

# Textomics: A Dataset for Genomics Data Summary Generation

Anonymous ACL submission

## Abstract

Summarizing biomedical discovery from genomics data using natural languages is an essential step in biomedical research but is mostly done manually. Here, we introduce Textomics, a novel dataset of genomics data description, which contains 22,273 pairs of genomics data matrix and its summary. Each summary is written by the researchers who generated the data and associated with a scientific paper. Based on this dataset, we study two novel tasks: generating textual summary from genomics data matrix and vice versa. Inspired by the successful applications of  $k$  nearest neighbors in modeling genomics data, We propose a  $k$ NN-Vec2Text model to address these tasks and observe substantial improvement on our dataset. We further illustrate how Textomics can be used to advance other applications, including evaluating scientific paper embeddings and generating masked templates for scientific paper understanding. Textomics serves as the first benchmark for generating textual summary for genomics data and we envision it will be broadly applied to other biomedical and natural language processing applications.

## 1 Introduction

Modern genomics research has become increasingly automated through being roughly divided into three sequential steps: next-generation sequencing technology produces a massive amount of genomics data, which are in turn processed by bioinformatics tools to identify key variants and genes, and, ultimately, analyzed by biologists to summarize the discovery (Goodwin et al., 2016; Kanehisa and Bork, 2003). In contrast to the first two steps that have been automated by new technologies and software, the last step of summarizing discovery is still largely performed manually, substantially slowing down the progress of scientific discovery

(Hwang et al., 2018). A plausible solution is to automatically summarize the discovery from genomics data using neural text generation, which has been successfully applied to radiology report generation (Wang et al., 2021; Yuan et al., 2019) and clinical notes generation (Melamud and Shivade, 2019; Lee, 2018; Miura et al., 2021).

In this paper, we study this novel task of generating sentences to summarize a genomics data matrix. There are several existing approaches that demonstrate encouraging results in generating short phrases to describe functions of a set of genes (Wang et al., 2018; Zhang et al., 2020; Kramer et al., 2014). However, our task is fundamentally different from these ones: the input of our task is a matrix that contains tens of thousands genes, which could be more noisy than a set of selected genes; the output of our task is sentences instead of short phrases or controlled vocabularies.

To study this task, we curate a novel dataset, Textomics, by integrating data from PMC, PubMed, and Gene Expression Omnibus (GEO) (Edgar et al., 2002) (Figure 1). GEO is the default database repository for researchers to upload their genomics data matrix, such as gene expression matrix and mutation matrix. Each genomics data matrix in GEO is a sample by feature matrix, where samples are often humans or mice that are sequenced together to study a specific biological problem and features are genes or variants. Each matrix is also associated with a few sentences that are written by researchers to summarize this data matrix. After pre-processing, we obtain 22,273 matrix summary pairs, spanning 9 sequencing technology platforms. Each matrix has on average 2,475 samples and 22,796 features. Each summary has on average 46 words.

We further propose a novel approach to automatically generate summary from a genomics data matrix, which is highly noisy and high-dimensional.  $k$

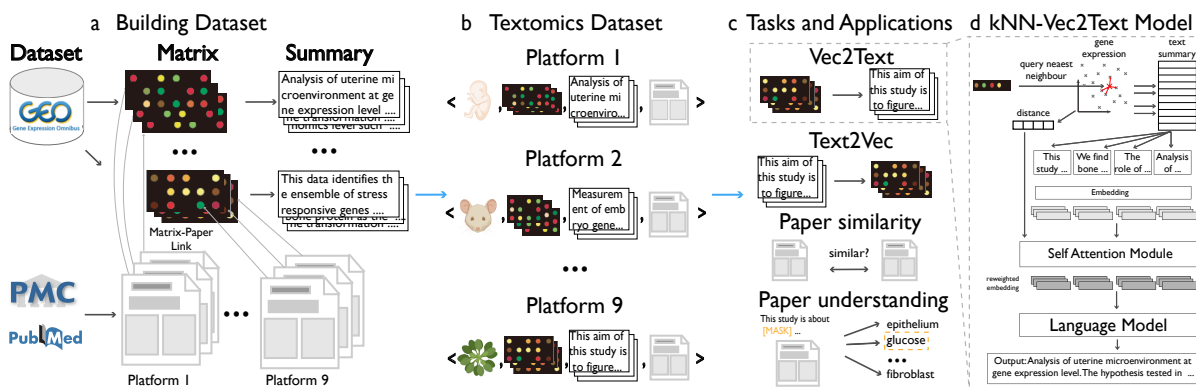


Figure 1: **Flow chart of Textomics.** a. Genomics data matrices and summaries are collected from GEO. Scientific papers are collected from PMC and PubMed. Each data matrix is associated with a unique summary and a unique scientific paper in Textomics. b. Textomics is divided into 9 sequencing platforms, spanning over various species. Data matrices in the same platforms share the same features and can therefore be used to train a machine learning model. c. Textomics can be used as the benchmark for a variety of tasks, including Vec2Text, Text2Vec, measuring paper similarity, and scientific paper understanding. d.  $k$ NN-Vec2Text is developed to address the task of Vec2Text, by first constructing a reference summary using similar genomics data matrix and then unifying these summaries to generate a new summary.

