

A Transformer for scATAC-scRNA Translation

by

Roger Jin

S.B. Computer Science and Engineering and Mathematics
Massachusetts Institute of Technology (2020)

Submitted to the Department of Electrical Engineering and Computer
Science

in partial fulfillment of the requirements for the degree of

Master of Engineering in Electrical Engineering and Computer Science

at the

MASSACHUSETTS INSTITUTE OF TECHNOLOGY

September 2022

© Massachusetts Institute of Technology 2022. All rights reserved.

Author
Department of Electrical Engineering and Computer Science
September 08, 2022

Certified by.....
Manolis Kellis
Professor of Computer Science and Electrical Engineering
Thesis Supervisor

Accepted by
Katrina LaCurts
Chair, Master of Engineering Thesis Committee

A Transformer for scATAC-scRNA Translation

by

Roger Jin

Submitted to the Department of Electrical Engineering and Computer Science
on September 08, 2022, in partial fulfillment of the
requirements for the degree of
Master of Engineering in Electrical Engineering and Computer Science

Abstract

scATAC-seq gives a comprehensive picture of the chromatin accessibility profile of a cell, covering not only protein-coding regions but also non-coding regulatory regions which are in theory missed by scRNA-seq. However, scATAC-seq data is high-dimensional and noisy, aspects which when compounded with data scarcity present challenges for modeling on even seemingly-simple downstream tasks such as cell-type prediction. As such, researchers may benefit from access to a large library of models to evaluate. While we do not demonstrate state of the art results in any capacity, we provide an implementation of a simple representation of sparse tabular data that allows it to be inputted into the popular transformer family of architectures, and use this representation to train a transformer that predicts scRNA-seq given scATAC-seq. Our code is made available [here](#).

Thesis Supervisor: Manolis Kellis

Title: Professor of Computer Science and Electrical Engineering

Acknowledgments

I would like to thank Professor Manolis Kellis for the opportunity to work in his lab, and for believing in and supporting this project.

I would like to thank my mentor Ben Lengerich for his continued support and encouragement, brilliant ideas, and fun discussions both about ML and about life. His excitement for ML is infectious, and some of the coolest ideas in ML I know about were first passed down to me by Ben.

I would like to thank Alvin Shi for being my first mentor at the Kellis lab, and for also giving me my first ML project my sophomore year of college.

I would also like to thank my co-mentor Jackie Yang, and fellow MEng student Doris Fu for their helpful discussions and continued intellectual guidance.

I would like to thank Na Sun and Xushen Xiong also for helpful discussions and for first introducing me to the scRNA-scATAC integration problem.

I would like to thank the NeurIPS Single Cell Competition team for not only organizing the competition and their amazing dataset, but also for gathering a community of researchers together to tackle these very challenging multiomics problems. I would like to thank Malte Luecken for showing me the SCiB [24] evaluation pipeline, and for answering my many questions about the dataset which he and his team put together.

I would like to thank Jehangir Amjad, Leslie Kaelbling, Ashia Wilson, Jacob White, and Yoon Kim for making my 6.036 TAship one of the most fun and rewarding experiences of my life.

I would like to thank Alex Churchill, Pranay Dighe, Vineet Garg, Erik Marchi, and Ahmed Tewfik at Apple who first introduced me to transformers and language modeling during two internships, edifying experiences which have shaped my interests for perhaps the rest of my career.

I would like to thank Ekin Akyürek for helpful discussion about this project and its framing and presentation.

I would like to thank Vera Sayzew of the Ugrad Admin office for her constant encouragement and understanding.

Lastly, I would like to thank my friends Shahul Alam, Varkey Alumootil, Michael Diao, Daniel Hu, Brice Huang, Steven Liu, Jason Lu, Sophia Xia, Jason Yang, Brendan Yap, and girlfriend YuanYuan Du who ran this race with me.

Contents

1	Introduction	13
1.1	The current paradigm	14
1.2	Why learn a translation instead of integrating directly?	15
1.2.1	Translation as a certificate	15
1.2.2	More fundamental questions	15
1.3	Existing methods	16
1.3.1	Seurat (v3): Canonical Correlation Analysis	16
1.3.2	LIGER: Nonnegative Matrix Factorization	16
1.3.3	Zero-knowledge integration: MultiMAP	17
1.4	Motivation for using transformers	17
1.5	The challenge of sparsity	18
2	Data	19
2.1	Data Format	19
2.2	Sparsity	19
2.3	Dataloading	20
2.4	Data Visualization	21
3	Methods	23
3.1	The squish embedding	23
3.2	The model	23
3.2.1	Connection to vanilla MLPs	24

4	Experiments	27
4.1	Experimental setup	27
4.2	Results and Discussion	29
5	Future work	31
5.1	Unsupervised training with CycleGAN	31
5.1.1	Cycle-consistency Loss	32
5.1.2	Limitation: GPU memory	33
5.2	Autoregressive decoding	33
5.2.1	Measuring enhancer-gene interaction with cross-attention	33
5.2.2	Limitation: sequence length	34
5.3	Biological awareness	34
5.3.1	3D chromatin structure	34
5.3.2	Cell-type awareness	34
6	Conclusion	35
A	Model details	37
A.1	DistilBERT	37
A.2	LongT5	37
A.3	Classification head	38
A.4	MLP	38
B	Training details	39
B.1	Hardware	39
B.2	Optimizer	39
C	Implementation Details	41
C.1	Scipy sparse data formats	41
C.2	Retrieving nonzero entries per row	41
C.3	Computing the squish index	41

List of Figures

1-1	Schematic of how our translator model can be used for scATAC-scRNA integration. UMAP on right taken from [1].	14
2-1	Histogram of number of nonzero counts per cell in train partition of scATAC-seq.	20
2-2	UMAP visualizations of scATAC data.	22
3-1	Figure adapted from [3]. The blank white box denotes a 0, which is equivalent to deleting the column W_2 from the linear combination. . .	24
3-2	Squish transformer architecture. Figure drawn based on design of [3].	25
5-1	Proposed architecture based on [45]. G_{AR} is a generator that performs ATAC-RNA translation, and G_{RA} is a generator that performs RNA-ATAC translation. D_R is a discriminator for RNA. We show the loss calculation for one training step of G_{AR} . D_A , the discriminator for ATAC, is omitted from the diagram. Figure adapted from [35]	32

List of Tables

2.1	Descriptive statistics of nonzero counts per cell in train partition of scATAC-seq.	20
4.1	MSE results from training	29

Chapter 1

Introduction

Single-cell sequencing allows interrogation of the intercellular heterogeneity within samples. In the past decade, single-cell RNA-sequencing technology has developed explosively, at a rate rivaling Moore's Law [37]. And in parallel, profiling methods for other omics modalities have improved their resolution to the single-cell level. It is now possible to collect single-cell data for not only genomics and transcriptomics, but also DNA accessibility, epigenetic modifications, and chromatin structure. Recently, technologies have been developed that can simultaneously profile multiple data modalities, such as SNARE-seq [8] for scRNA-scATAC, and Slide-seq [32] for spatial transcriptomics. However, much of the readily available data remains unimodal, and a challenge that has arisen over recent years is that of integrating multiple data modalities collected from the same sample - that is, matching corresponding cell states across modalities to garner biological insights across multiple modalities. One long sought-after goal of such integration is the ability to infer causal gene regulatory networks from a combination of epigenomic and transcriptomic data, but to our knowledge this has not yet been accomplished. Our work focuses on scRNA-scATAC integration, though our methodology can potentially be adapted to other data modalities.

1.1 The current paradigm

scRNA-seq and scATAC-seq are complementary methods for profiling the internal state of a cell. Whereas scRNA-seq measures gene expression via sequencing of mRNA transcripts, scATAC-seq measures chromatin accessibility via sequencing of DNA segments excised by Tn5 transposase [19]. To what extent one modality can be inferred from the other remains an open question, and is the subject of our investigation.

