1563-1577.
150. Schaar, A., et al., *Nicheformer: a foundation model for single-cell and spatial omics*. 2024. Preprint at bioRxiv, 2024. **4**: p. 589472.
151. Yao, Z., et al., *A high-resolution transcriptomic and spatial atlas of cell types in the whole mouse brain*. Nature, 2023. **624**(7991): p. 317-332.
152. He, S., et al., *High-plex multiomic analysis in FFPE at subcellular level by spatial molecular imaging*. bioRxiv 467020. 2021.
153. Wen, H., et al., *Single cells are spatial tokens: Transformers for spatial transcriptomic data imputation*. arXiv preprint arXiv:2302.03038, 2023.