nearest neighbor ( $k$ NN) approaches have obtained great success in genomics data by capturing the hidden modules within it (Levine et al., 2015; Baran et al., 2019). The key idea of our method is to find  $k$  nearest summaries according to the genomics data similarity and then exploit attention mechanism to convert these  $k$  nearest summaries to a new summary. Our method obtained substantial improvement in comparison to baseline approaches. We further illustrated how we can generate a genomics data matrix from a given summary, offering the possibility to simulate genomics data from textual description. We then introduced how Textomics can be used as a novel benchmark for measuring scientific paper similarity and evaluating scientific paper understanding. To the best of our knowledge, Textomics and  $k$ NN-Vec2Text together build up the first large-scale benchmark for genomics data summary generation, and can be broadly applied to a variety of natural language processing tasks.

## 2 Textomics Dataset

We collected genomics data matrices from Gene Expression Omnibus (GEO) (Edgar et al., 2002). The feature of each data matrix is a gene or a variant and the sample of each matrix is an experimental subject, such as an experimental animal or a patient. Each data matrix is associated with an expert-written summary, describing this data matrix. We obtained in total 164,667 matrix-summary pairs, spanning 12,219 sequencing platforms. We

truncated the summary that is longer than 64 words.

Data matrices belonging to the same sequencing platform share the same set of features, and can thus be used together to train the model. To this end, we first selected 20,000 features that have the largest standard deviation and lower missing rate for each platform and excluded samples that have a substantially higher missing rate. We then selected 9 platforms with the lowest rate of missing values and the largest number of matrix-summary pairs. We imputed the resulted data matrix using averaging imputation and excluded outliers and non-informative summary (e.g., “Please see our data below”) through both manual inspection and an automated approach that excluded the summary that is substantially different from all other summaries based on pairwise BLEU scores. Finally, each of the 9 platforms contains 471 matrix-summary pairs on average, presenting a desirable number of training samples to develop data summary generation models. We summarized the statistics of these 9 platforms in **Supplementary Table S1**.

Data matrices belonging to the same platform have distinct samples (e.g., patient samples collected from two hospitals). In order to make them comparable and provide fixed-size features for machine learning models, we used a five-number summary to represent each data matrix. In particular, we calculated the smallest, the first quartile, the median, the third quartile, and the largest value of each feature across samples in a specific data

matrix. We then concatenated these values of all features, resulting in a 100k-dimensional feature vector for each data matrix. This vector will be used as the input to the machine learning model. We used the original summary written by the author as the output of the machine learning model.

Each data matrix is associated with a scientific paper, which describes how the authors generated and used the data. Therefore, the data matrix and the summary can be used to help embed these papers. We additionally retrieved these papers from PubMed and PMC databases according to the paper titles enclosed in GEO. We obtained the full text for those 7,691 freely accessible ones. We will introduce two applications that jointly use scientific papers and matrix-summary pairs in Section 6.

### 3 Task Description

We aim to accelerate genomics discovery by generating a textual summary given the five-number summary-based vector of a genomics data matrix. We refer to the five-number summary-based vector as gene feature vector for simplicity. Specifically, consider textual summary domain  $\mathcal{D}$  and gene feature vector domain  $\mathcal{V}$ , let  $\mathbf{D} = \{\mathbf{D}_{\mathcal{D}}, \mathbf{D}_{\mathcal{V}}\} = \{(d_i, v_i)\}_{i=1}^N \stackrel{dist}{\sim} \mathbb{P}(\mathcal{D}, \mathcal{V})$  be a dataset contains  $N$  summary-vector pairs sampled from the joint distribution of these two domains, where  $d_i \triangleq \langle d_i^1, d_i^2, \dots, d_i^{n_{d_i}} \rangle$  denotes a token sequence and  $v_i \in R^{l_v}$  denotes the gene feature vector. Here  $d_i^j \in C$ ,  $C$  is the vocabulary. We now formally de-