The current paradigm for integrating scATAC and scRNA is to effectively translate ATAC to RNA by converting ATAC counts into a gene activity matrix by for each gene, binning ATAC counts near its transcription start site (TSS). This method uses only ATAC signal local to a gene, and does not account for regulatory factors or processes which may alter gene expression by means other than controlling chromatin density.

Upon translating scATAC into scRNA, methods traditionally used for integrating multiple scRNA datasets (e.g. for batch correction) can then be used to integrate scATAC with scRNA datasets (Figure 1-1).

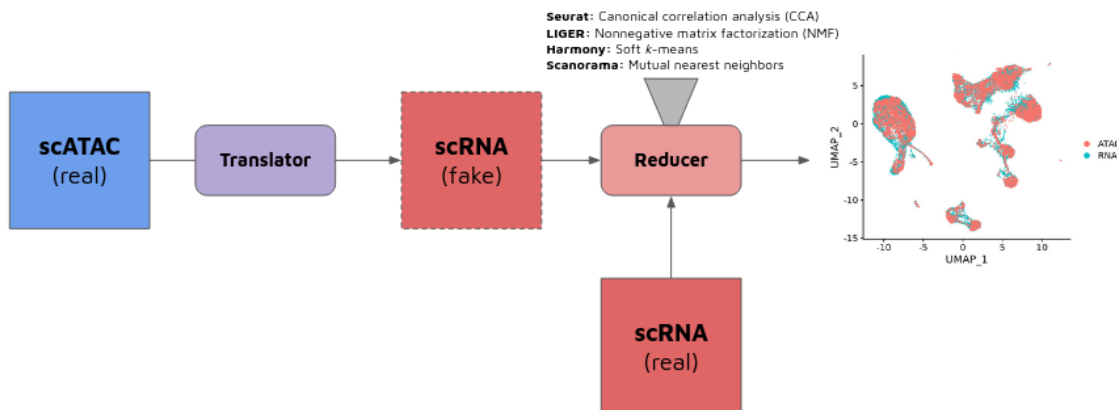


Figure 1-1: Schematic of how our translator model can be used for scATAC-scRNA integration. UMAP on right taken from [1].

Our goals are to

- (a) Find some translation between scATAC and scRNA which has comparable or better performance than existing methods.

- (b) In training such a mapping between scATAC and scRNA, learn something about the relationship between ATAC and RNA.

1.2 Why learn a translation instead of integrating directly?

1.2.1 Translation as a certificate

Although there exist methods which can directly align scATAC-seq and scRNA-seq measurements in a shared latent space [17] without an explicit translation between them, these methods in principle suffer from a mutual exclusion between user trust and user benefit. For example, a researcher may want to transfer cell type annotations from cells whose scRNA-seq profiles are known but scATAC-seq profiles are unknown to cells whose scATAC-seq profiles are known but scRNA-seq profiles are unknown. A reasonable approach may then be to co-embed the cells in a shared latent space and label each scATAC-known cell with the cell type most abundant among its 5-nearest scRNA-known neighbors. In order to validate such a method on the researcher’s dataset of interest however, the researcher must already have access to some method of generating “ground truth” cell types for the scATAC cells, in which case the original integration method becomes unnecessary for cell-type prediction.

Although the predicted gene expression levels themselves may not be the desired downstream target, a full translation from scATAC to scRNA is one form of alternative certificate to ground truth targets that an integration method is trustworthy.

1.2.2 More fundamental questions

If all of the information in scRNA is encapsulated within scATAC or vice versa, then there is no need to perform both measurements. And if it is difficult to learn a translation between scATAC and scRNA, this raises the question of what information is contained in one that is impossible to infer from the other, and what additional data could be leveraged to augment the original input to make this impossible inference

possible. For example, scATAC-seq is in some sense more general than scRNA-seq because it captures the accessibility of not just genes, but also regulatory regions. However, it is unknown whether accessibility is sufficient to predict gene expression - that is, whether there are multiple plausible scRNA state for a single scATAC state. Even with a fixed set of accessible regions, factors such as 3D proximity, epigenomic markers, and environmental factors could generate some variability in gene expression. Augmenting scATAC input with scHiC [27], single-cell methylation [21], or single-cell metabolomics [34] might shed some insight into these questions.

1.3 Existing methods

There are a host of methods of translating between scRNA-seq and scATAC-seq, and we review a few of the popular methods, as well as a few recent methods which have demonstrated state of the art performance.

1.3.1 Seurat (v3): Canonical Correlation Analysis

Seurat v3 [36] introduces the Signac package for scRNA-scATAC integration. After deriving a gene by cell matrix from scATAC, Signac uses canonical correlation analysis (CCA) to find a linear projection of the modalities into a low dimensional subspace that maximally correlates gene activities across cells. Seurat then constructs a nearest neighbors graph and applies a batch transformation to bring mutual nearest neighbors in the graph closer together. One advantage of Seurat is that CCA has a unique closed-form solution (given by the left singular vectors of the SVD), and so one can expect consistent results across multiple applications of the method. However, only two datasets can be integrated at a time, and the dimensionality reduction is linear.

1.3.2 LIGER: Nonnegative Matrix Factorization

LIGER [41] also begins by deriving a gene by cell matrix from scATAC, using Seurat. It then decomposes each inputted gene by cell matrix by nonnegative matrix factor-

ization (NMF) as a sum of outer products, where one outer product contains a matrix shared across all modalities. Cell-representations derived from this outer product are then used to construct a neighborhood graph where clustering is performed. One advantage of LIGER is that it can accept multiple data modalities, and nonnegativity gives some potential for interpretability. However, NMF in general does not have a unique solution, and so its application may lead to inconsistent results.

1.3.3 Zero-knowledge integration: MultiMAP

MultiMAP [17] is an extension of UMAP [25] for multimodal data integration, and unlike translation-based methods, does not require inputted datasets to have the same feature dimensions. Using similar methods from algebraic topology as UMAP, MultiMAP integrates distance metrics across the multiple datasets to embed them on a shared latent manifold. MultiMAP is extremely flexible and computationally efficient, and produces a nonlinear transformation of the data. However, one limitation is its lack of interpretability, though one may argue that none of the existing methods are very interpretable.

1.4 Motivation for using transformers

Transformers [38] currently sit at the top of many benchmarks spanning language [16, 10], vision [12, 43], and recently even reinforcement learning [7, 18]. Models are moving from specialist to generalist [30], and there has been some evidence that models are already capable of generalization across data modalities [29]. While we make no claims that transformers are the best architecture suited for scATAC-seq, undoubtedly much of the research and progress in machine learning today involves some kind of transformer, and having the option to use a transformer may prove powerful in the future. Our main contribution is a representation of scATAC data which is easily inputted into a transformer.

1.5 The challenge of sparsity

One computational challenge is that ATAC-data is high-dimensional - for each cell, we have 116,490 observations (peaks). If we were to naively design a feedforward net that consumes this data as input, However, this data is usually sparse, and on average only around 3,400 peaks have nonzero count. As such, we might expect the data to live in some lower-dimensional subspace of the full feature space, and we can exploit this via dimensionality reduction or tensor decomposition. However, in applying dimensionality reduction, we often lose interpretability along the feature axis - we lose sense of individual peaks and correlations within sets of peaks in exchange for a lower-dimensional embedding of cells. We propose a simple representation of scATAC data that is

1. Compact - the size of the representation is variable and scales linearly with the number of peaks with nonzero counts.
2. Extensible - there is a natural way to add support for new peaks without performance degradation on the set of peaks currently supported (this is guaranteed because we do not need to alter any of the original parameters of the model when adding support for a new peak).
3. Informative - there is a natural way to produce representations for both peaks and cells.

Admittedly these properties are not unique to our setup, and in fact one could argue that the weights from a simple MLP or OLS regression could be interpreted in similar fashion.

Chapter 2

Data

We use the ATAC2GEX data from the NeurIPS 2021 Multimodal Single-Cell Data Integration Competition [23]. The dataset consists of 42,491 cells, and for each cell the authors provide counts for 116,490 ATAC peaks and expression levels for 13,491 genes. We partition this data in an 80%/10%/10% split for train, validation, and testing respectively.

2.1 Data Format

The data is provided in `h5ad` files, an `hdf5` format specific to the `anndata` package [39], which provides an abstraction for annotated dataframes with support for sparse data. Both ATAC and RNA are stored in compressed sparse column (CSC) format.

2.2 Sparsity

We provide in [Figure 2-1](#) and [Table 2.1](#) descriptive statistics of the number of nonzero entries per cell in the ATAC training data. We notice that there is quite a bit of variance in the number of nonzero entries, which corresponds to the sequence length in our transformer. Although we did not have time to implement this, an optimization to our dataloader that would've sped up training significantly would've been to sort the data by sequence length, sample batches of sequences of similar length, and use

larger batch sizes for shorter sequences.

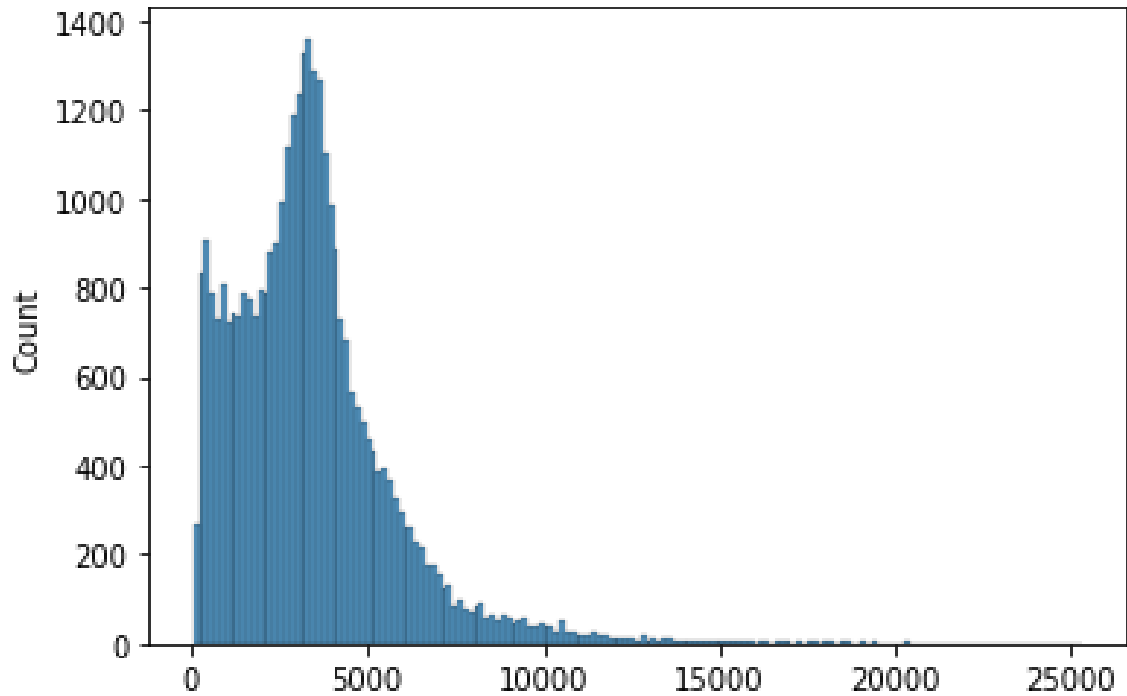


Figure 2-1: Histogram of number of nonzero counts per cell in train partition of scATAC-seq.

Count	33,993
Mean	3379
Std	2259
Min	62
25%	1827
50%	3155
75%	4306
max	25301

Table 2.1: Descriptive statistics of nonzero counts per cell in train partition of scATAC-seq.

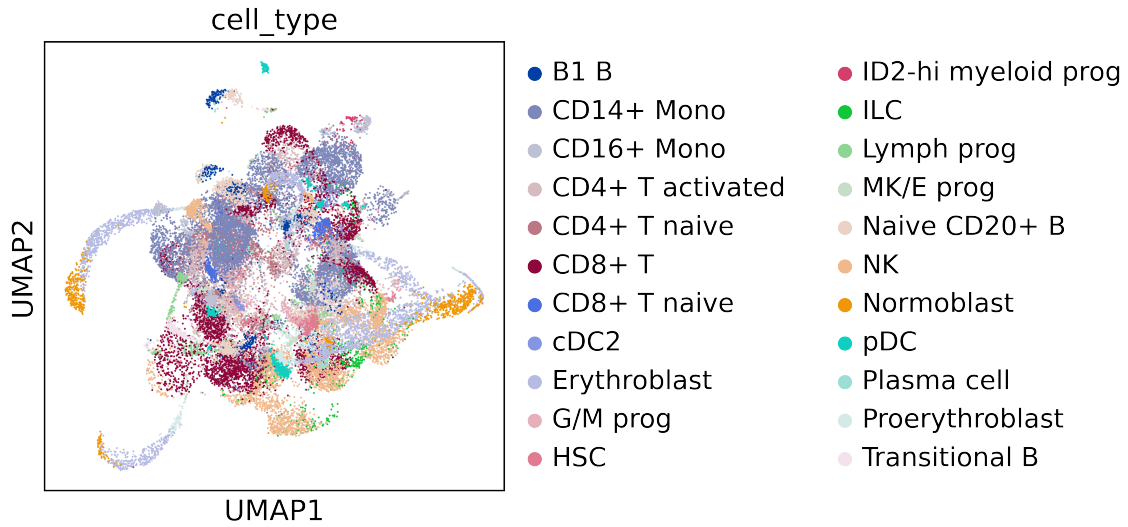
2.3 Dataloading

`anndata`'s dataloader converts batches to dense matrices by default, but we found that there was a relatively straightforward vectorized way to convert a batch of ATAC-

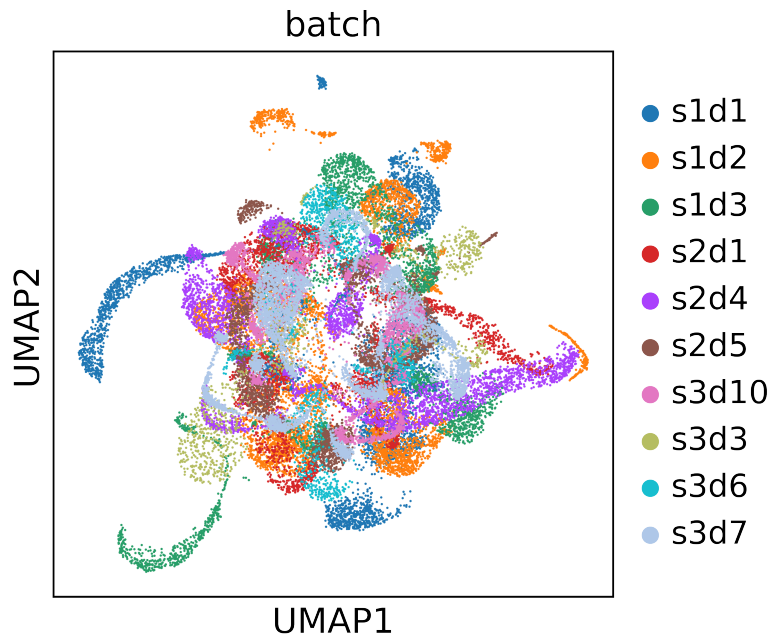
counts in CSC format to our squish index format, and we were able to write a dataloader that avoids ever densifying the ATAC batch, saving some memory.