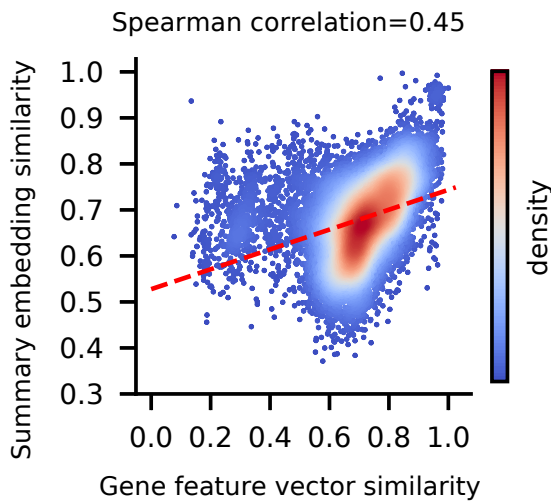


Figure 2: Density plot showing the Spearman correlation between text-based similarity (y-axis) and vector-based similarity (x-axis) on sequencing platform GPL6246. Each dot is a pair of data samples.

fine two cross-domain generation tasks, Vec2Text and Text2Vec, based on our dataset. Given a gene feature vector  $v_i$ , Vec2Text aims to generate a summary  $d_i$  that could best describe this vector  $v_i$ ; given a textual summary  $d_i$ , Text2Vec aims to generate the gene feature vector  $v_i$  that  $d_i$  describes. Since we are studying a novel task on a novel dataset, we first examined the feasibility of this task. To this end, we obtained the dense representation of each textual summary using the pre-trained SPECTER model (Cohan et al., 2020) and use these representations to calculate a summary-based similarity between each pair of summaries. We also calculated a vector-based similarity based on the gene feature vector using the cosine similarity. We found that these two similarity measurements show a substantial agreement (Figure 2, Supplementary Table S2). All 9 platforms achieved a Spearman correlation greater than 0.2, suggesting the possibility to generate textual summary from the gene feature vector and vice versa.

## 4 Methods

### 4.1 Vec2Text

We first introduce a base model that tries to encode gene expression vectors into the semantic embedding space and then decodes it to generate texts. The base model contains a word embedding function  $\text{Emb}(\cdot)$ , a gene feature vector encoder  $\text{Enc}_v(\cdot)$  and a decoder  $\text{Dec}_v(\cdot)$ . Given a gene feature vector  $v_i$ , the encoder will first embed the data into a semantic representation space  $s_i^{(0)} = \text{Enc}_v(v_i)$ , and then the decoder will start from this representation for the text generation. The generation process is autoregressive. It generates  $j$ -th word  $\hat{d}_i^{(j)}$  and its embedding  $s_i^{(j)}$  as:

$$P(\hat{d}_i^{(j)} | s_i^{(<j)}) = \text{Dec}_v(s_i^{(<j)}), j = 1, \dots, n_{d_i}. \quad (1)$$

Then we sample the next word and obtain its embedding as:

$$s_i^{(j)} = \text{Emb}(\hat{d}_i^{(j)}), \hat{d}_i^{(j)} \stackrel{sample}{\sim} P(\hat{d}_i^{(j)} | s_i^{(<j)}). \quad (2)$$

This model is trained using the following loss function:

$$\mathcal{L}_{\text{base}} = -\frac{1}{|\mathbf{D}_{\mathcal{V}}|} \sum_{i=1}^{|\mathbf{D}_{\mathcal{V}}|} \sum_{j=1}^{n_{d_i}} \log P(\hat{d}_i^{(j)} | s_i^{(<j)}). \quad (3)$$

#### 4.1.1 $k$ NN-Vec2Text Model

The base model attempts to learn an encoder that projects a gene feature vector to a semantic representation. However, the substantial noise and the high-dimensionality of the gene feature vector pose

great challenges to effectively learn that projection.  $k$ -nearest neighbors models have been extensively used as the solution to overcome such issues in genomics data analysis (Levine et al., 2015; Baran et al., 2019). Therefore, one plausible solution is to explicitly leverage summaries from similar gene feature vectors to improve the generation. Inspired by the encouraging performance in using  $k$ -nearest neighbors ( $k$ NN) in seq2seq models (Khandelwal et al., 2019, 2021) and genomics data analysis (Levine et al., 2015; Baran et al., 2019), we propose to convert the Vec2Text problem to a Text2Text problem according to the  $k$ -nearest neighbor of each vector.

For a given gene feature vector  $g$ , we use  $e_i$  to denote its Euclidean distance to another gene feature vectors  $v_i$  in  $\mathbf{D}$ . We then select the summaries of  $k$  samples that have the minimum Euclidean distances as the reference summary list  $\tilde{\mathbf{t}} = [d_{j_1}, \dots, d_{j_k}]$ , where  $j_m \in \{1, 2, \dots, |\mathbf{D}|\}$  denotes the index of ordered summaries w.r.t the Euclidean distance, i.e.,  $e_{j_1} \leq e_{j_2} \leq \dots \leq e_{j_{|\mathbf{D}|}}$ .