2.4 Data Visualization

We provide in [Figure 2-2](#) uniform manifold approximation and projection (UMAP) visualizations of the ATAC data, labeled both by cell type and experimental batch. We notice that there are clear batch effects, as the cells cluster by batch. We also notice common artifacts between batches, such as erythroblasts forming outward-emanating tails that terminate in a cluster of normoblasts.



(a) UMAP of the ATAC-data, labeled by cell type.



(b) UMAP of the ATAC-data, labeled by experimental batch.

Figure 2-2: UMAP visualizations of scATAC data.

Chapter 3

Methods

3.1 The squish embedding

What is the most naive way to compress a row of a sparse matrix? If we have 5,000 nonzero counts dispersed throughout a sea of 116,490 zeros, we can wishfully delete all the zeros and represent the row as a sequence of just 5,000 counts. However, when disassociated from their indices, the counts lose their meaning. We repair this by learning an embedding vector for each of the 116,490 peaks and scaling it by the corresponding count. Each cell is then represented by a sequence of embedding vectors, scaled by their corresponding count. We choose not to normalize or softmax these scale factors because LayerNorm [5] will prevent them from shifting the mean or variance of the post-multi-head-attention activations.

3.2 The model

Since we have formulated the input as a sequence (or more accurately, a set), it is then straightforward to train a transformer to map an input ATAC squish embedding representation to an output gene expression vector. One key difference between the transformer parameterization and the naive feedforward net parameterization which enables this setup is that aside from the positional embedding matrix, which scales linearly with the maximum sequence length, the transformer parameterization is

independent of the sequence length. In our case our embedding matrix is sort of like a positional embedding matrix, as each index corresponds to a peak. Thus given a fixed set of possible ATAC peak measurements, we are able to input sequences of arbitrary length up to the number of embeddings in the embedding matrix (currently 116,490). [Figure 3-2](#) gives a simple schematic of the model.

3.2.1 Connection to vanilla MLPs

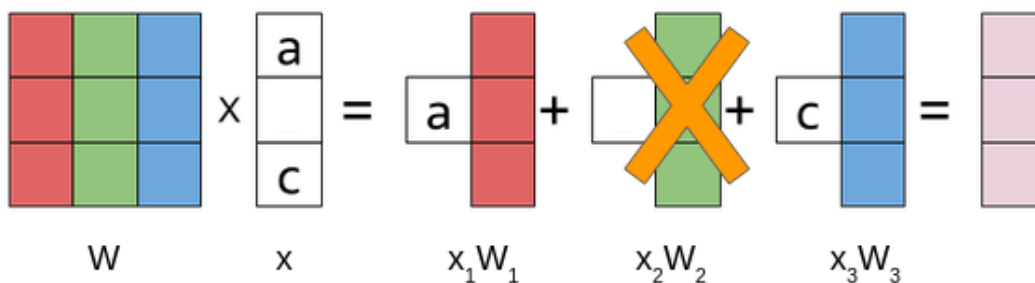


Figure 3-1: Figure adapted from [3]. The blank white box denotes a 0, which is equivalent to deleting the column W_2 from the linear combination.

The method we propose is essentially matrix multiplication with a slight modification, and explaining this modification requires a bit of boring notational setup for which we apologize to readers in advance. Let W be the weights of the first linear layer of a neural network \mathcal{M} and let x be an input to \mathcal{M} . Then the first layer implements the function $f : \mathbb{R}^{\dim(x)} \rightarrow \mathbb{R}^{\dim(Wx)}$

$$\begin{aligned}
 f(x) &= Wx \\
 &= x_1 W_1 + \dots + x_n W_n,
 \end{aligned}$$

where we denote the k th column of W by W_k . We can then view matrix multiplication as a two step process:

- 1) Scale the columns of W by the corresponding entries in x .
- 2) Sum the scaled columns together.

Our method simply replaces the sum in step 2) with a transformer followed by a pooling function. The transformer is a function from $\mathbb{R}^{\dim(x) \times d_1} \rightarrow \mathbb{R}^{\dim(x) \times d_2}$ for hidden dimensions d_1, d_2 , and the pooling function is a function from $\mathbb{R}^{\dim(x) \times d_2} \rightarrow \mathbb{R}^{d_2}$, so taking $d_2 = \dim(Wx)$, their composition has the same domain and range as $f(x) = Wx$.

However, one important implementation detail is that instead of feeding in 0 vectors to the transformer, we delete them from the sequence¹. To support batch processing, we use zero vectors as padding and use an attention mask to avoid attending to padding.

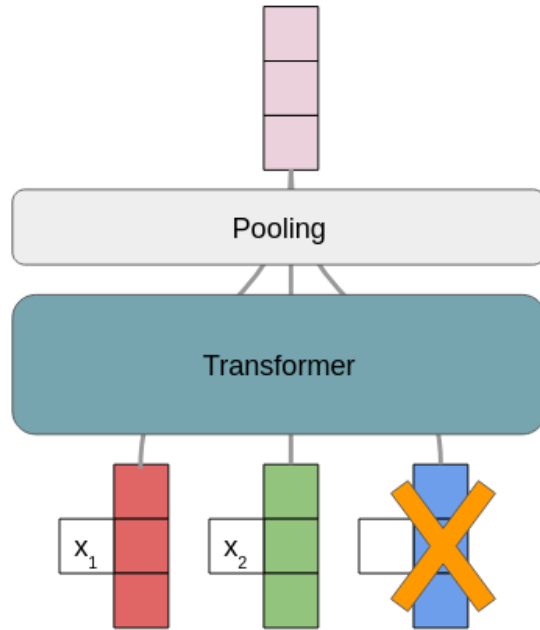


Figure 3-2: Squish transformer architecture. Figure drawn based on design of [3].

¹More accurately, we construct the squish embedding without the zero vectors in the first place.

Chapter 4

Experiments

Having specified a representation of scATAC as a sequence of scaled embeddings, we can in theory feed these into any encoder model. In practice we are limited by GPU-memory, and as the original transformer [38] on which many subsequent models are based uses full global self-attention, for which in theory time and memory scale quadratically in sequence length, we are somewhat limited in our choices to either small lightweight models or models which replace full self-attention with another mapping from $\mathbb{R}^{\text{seq len} \times d_1} \rightarrow \mathbb{R}^{\text{seq len} \times d_2}$, for hidden dimensions d_1, d_2 . We choose an example from either family. As an example of a small lightweight encoder, we choose DistilBert [33], which is a compression of BERT [11] trained via teacher-student knowledge-distillation, and as an example of a larger model with a memory-efficient drop-in replacement for attention, we take the encoder portion of LongT5 [15], which uses the hybrid detailed local attention/compressed distant attention method from ETC [4].

4.1 Experimental setup

We run training for a fixed number of epochs, and report lowest validation loss achieved. Regrettably, I had a bug where I was using the final model state rather than the best-performing model state in testing, so these are validation numbers which should not be regarded as a fair measure of the generalization capability of the

model.

Running on A100 GPUs, I found I could fit a sequence length of around 1600 into DistilBERT with batch size 32, and ran training under the following conditions:

1. `distilbert_counts`: I use the last layer embedding of the first token as input to the regression head, and ATAC counts as input. No [CLS] token introduced.
2. `distilbert_binary`: I use the last layer embedding of the first token as input to the regression head, and binary ATAC counts as input (1 if count > 0, else 0).
3. `distilbert_cls_counts`: I prepend all inputs with a [CLS] token and use the last layer embedding of the [CLS] token as input to the regression head, and ATAC counts as input.
4. `distilbert_cls_counts_binary`: Same as `distilbert_cls_counts`, but with binarized counts as input.

As the DistilBERT results were unpromising and only 20% of cells had fewer than 1600 nonzero peaks, I ran LongT5 with sequence length 9633¹. On V100 I could only fit a batch size of 1, while on A100 I could fit a batch size of 4. Training was much slower with the longer sequence length and smaller batch size, and I was not able to train a model until validation loss reach a local minimum, so I report the final loss.