In addition to alleviating the noise in genomics data using the reference summary list (Levine et al., 2015; Baran et al., 2019), our method explicitly converts the Vec2Text problem to a Text2Text problem, and can thus seamlessly incorporate many advanced pre-trained language models into our framework. The resulted problem we need to solve is a  $k$  sources to one target generation problem. One naive solution is to concatenate the  $k$  reference summaries together. However, this concatenation will make the source text much longer than the target text and how to order each summary during concatenation also remains unclear. Instead, we propose to transform this problem into  $k$  one-to-one generation problem and then use attention-based strategy to fuse them. Concretely, let  $\mathbf{n}_j = \max\{n_{j_1}, \dots, n_{j_k}\}$  be the maximum length among all the reference summaries. We first get the representation of summaries  $x_{j_m} = \text{Emb}(d_{j_m}) = \langle x_{j_m}^{(1)}, \dots, x_{j_m}^{(\mathbf{n}_j)} \rangle$  for  $m = 1, \dots, k$ . We construct fixed-length reference summaries by padding after the end of each summary with length less than  $\mathbf{n}_j$ . We then utilize self-attention module (SA) (Vaswani et al., 2017) to get the aggregated embedding of each reference with their embeddings as well as the gene feature vector distance  $e_i$ . Let  $Q_r, K_r, V_r$  be the query, key, value matrix of embedding sequence  $r = \langle r^{(1)}, \dots, r^{(l_r)} \rangle$ , we have:

$$\text{SA}(r) = \text{Attention}(Q_r, K_r, V_r). \quad (4)$$

We then calculate the attention score as following:

$$a_{j_m} = \text{SA}(\langle x_{j_m}^{(1)}, \dots, x_{j_m}^{(\mathbf{n}_{j_k})} \rangle), \quad (5)$$

$$\text{sc}_j = \text{SA}(\langle e_{j_1} \cdot a_{j_1}, \dots, e_{j_k} \cdot a_{j_k} \rangle), \quad (6)$$

where  $\text{sc}_j = [\text{sc}_{j_1}, \dots, \text{sc}_{j_k}] \in R^k$ . The final score is then calculated based on the attention scores and temperature  $\tau$  as:

$$w_{j_m} = \frac{\exp(\tau \cdot \text{sc}_{j_m})}{\sum_{l=1}^k \exp(\tau \cdot \text{sc}_{j_l})}. \quad (7)$$

Then, we aggregate embedding sequences by taking weighted averages:

$$\tilde{x}_j^{(l)} = \sum_{m=1}^k w_{j_m} x_{j_m}^{(l)}, l = 1, \dots, \mathbf{n}_j. \quad (8)$$

Let  $P_{<l,x}(d) = P_{\theta_{LM}}(d^{(l)} | d^{(<l)}, x)$ ,  $0 < l < n_d$  be the probability distribution of  $d^{(l)}$  output by the language model  $\theta_{LM}$  conditioned on the sequences of the embedding vectors  $x$  and the first  $l-1$  sequence tokens. We feed the aggregated embedding sequences into the language model to reconstruct the summary  $d$  using an autoregressive-based loss function:

$$\mathcal{L}_{k\text{NN-Vec2Text}} = -\frac{1}{|\mathbf{D}_{\mathcal{D}}|} \sum_{d \in \mathbf{D}_{\mathcal{D}}} \sum_{l=1}^{n_d} \log P_{<l, \tilde{x}_j}(d). \quad (9)$$

## 4.2 Text2Vec

We model the reverse problem of generating the gene feature vector  $v$  from a textual summary  $d$  as a regression problem. Our model is composed with a semantic encoder  $\text{Enc}_d(\cdot)$  and a readout head  $\text{MLP}(\cdot)$ . Specifically, the encoder will embed the textual summary into dense representation  $x = \text{Enc}_d(d)$ , and the readout head will map the representation to the gene feature vector  $\hat{v} = \text{MLP}(x)$ . Then we train this model by minimizing the mean square errors:

$$\mathcal{L}_v = \sqrt{\frac{1}{|\mathbf{D}_{\mathcal{V}}|} \sum_{v_i \in \mathbf{D}_{\mathcal{V}}} \frac{1}{l_d} \sum_{j=1}^{l_d} (\hat{v}_i^{(j)} - v_i^{(j)})^2}. \quad (10)$$

## 5 Results

### 5.1 Vec2Text

To evaluate the performance of  $k$ NN-Vec2Text on the task of Vec2Text, we compared it to the base model based on Transformer (Vaswani et al., 2017) and GPT-2 (Radford et al., 2019), as well as Sent-VAE (Bowman et al., 2016). For  $k$ NN-Vec2Text, we set  $k = 4$  and  $\tau = 0.1$ , and used T5