1. `longt5_cls_counts`: I use the last layer embeddings of the [CLS] token, and ATAC counts as input.

As a baseline I fed the original ATAC data into a 2-layer MLP.

1. `mlp_counts`: I feed the ATAC counts into a two-layer MLP. No squish embedding.
2. `mlp_binary`: Same as `mlp_counts` but with binary input.

¹A sequence length of 10,000 gives OOM on V100 so this is close to the longest possible

4.2 Results and Discussion

We present results from training in [Table 4.1](#). As we do not outperform a simple two-layer MLP, it is evident that our method does not improve MSE. As consolation we discuss some of the possible future work that may make this method worthwhile even if it is outperformed by MLP.

Experimental Condition	Min Val MSE
distilbert_counts	0.094
distilbert_binary	0.095
distilbert_cls_counts	0.096
distilbert_cls_counts_binary	0.096
longt5_cls_counts	0.100 ²
mlp_counts	0.093
mlp_binary	0.090

Table 4.1: MSE results from training

Chapter 5

Future work

5.1 Unsupervised training with CycleGAN

Simultaneously profiling scATAC and scRNA per cell is expensive and noisy, and there are currently only a few publicly available multiomic datasets available. In contrast, there are large single-cell atlases available for both scRNA-seq and scATAC-seq [31] individually, where for cells of the same type, both measurements may be available, but only one is known per cell. This is analogous to the problem setting of unsupervised image to image translation in computer vision, and so it may be possible to adopt an adversarial cycle-consistent training scheme such as CycleGAN's to leverage unpaired data in either modality.

As depicted in [Figure 5-1](#), the architecture consists of the following components:

1. A generator which takes ATAC as input and outputs a translation to RNA, which we denote G_{AR} .
2. A generator which takes RNA as input and outputs a translation to ATAC, which we denote G_{RA} .
3. Discriminators D_R for RNA and D_A for ATAC.

5.1.1 Cycle-consistency Loss

We can define the cycle-consistency loss as follows:

$$\mathcal{L}_{\text{cyc}}(G_{AR}, G_{RA}) = \mathbb{E}_{a \sim \text{ATAC}} [\|G_{RA}(G_{AR}(a)) - a\|_1] + \mathbb{E}_{r \sim \text{RNA}} [\|G_{AR}(G_{RA}(r)) - r\|_1].$$

We can thus write the loss function for the full architecture as

$$\mathcal{L}(G_{AR}, G_{RA}, D_A, D_R) = \mathcal{L}_{\text{cyc}}(G_{AR}, G_{RA}) + \mathcal{L}_{\text{WGAN}}(G_{AR}, D_R) + \mathcal{L}_{\text{WGAN}}(G_{RA}, D_A),$$

where $\mathcal{L}_{\text{WGAN}}$ denotes the Wasserstein loss.

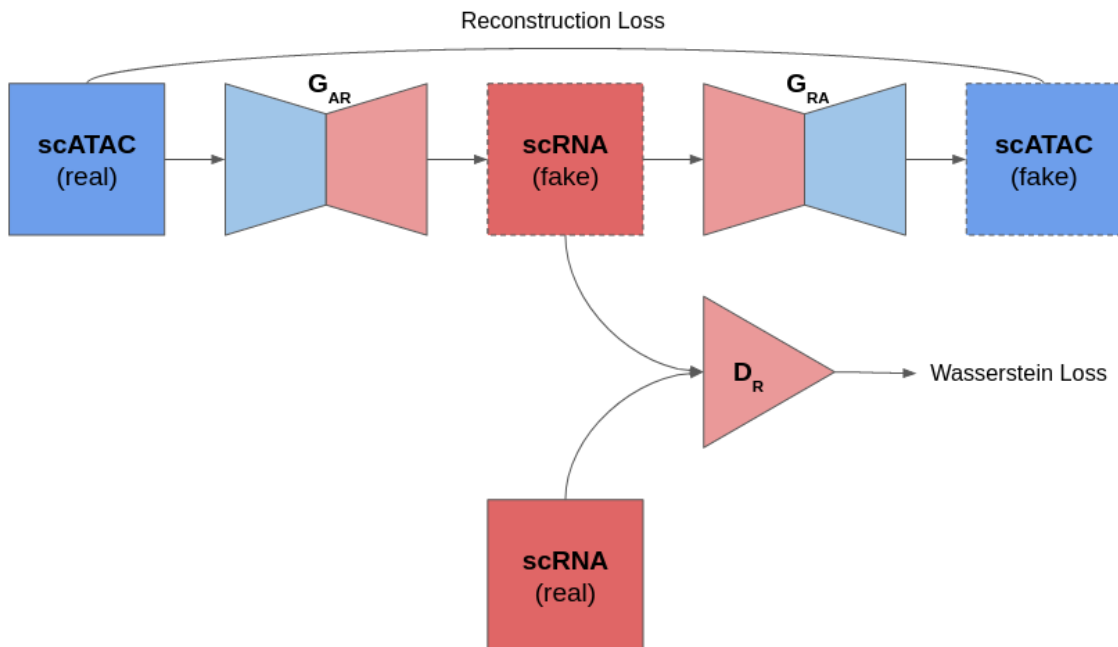


Figure 5-1: Proposed architecture based on [45]. G_{AR} is a generator that performs ATAC-RNA translation, and G_{RA} is a generator that performs RNA-ATAC translation. D_R is a discriminator for RNA. We show the loss calculation for one training step of G_{AR} . D_A , the discriminator for ATAC, is omitted from the diagram. Figure adapted from [35]

5.1.2 Limitation: GPU memory

We elected to stop pursuing this direction because at least with a naive implementation, all four models are hosted on the same GPU (although we only need to track gradients for two at a time). This severely limits the size of the models that can be tested. However, recent advances in parallelization [6] have demonstrated the potential to reduce model memory footprint significantly, and perhaps this direction should be revisited.

5.2 Autoregressive decoding

There is currently no popular method which decodes either a piece of text or an entire image in one step from encoder embeddings. Conditioned on the input context, the distribution of outputs in each position could still be multimodal, and taking a sample from each distribution independently could result in a generation from a low-density portion of the true joint output distribution. Autoregressive decoding [38, 2] iteratively chooses tokens and adds them to the context, so that subsequently chosen tokens are more likely to be consistent with previously chosen ones.

5.2.1 Measuring enhancer-gene interaction with cross-attention

Performance gains in terms of realism of output aside, adding a step-wise decoder will give us access to cross-attentions between queries from RNA and keys from ATAC, which we may be able to interpret as some measure of interactivity between genes and peaks. This type of analysis is popular for language models and their fine-tuned derivatives [20].

A step-wise decoder will also necessarily force us to learn a joint embedding of ATAC-peaks and genes, which may allow us to draw analogies [26] between peaks and genes.

5.2.2 Limitation: sequence length

In our case we would be interested in decoding an output of size 13,491 from an input of size 116,490. We saw earlier that squish embedding lowered the sequence length from 116,490 to around 3400 on average. We can apply a similar trick in the decoding stage, where if we generate an estimate of an expression level near 0, we do not add it to the context. However, as we are running into OOM errors even now in our encoder-only setup, we anticipate actual implementation of this to be very tricky.

5.3 Biological awareness

5.3.1 3D chromatin structure

Thus far we have mostly considered generic transformer architectures that are suitable across many different tasks, but there may be small adaptations we can make to specialize transformers to biology. For example, while sparse transformers [9, 44, 4] often employ some form of sliding-window attention with the thought that local information probably matters more than global information, this form of attention is not able to capture the 3D structure of chromatin. It may be that an additional set of positional embeddings is needed to convey at least generic information about 3D proximity, perhaps based on a HiC contact map [14].

5.3.2 Cell-type awareness

If one has access to ground-truth cell type labels for both scATAC and scRNA, a straightforward way to do make integration more reasonable is to partition cells by type and integrate within each type. However, currently there is no widely accepted method to attain ground-truth cell-type labels from scATAC data alone, suggesting that we could allow these types to be determined by the model. The clustered structure of single-cell data may then be well-suited to modeling by a mixture of experts model [13].

Chapter 6

Conclusion

We provide an implementation of an interface that allows researchers to input sparse data into a transformer, giving researchers access to tens of thousands of open-source transformer models [42]. We evaluate our method on the task of scATAC-scRNA translation, unfortunately demonstrating no better performance than a two-layer MLP. We posit that despite this, our work still has some downstream benefits in terms of interpretability by providing the potential to compute cross-attention between scRNA genes and scATAC peaks by training an autoregressive decoder for scRNA. We propose extensive future work, including a mixture of experts model which may be able to learn specialized networks for each cell type.

Appendix A

Model details

Experiment configuration files are available [here](#).

A.1 DistilBERT

We use the pre-trained weights of HuggingFace’s [distilbert-base-uncased](#) for the transformer body, and overwrite the default embedding matrix with a randomly initialized one of size $116,490 \times 768$, the number of ATAC peaks by the hidden dimension of the model.

We also set the positional encodings to a custom Zero PyTorch module which returns 0 given any input.

We use a sequence length of 1600 and a batch size of 32. Due to the short sequence length, we sort the ATAC peaks in decreasing order by count variance and select the top 1600 most variable nonzero peaks for each cell we input.

A.2 LongT5

We use the pre-trained weights of the encoder portion of Huggingface’s [long-t5-local-base](#). As far as I can tell, LongT5 does not use positional embeddings (which I suppose is unnecessary because this version implements local attention to give some positional awareness), so there is no need to zero them out.

We use a sequence length of 9663 and the maximum batch size possible before OOM errors. For a 32GB V100 this is 1 and for an 80GB A100 this is 4.

A.3 Classification head

For both DistilBERT and LongT5, we add a classification head which implements the function

$$f(T(x)) = \text{ReLU}(W_2 \text{ReLU}(W_1 T(x) + b_1) + b_2),$$

where we include the outer ReLU because RNA expression is nonnegative. W_1 is 768×768 and W_2 is $13,491 \times 768$.

A.4 MLP

Similarly, our MLP baseline implements the function

$$f(T(x)) = \text{ReLU}(W_2 \text{ReLU}(W_1 T(x) + b_1) + b_2),$$

but where W_1 is $768 \times 116,490$ and W_2 is $13,491 \times 768$.

Appendix B

Training details

B.1 Hardware

We primarily train our models on A100 and V100 GPUs on the [MIT OpenMind Cluster](#). I would like to thank Professor Kellis for generously sponsoring a lab subscription to the cluster.

I also have a setup on [PaperSpace](#) which I use mostly for development purposes.

B.2 Optimizer

I use the PyTorch [\[28\]](#) implementation of the AdamW optimizer [\[22\]](#) with default $[\beta_1, \beta_2]$ and learning rate between 10^{-5} and 10^{-3} .

Appendix C

Implementation Details

C.1 Scipy sparse data formats

Scipy [40] supports sparse matrix storage in **COO**, **CSR**, and **CSC** formats, which are relatively efficient to convert between. The data from the NeurIPS 2021 competition [23] is stored in CSC format.

C.2 Retrieving nonzero entries per row

We measure sparsity by retrieving the number of nonzero entries per row. There is an efficient trick to do this with the CSC format, by taking incremental differences of the `indptr` attribute as shown below.

```
import numpy as np
```

```
nonzeros = np.diff(my_csc_matrix.indptr)
```

C.3 Computing the squish index

Although the squish index is not necessarily sparse, it is convenient to first compute its sparse COO format and then densify it. At a high level, we follow approximately the following steps:

1. Convert the input matrix to CSR format and then COO format. The CSR format intermediate is important for correct sorting of the rows. The `row` and `data` attributes of the COO format will be the same as the final matrix. All that remains is to compute the `col` attribute.
2. Compute the number of nonzero entries per row, either with the method from [Appendix C.2](#) or by using the `torch.unique_consecutive` function on the `row` attribute.
3. Compute a “cyclic arrange” of the number of nonzero entries. That is - for a nonzero entry count array $[x_1, x_2, \dots, x_n]$, return

$$[0, 1, \dots, x_1 - 1, 0, 1, \dots, x_2 - 1, \dots, 0, 1, \dots, x_n - 1].$$

This cyclic arrange is the `col` attribute.

4. The `data`, `row`, `col` can then be used to construct the squish index in COO format, which can then directly be densified into the desired format.

The full implementation is available [here](#).

Bibliography

- [1] Joint RNA and ATAC analysis: 10x multiomic. September 2022. [Online; accessed 8. Sep. 2022]. URL: https://satijalab.org/signac/articles/pbmc_multiomic.html.
- [2] The Unreasonable Effectiveness of Recurrent Neural Networks. March 2022. [Online; accessed 8. Sep. 2022]. URL: <http://karpathy.github.io/2015/05/21/rnn-effectiveness>.
- [3] Visualizing matrix multiplication as a linear combination - Eli Bendersky's website. September 2022. [Online; accessed 8. Sep. 2022]. URL: <https://eli.thegreenplace.net/2015/visualizing-matrix-multiplication-as-a-linear-combination>.
- [4] Joshua Ainslie, Santiago Ontanon, Chris Alberti, Vaclav Cvicek, Zachary Fisher, Philip Pham, Anirudh Ravula, Sumit Sanghai, Qifan Wang, and Li Yang. ETC: Encoding Long and Structured Inputs in Transformers. *arXiv*, April 2020. [arXiv:2004.08483](https://arxiv.org/abs/2004.08483), [doi:10.48550/arXiv.2004.08483](https://doi.org/10.48550/arXiv.2004.08483).
- [5] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. Layer Normalization. *arXiv*, July 2016. [arXiv:1607.06450](https://arxiv.org/abs/1607.06450), [doi:10.48550/arXiv.1607.06450](https://doi.org/10.48550/arXiv.1607.06450).
- [6] Zhengda Bian, Hongxin Liu, Boxiang Wang, Haichen Huang, Yongbin Li, Chuanrui Wang, Fan Cui, and Yang You. Colossal-AI: A Unified Deep Learning System For Large-Scale Parallel Training. *arXiv*, October 2021. [arXiv:2110.14883](https://arxiv.org/abs/2110.14883), [doi:10.48550/arXiv.2110.14883](https://doi.org/10.48550/arXiv.2110.14883).
- [7] Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Michael Laskin, Pieter Abbeel, Aravind Srinivas, and Igor Mordatch. Decision Transformer: Reinforcement Learning via Sequence Modeling. *arXiv*, June 2021. [arXiv:2106.01345](https://arxiv.org/abs/2106.01345), [doi:10.48550/arXiv.2106.01345](https://doi.org/10.48550/arXiv.2106.01345).
- [8] Song Chen, Blue B. Lake, and Kun Zhang. High-throughput sequencing of the transcriptome and chromatin accessibility in the same cell. *Nat. Biotechnol.*, 37(12):1452–1457, Dec 2019. [doi:10.1038/s41587-019-0290-0](https://doi.org/10.1038/s41587-019-0290-0).
- [9] Rewon Child, Scott Gray, Alec Radford, and Ilya Sutskever. Generating Long Sequences with Sparse Transformers. *arXiv*, April 2019. [arXiv:1904.10509](https://arxiv.org/abs/1904.10509), [doi:10.48550/arXiv.1904.10509](https://doi.org/10.48550/arXiv.1904.10509).