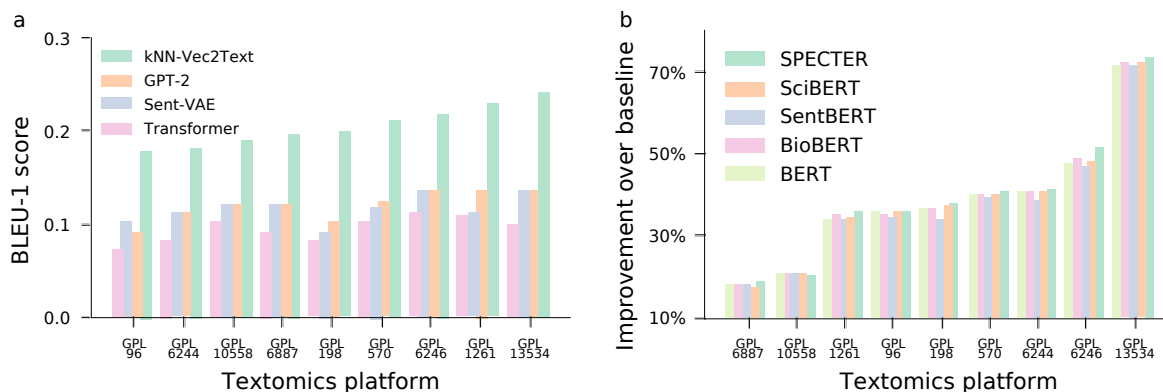


Figure 3: Performance on Vex2Text and Text2Vec using Textomics as the benchmark. a. Bar plot comparing our method  $k$ NN-Vec2Text with existing approaches on the task of Vec2Text across 9 platforms in Textomics. b. Bar plot comparing the performance of different scientific paper embedding methods across 9 platforms in Textomics.

Table 1: A case study of the generated text by  $k$ NN-Vec2Text. Summaries of the four nearest neighbors in the input space are shown. The generated text is composed of short spans from four different neighbors (colored in red).

Neighbor 1:	Analysis of B16 tumor microenvironments at gene expression level. The hypothesis tested in the present study was that Tregs orchestrated the immune response triggered in presence of tumors.
Neighbor 2:	This study aims to look at gene expression profiles between wildtype and Bapx1 knockout cells of the gut in a E12.5 mouse embryo.
Neighbor 3:	The role of bone morphogenetic protein2 in regulating transformation of the uterine stroma during embryo implantation in mice was investigated by the conditional ablation of Bmp2 in the uterus using the mouse.
Neighbor 4:	Measurement of specific gene expression in clinical samples is a promising approach for monitoring the recipient immune status to the graft in organ transplantation.
Generated:	Analysis of uterine microenvironment at gene expression level. The hypothesis tested in the present study was that Tregs orchestrated the immune response triggered in presence of embryo.
Truth:	Analysis of uterine microenvironment at gene expression level. The hypothesis tested in the present study was that Tregs orchestrated the immune response triggered in presence of embryo.

(Raffel et al., 2020) as the language model. For all 9 platforms, we reported the average performance under 5-fold cross validation. The results of BLEU-1 score are summarized in **Figure 3a**. We found that  $k$ NN-Vec2Text substantially outperformed other methods by a large margin. Specifically,  $k$ NN-Vec2Text obtained a 0.206 BLEU-1 score on average while none of the other three methods achieved an average BLEU-1 score greater than 0.150. The prominent performance of our method demonstrates the effectiveness of using a  $k$ -nearest neighbor approach to convert the Vec2Text problem to a Text2Text problem.

To further understand the superior performance of the  $k$ NN-Vec2Text model, we presented a case study in **Table 1**. In this case study, the generated summary is highly accurate compared to the ground truth summary. By examining the summaries of the 4 nearest neighbors in the gene feature vector space, we found that the generated summary is composed of short spans from each individual neighbor, again indicating the advantage of using a  $k$ -nearest neighbor for this task. Our method

leveraged an attention mechanism to unify these four neighbors, thus offering an accurate generation. We also observed consistent improvement of our method over comparison approaches on other metrics and summarized the results in **Supplementary Table S3**.

## 5.2 Text2Vec

We next used the Text2Vec task to illustrate how our dataset can be used to compare the performance of different pre-trained language models. In particular, we compared a recently proposed scientific paper embedding method SPECTER (Cohan et al., 2020), which has demonstrated prominent performance in a variety of scientific paper analysis tasks, with SciBERT (Beltagy et al., 2019), BioBERT (Lee et al., 2020) and SentBERT (Wang and Kuo, 2020) and the vanilla BERT (Devlin et al., 2019). While the other language models directly take the token sequence as the input, SPECTER model needs to take both the abstract and the title. To make a fair comparison, we concatenated the title and the summary as the input for models other than SPECTER. For all 9 platforms, we re-

ported the average performance under 5-fold cross validation. We further implemented a simple averaging baseline approach that predicts the vector for a test summary according to the average vectors of training samples. This baseline does not utilize any textual summary and can thus help us assess the effect of using textual summary information in this task. We used RMSE to evaluate the performance of all methods. We reported the RMSE improvement of each method over the averaging baseline model in **Figure 3b**. We found that all methods outperform the baseline approaches by gaining at least 15% improvement, indicating the importance of considering textual summary in this task. SPECTER achieved the best overall performance among all five methods, suggesting the advantage to separately model the title and the abstract when embedding scientific papers.

## 6 Applications

### 6.1 Evaluate paper embedding via Textomics

Embedding scientific papers is crucial to effectively identify emerging research topics and new knowledge from scientific literature. To this end, many machine learning models have been proposed to embed scientific papers into dense embeddings and then applied these embeddings for a variety of downstream applications (Cohan et al., 2020; Lee et al., 2020; Wang and Kuo, 2020; Beltagy et al., 2019; Devlin et al., 2019). However, there is currently limited golden standard that can measure the similarity between two papers. As a result, existing approaches use surrogate metrics such as citation relationship, keywords, and user activities to evaluate their paper embeddings (Cohan et al., 2020; Chen et al., 2019; Wang et al., 2019).

Textomics can be used to measure these paper embedding approaches by examining the consistency between the embedding-based paper similarity and the embedding-based summary similarity since both the paper and the summary are written by the same authors. In particular, for a pair of summaries  $d_i, d_j \in \mathbf{D}_{\mathcal{D}}$ , let  $t_i, t_j$  be the text (e.g., abstracts) extracted from their corresponding scientific papers. Let  $\text{Enc}_d$  be the encoder of the paper embedding method we want to evaluate. We first get their embeddings as:

$$s_{d_i}, s_{d_j} = \text{Enc}_d(d_i), \text{Enc}_d(d_j) \in R^{l_s}, \quad (11)$$

$$s_{t_i}, s_{t_j} = \text{Enc}_d(t_i), \text{Enc}_d(t_j) \in R^{l_s}. \quad (12)$$

We then compute the pairwise Euclidean distance

between all pairs of summaries and all pairs of paper text as:

$$s_{d_i, j} = \sqrt{\sum_{k=1}^{l_s} (s_{d_i}^{(k)} - s_{d_j}^{(k)})^2} \in R, \quad (13)$$

$$s_{t_i, j} = \sqrt{\sum_{k=1}^{l_s} (s_{t_i}^{(k)} - s_{t_j}^{(k)})^2} \in R. \quad (14)$$

To evaluate the quality of the encoder  $\text{Enc}_d$ , we can calculate the Spearman correlation between the pairwise summary similarity and the pairwise text similarity. A larger Spearman correlation indicates this  $\text{Enc}_d$  is more accurate in embedding scientific papers. As a proof-of-concept, we obtained the full text of 7,691 papers in our dataset from the freely accessible PubMed Central. We segmented each paper into five sections of abstract, introduction, method, result and conclusion. We first compared different paper embedding methods using the abstract of a paper. The five embedding methods we considered are introduced in section 5.1. Since SPECTER takes both the title and paragraph as the input we used the first sentence of the summary as a pseudo-title when encoding the summary. The results are summarized in **Figure 4a**. We found that SPECTER was substantially better than other methods on 8 out of the 9 platforms. SPECTER is specifically developed to embed scientific papers by processing the title and the abstract separately, whereas other pre-trained language models simply concatenated the title and the abstract. The superior performance of SPECTER suggests the importance of separately modeling paper title and abstract when embedding scientific papers. SentBERT obtained the best performance among four pre-trained language models, partially due to its prominent performance in sentence-level embedding. We further noticed that the relative performance among different methods is largely consistent with the previous work evaluated on other metrics (Cohan et al., 2020), demonstrating the high-quality of Textomics.