- [10] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Ben Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier Garcia, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayanan Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Diaz, Orhan Firat, Michele Catasta, Jason Wei, Kathy Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. PaLM: Scaling Language Modeling with Pathways. *arXiv*, April 2022. [arXiv:2204.02311](https://arxiv.org/abs/2204.02311), [doi:10.48550/arXiv.2204.02311](https://doi.org/10.48550/arXiv.2204.02311).
- [11] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. *arXiv*, October 2018. [arXiv:1810.04805](https://arxiv.org/abs/1810.04805), [doi:10.48550/arXiv.1810.04805](https://doi.org/10.48550/arXiv.1810.04805).
- [12] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. *arXiv*, October 2020. [arXiv:2010.11929](https://arxiv.org/abs/2010.11929), [doi:10.48550/arXiv.2010.11929](https://doi.org/10.48550/arXiv.2010.11929).
- [13] William Fedus, Barret Zoph, and Noam Shazeer. Switch Transformers: Scaling to Trillion Parameter Models with Simple and Efficient Sparsity. *arXiv*, January 2021. [arXiv:2101.03961](https://arxiv.org/abs/2101.03961), [doi:10.48550/arXiv.2101.03961](https://doi.org/10.48550/arXiv.2101.03961).
- [14] Fan Feng, Yuan Yao, Xue Qing David Wang, Xiaotian Zhang, and Jie Liu. Connecting high-resolution 3D chromatin organization with epigenomics. *Nat. Commun.*, 13(2054):1–10, April 2022. [doi:10.1038/s41467-022-29695-6](https://doi.org/10.1038/s41467-022-29695-6).
- [15] Mandy Guo, Joshua Ainslie, David Uthus, Santiago Ontanon, Jianmo Ni, Yun-Hsuan Sung, and Yinfei Yang. LongT5: Efficient Text-To-Text Transformer for Long Sequences. *arXiv*, December 2021. [arXiv:2112.07916](https://arxiv.org/abs/2112.07916), [doi:10.48550/arXiv.2112.07916](https://doi.org/10.48550/arXiv.2112.07916).
- [16] Gautier Izacard, Patrick Lewis, Maria Lomeli, Lucas Hosseini, Fabio Petroni, Timo Schick, Jane Dwivedi-Yu, Armand Joulin, Sebastian Riedel, and Edouard Grave. Few-shot Learning with Retrieval Augmented Language Models. *arXiv*, August 2022. [arXiv:2208.03299](https://arxiv.org/abs/2208.03299), [doi:10.48550/arXiv.2208.03299](https://doi.org/10.48550/arXiv.2208.03299).
- [17] Mika Sarkin Jain, Cecilia Dominguez Conde, Krzysztof Polanski, Xi Chen, Jongeun Park, Rachel A. Botting, Emily Stephenson, Muzlifah Haniffa, Austen

- Lamacraft, Mirjana Efremova, and Sarah A. Teichmann. MultiMAP: Dimensionality Reduction and Integration of Multimodal Data. *bioRxiv*, page 2021.02.16.431421, Feb 2021. URL: <https://doi.org/10.1101/2021.02.16.431421>, [arXiv:2021.02.16.431421](https://arxiv.org/abs/2021.02.16.431421).
- [18] Kuang-Huei Lee, Ofir Nachum, Mengjiao Yang, Lisa Lee, Daniel Freeman, Winnie Xu, Sergio Guadarrama, Ian Fischer, Eric Jang, Henryk Michalewski, and Igor Mordatch. Multi-Game Decision Transformers. *arXiv*, May 2022. [arXiv:2205.15241](https://arxiv.org/abs/2205.15241), [doi:10.48550/arXiv.2205.15241](https://doi.org/10.48550/arXiv.2205.15241).
- [19] Niannian Li, Kairang Jin, Yanmin Bai, Haifeng Fu, Lin Liu, and Bin Liu. Tn5 Transposase Applied in Genomics Research. *Int. J. Mol. Sci.*, 21(21), November 2020. [doi:10.3390/ijms21218329](https://doi.org/10.3390/ijms21218329).
- [20] Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, Rémi Leblond, Tom Eccles, James Keeling, Felix Gimeno, Agustin Dal Lago, Thomas Hubert, Peter Choy, Cyprien de Masson D’Autume, Igor Babuschkin, Xinyun Chen, Po-Sen Huang, Johannes Welbl, Sven Gowal, Alexey Cherepanov, James Molloy, Daniel J. Mankowitz, Esme Sutherland Robson, Pushmeet Kohli, Nando de Freitas, Koray Kavukcuoglu, and Oriol Vinyals. Competition-Level Code Generation with AlphaCode. *arXiv*, February 2022. [arXiv:2203.07814](https://arxiv.org/abs/2203.07814), [doi:10.48550/arXiv.2203.07814](https://doi.org/10.48550/arXiv.2203.07814).
- [21] Hanqing Liu, Jingtian Zhou, Wei Tian, Chongyuan Luo, Anna Bartlett, Andrew Aldridge, Jacinta Lucero, Julia K. Osteen, Joseph R. Nery, Huaming Chen, Angeline Rivkin, Rosa G. Castanon, Ben Clock, Yang Eric Li, Xiaomeng Hou, Olivier B. Poirion, Sebastian Preissl, Antonio Pinto-Duarte, Carolyn O’Connor, Lara Boggeman, Conor Fitzpatrick, Michael Nunn, Eran A. Mukamel, Zhuzhu Zhang, Edward M. Callaway, Bing Ren, Jesse R. Dixon, M. Margarita Behrens, and Joseph R. Ecker. DNA methylation atlas of the mouse brain at single-cell resolution. *Nature*, 598:120–128, October 2021. [doi:10.1038/s41586-020-03182-8](https://doi.org/10.1038/s41586-020-03182-8).
- [22] Ilya Loshchilov and Frank Hutter. Decoupled Weight Decay Regularization. *arXiv*, November 2017. [arXiv:1711.05101](https://arxiv.org/abs/1711.05101), [doi:10.48550/arXiv.1711.05101](https://doi.org/10.48550/arXiv.1711.05101).
- [23] Malte Luecken, Daniel Burkhardt, Robrecht Cannoodt, Christopher Lance, Aditi Agrawal, Hananeh Aliee, Ann Chen, Louise Deconinck, Angela Detweiler, Alejandro Granados, Shelly Huynh, Laura Isacco, Yang Kim, Dominik Klein, Bony De Kumar, Sunil Kuppasani, Heiko Lickert, Aaron McGeever, Joaquin Melgarejo, Honey Mekonen, Maurizio Morri, Michaela Müller, Norma Neff, Sheryl Paul, Bastian Rieck, Kaylie Schneider, Scott Steelman, Michael Sterr, Daniel Treacy, Alexander Tong, Alexandra-Chloe Villani, Guilin Wang, Jia Yan, Ce Zhang, Angela Pisco, Smita Krishnaswamy, Fabian Theis, and Jonathan M. Bloom. A sandbox for prediction and integration of DNA, RNA, and proteins in single cells. *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks*, 1, December 2021.

URL: <https://datasets-benchmarks-proceedings.neurips.cc/paper/2021/hash/158f3069a435b314a80bdcb024f8e422-Abstract-round2.html>.