After observing the superior performance of SPECTER, we next investigated which section of the paper can be best used to assess paper similarity. Although existing paper embedding approaches often leverage the abstract for embedding, other sections, such as introduction and results might also be informative, especially for paper describing a specific dataset or method. We thus applied SPECTER to embed five different sections of each scientific

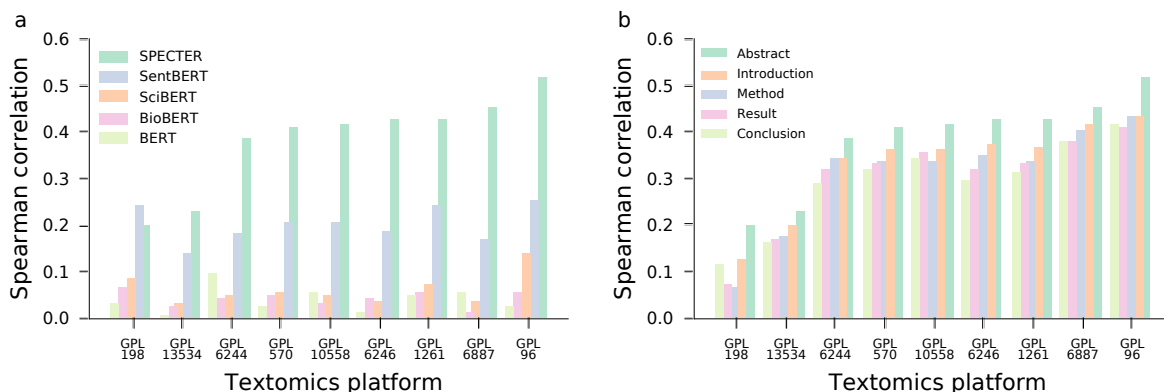


Figure 4: Performance on using Textomics as the benchmark to evaluate scientific paper embeddings. (A). Bar plot showing the comparison on embedding scientific papers using Textomics as the benchmark. (B). Bar plot showing the comparison on SPECTER embedding of different paper sections using Textomics as the benchmark.

paper and used Textomics to evaluate which section can best reflect paper similarity. We observed a consistent improvement of using the abstract section in comparison to other paper sections (Figure 4B), which is consistent with the intuition that the abstract represents a good summary of the scientific paper, again indicating the reliability of using Textomics to evaluate paper embedding methods.

## 6.2 Scientific paper understanding

Creating masked sentences and then filling in these masks can examine whether the machine learning model has properly understood a scientific paper. However, one challenge in such research is how to generate masked sentences that are relevant to a given paper while also ensuring the answer is enclosed in the paper. Our dataset could be used to automatically generate such masked sentences using the summary, which is highly relevant to the paper but also not overlapped with the paper. In particular, we can mask out keywords from the summary and then use this masked summary as the question and let a machine learning model to find the answer from the non-overlapping scientific paper. Let  $C_{\text{bio}}$  be a dictionary that contains biological keywords we want to mask out from the summary,  $(d_i, t_i)$  be a pair of textual summary and paragraph text extracted from its corresponding scientific paper. If the  $j$ -th word  $w_i = d_i^{(j)} \in C_{\text{bio}}$  in the summary belongs to  $C_{\text{bio}}$ , our proposed task is to predict which word in  $C_{\text{bio}}$  is the missing word in  $d_{\text{masked}}$  given  $t_i$ . The masked summary  $d_{\text{masked}}$  is the same as  $d_i$  except its  $j$ -th word is substituted with [PAD]. For simplicity, we only mask at most one token in  $d_i$ . We therefore form our task as a multi-class classification problem. Sim-

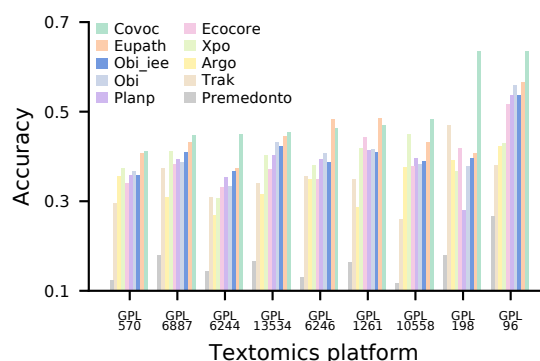


Figure 5: Bar plot showing the accuracy of filling the masked sentences of ten biomedical categories across 9 platforms using Textomics as the benchmark.

ilar to section 6.1, we used the paper abstract as the paragraph text  $t_i$ . To generate  $C_{\text{bio}}$ , we leveraged a recently developed biological terminology dataset Graphine (Liu et al., 2021), which provides the biological phrases spanning 227 categories. We selected 10 categories that can produce the largest number of masked sentences in Textomics. We manually filtered ambiguous words and stop words. On average, each category contains 317 keywords. We used a fully connected neural network to perform the multi-class classification task. The input feature is the concatenation of the masked summary embedding and the paragraph embedding. We used SPECTER to derive these embeddings as it has obtained the best performance in our previous analysis. The results are summarized in Figure 5. We observed high accuracy on all ten categories, which are much better than the 0.4% accuracy by random guessing, indicating the usefulness of our benchmark in scientific paper understanding. Finally, we found that the performance of each category varied

700 across different platforms, suggesting the possibil-  
701 ity to further improve the performance by jointly  
702 learning from all platforms.