- [24] Malte D. Luecken, M. Büttner, K. Chaichoompu, A. Danese, M. Interlandi, M. F. Mueller, D. C. Strobl, L. Zappia, M. Dugas, M. Colomé-Tatché, and Fabian J. Theis. Benchmarking atlas-level data integration in single-cell genomics. *Nat. Methods*, 19:41–50, January 2022. doi:10.1038/s41592-021-01336-8.
- [25] Leland McInnes, John Healy, and James Melville. UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction. *arXiv*, Feb 2018. URL: <https://arxiv.org/abs/1802.03426v3>, arXiv:1802.03426.
- [26] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient Estimation of Word Representations in Vector Space. *arXiv*, January 2013. arXiv:1301.3781, doi:10.48550/arXiv.1301.3781.
- [27] Takashi Nagano, Yaniv Lubling, Tim J. Stevens, Stefan Schoenfelder, Eitan Yaffe, Wendy Dean, Ernest D. Laue, Amos Tanay, and Peter Fraser. Single-cell Hi-C reveals cell-to-cell variability in chromosome structure. *Nature*, 502:59–64, October 2013. doi:10.1038/nature12593.
- [28] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019. URL: <http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf>.
- [29] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning Transferable Visual Models From Natural Language Supervision. *arXiv*, February 2021. arXiv:2103.00020, doi:10.48550/arXiv.2103.00020.
- [30] Scott Reed, Konrad Zolna, Emilio Parisotto, Sergio Gomez Colmenarejo, Alexander Novikov, Gabriel Barth-Maron, Mai Gimenez, Yury Sulsky, Jackie Kay, Jost Tobias Springenberg, Tom Eccles, Jake Bruce, Ali Razavi, Ashley Edwards, Nicolas Heess, Yutian Chen, Raia Hadsell, Oriol Vinyals, Mahyar Bordbar, and Nando de Freitas. A Generalist Agent. *arXiv*, May 2022. arXiv:2205.06175, doi:10.48550/arXiv.2205.06175.

- [31] Aviv Regev, Sarah A. Teichmann, Eric S. Lander, Ido Amit, Christophe Benoist, Ewan Birney, Bernd Bodenmiller, Peter Campbell, Piero Carninci, Menna Clatworthy, Hans Clevers, Bart Deplancke, Ian Dunham, James Eberwine, Roland Eils, Wolfgang Enard, Andrew Farmer, Lars Fugger, Berthold Göttgens, Nir Hacohen, Muzlifah Haniffa, Martin Hemberg, Seung Kim, Paul Klenerman, Arnold Kriegstein, Ed Lein, Sten Linnarsson, Emma Lundberg, Joakim Lundberg, Partha Majumder, John C. Marioni, Miriam Merad, Musa Mhlanga, Martijn Nawijn, Mihai Netea, Garry Nolan, Dana Pe'er, Anthony Phillipakis, Chris P. Ponting, Stephen Quake, Wolf Reik, Orit Rozenblatt-Rosen, Joshua Sanes, Rahul Satija, Ton N. Schumacher, Alex Shalek, Ehud Shapiro, Padmanee Sharma, Jay W. Shin, Oliver Stegle, Michael Stratton, Michael J. T. Stubbington, Fabian J. Theis, Matthias Uhlen, Alexander van Oudenaarden, Allon Wagner, Fiona Watt, Jonathan Weissman, Barbara Wold, Ramnik Xavier, Nir Yosef, and Human Cell Atlas Meeting Participants. The Human Cell Atlas. *eLife*, 6, 2017. doi:10.7554/eLife.27041.
- [32] Samuel G. Rodriques, Robert R. Stickels, Aleksandrina Goeva, Carly A. Martin, Evan Murray, Charles R. Vanderburg, Joshua Welch, Linlin M. Chen, Fei Chen, and Evan Z. Macosko. Slide-seq: A scalable technology for measuring genome-wide expression at high spatial resolution. *Science*, 363(6434):1463–1467, Mar 2019. doi:10.1126/science.aaw1219.
- [33] Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter. *arXiv*, October 2019. arXiv:1910.01108, doi:10.48550/arXiv.1910.01108.
- [34] Bindesh Shrestha. Single-Cell Metabolomics by Mass Spectrometry. *Methods Mol. Biol.*, 2064(1-8.):, 2020. arXiv:31565762, doi:10.1007/978-1-4939-9831-9_1.
- [35] Jay Sinha. CycleGAN for Image to Image Translation. *Jay Sinha's Blog*, April 2021. URL: <https://blog.jaysinha.me/train-your-first-cyclegan-for-image-to-image-translation>.
- [36] Tim Stuart, Andrew Butler, Paul Hoffman, Christoph Hafemeister, Efthymia Papalexi, William M. Mauck, Yuhao Hao, Marlon Stoeckius, Peter Smibert, and Rahul Satija. Comprehensive Integration of Single-Cell Data. *Cell*, 177(7):1888–1902.e21, Jun 2019. doi:10.1016/j.cell.2019.05.031.
- [37] Valentine Svensson, Roser Vento-Tormo, and Sarah A. Teichmann. Exponential scaling of single-cell RNA-seq in the last decade. *arXiv*, Apr 2017. URL: <https://arxiv.org/abs/1704.01379v2>, arXiv:1704.01379.
- [38] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention Is All You Need. *arXiv*, June 2017. arXiv:1706.03762, doi:10.48550/arXiv.1706.03762.

- [39] Isaac Virshup, Sergei Rybakov, Fabian J. Theis, Philipp Angerer, and F. Alexander Wolf. anndata: Annotated data. *bioRxiv*, page 2021.12.16.473007, December 2021. URL: <https://doi.org/10.1101/2021.12.16.473007>, [arXiv:2021.12.16.473007](https://arxiv.org/abs/2021.12.16.473007).
- [40] Pauli Virtanen, Ralf Gommers, Travis E. Oliphant, Matt Haberland, Tyler Reddy, David Cournapeau, Evgeni Burovski, Pearu Peterson, Warren Weckesser, Jonathan Bright, Stéfan J. van der Walt, Matthew Brett, Joshua Wilson, K. Jarrod Millman, Nikolay Mayorov, Andrew R. J. Nelson, Eric Jones, Robert Kern, Eric Larson, C. J. Carey, İlhan Polat, Yu Feng, Eric W. Moore, Jake VanderPlas, Denis Laxalde, Josef Perktold, Robert Cimrman, Ian Henriksen, E. A. Quintero, Charles R. Harris, Anne M. Archibald, Antônio H. Ribeiro, Fabian Pedregosa, and Paul van Mulbregt. SciPy 1.0: fundamental algorithms for scientific computing in Python. *Nat. Methods*, 17:261–272, March 2020. [doi:10.1038/s41592-019-0686-2](https://doi.org/10.1038/s41592-019-0686-2).
- [41] Joshua D. Welch, Velina Kozareva, Ashley Ferreira, Charles Vanderburg, Carly Martin, and Evan Z. Macosko. Single-Cell Multi-omic Integration Compares and Contrasts Features of Brain Cell Identity. *Cell*, 177(7):1873–1887.e17, 2021. [doi:10.1016/j.cell.2019.05.006](https://doi.org/10.1016/j.cell.2019.05.006).
- [42] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. HuggingFace’s Transformers: State-of-the-art Natural Language Processing. *arXiv*, October 2019. [arXiv:1910.03771](https://arxiv.org/abs/1910.03771), [doi:10.48550/arXiv.1910.03771](https://doi.org/10.48550/arXiv.1910.03771).
- [43] Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gunjan Baid, Zirui Wang, Vijay Vasudevan, Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, Ben Hutchinson, Wei Han, Zarana Parekh, Xin Li, Han Zhang, Jason Baldridge, and Yonghui Wu. Scaling Autoregressive Models for Content-Rich Text-to-Image Generation. *arXiv*, June 2022. [arXiv:2206.10789](https://arxiv.org/abs/2206.10789), [doi:10.48550/arXiv.2206.10789](https://doi.org/10.48550/arXiv.2206.10789).
- [44] Manzil Zaheer, Guru Guruganesh, Avinava Dubey, Joshua Ainslie, Chris Alberti, Santiago Ontanon, Philip Pham, Anirudh Ravula, Qifan Wang, Li Yang, and Amr Ahmed. Big Bird: Transformers for Longer Sequences. *arXiv*, July 2020. [arXiv:2007.14062](https://arxiv.org/abs/2007.14062), [doi:10.48550/arXiv.2007.14062](https://doi.org/10.48550/arXiv.2007.14062).
- [45] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A. Efros. Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks. *arXiv*, Mar 2017. URL: <https://arxiv.org/abs/1703.10593v7>, [arXiv:1703.10593](https://arxiv.org/abs/1703.10593).