## 704 7 Related work

705 Our task is related to existing works that take a  
706 structured data as the input and then generate the  
707 unstructured text. Different input data modalities  
708 and related datasets have been considered in the  
709 literature, including text triplets in RDF graphs  
710 (Gardent et al., 2017; Ribeiro et al., 2020; Song  
711 et al., 2021; Chen et al., 2020)), text-data tables  
712 (Lebret et al., 2016; Rebuffel et al., 2021; Dusek  
713 et al., 2019; Rebuffel et al., 2019; Puduppully and  
714 Lapata, 2021; Chen et al., 2020), electronic medical  
715 records (Lee, 2018; Guan et al., 2018), radiology  
716 reports (Wang et al., 2021; Yuan et al., 2019; Miura  
717 et al., 2021), and other continuous data modalities  
718 without explicit textual structures such as image  
719 (Lin et al., 2015; Cornia et al., 2020; Ke et al.,  
720 2019; Radford et al., 2021), audio (Drossos et al.,  
721 2019; Manco et al., 2021; Wu et al., 2021; Mei  
722 et al., 2021), and video (Li et al., 2021; Ging et al.,  
723 2020; Zhou et al., 2018; Li et al., 2020). Different  
724 from these structures, our dataset takes a high di-  
725 mensional genomics feature matrix as input, which  
726 doesn't exhibit structure and thus substantial differ-  
727 ent from other modalities. Moreover, our dataset  
728 is the first dataset that aims to convert genomics  
729 feature vector to textual summary. The substantial  
730 noise and high-dimensionality of genomics data  
731 matrix further pose unique challenges in text gen-  
732 eration.

733 Our  $k$ NN-Vec2Text model is inspired by the re-  
734 cent success in applying  $k$ NN-based language mod-  
735 els to machine translation (Khandelwal et al., 2021)  
736 and language models (Khandelwal et al., 2019; He  
737 et al., 2021; Ton et al., 2021). The main differ-  
738 ence between our methods and their approaches is  
739 that while we try to leverage  $k$ NN in the genomics  
740 vector space to construct reference texts, they use  
741  $k$ NN in the text embedding space during the au-  
742 toregressive generation process to help adjust the  
743 sample distribution. There are some other methods  
744 that can be used to generate text from vectors, such  
745 as (Bowman et al., 2016; Song et al., 2019; Miao  
746 and Blunsom, 2016; Montero et al., 2021; Zhang  
747 et al., 2019). Their inputs are latent vectors that  
748 need to be inferred from the data and do not have  
749 specific meanings, which are different from our  
gene feature vectors.

## 750 8 Conclusion and future work

751 In this paper, we have proposed a novel dataset  
752 Textomics, containing 22,273 pairs of genomics  
753 matrix and its corresponding textual summary. We  
754 then introduce a novel task of Vec2Text based on  
755 our dataset. This task aims to generate the tex-  
756 tual summary based on the gene feature vector.  
757 To address this task, we propose a novel method  
758  $k$ NN-Vec2Text, which constructs the reference text  
759 using nearest neighbours in the gene feature vector  
760 space and then generates a new summary accord-  
761 ing to this reference text. We further introduce  
762 two applications that can be advanced using our  
763 dataset. One application aims at evaluating sci-  
764 entific paper similarity according to the similarity  
765 of its corresponding data summary, and the other  
766 application leverages our dataset to automatically  
767 generate masked sentences for scientific paper un-  
768 derstanding.

769 Our method searches for the nearest neighbours  
770 by calculating the Euclidean distance between five-  
771 number summary vectors of the genomics feature  
772 matrix. However, this might lose useful informa-  
773 tion lied in the original matrix. It's worth exploring  
774 end-to-end approaches that can learn embeddings  
775 from the genomics feature matrix instead of repre-  
776 senting them as five-number summary vectors. On  
777 the Text2Vec side, we are interested in extending  
778 our work to directly generate the whole genomics  
779 feature matrix instead of the five-number summary  
780 vectors. Also, it would be interesting to jointly  
781 learn the Text2Vec and the Vec2Text tasks, and  
782 one potential solution is to further decode the gen-  
783 erated vector to reconstruct the embedding of the  
784 summaries in Text2Vec, and leverage the resulted  
785 decoder to predict the embedding of text by using  
786  $k$ NN method in the text embedding space.

787 To the best of our knowledge, Textomics and  
788  $k$ NN-Vec2Text serves as the first large-scale ge-  
789 nomics data description benchmark, and we en-  
790 vision it will be broadly applied to other natural  
791 language processing and biomedical tasks. On  
792 the biomedical side, summaries in the Textomics  
793 dataset could be used to impute experimentally  
794 measured gene expression data matrix and serve as  
795 additional features in classifying these genomics  
796 feature data. On the NLP side, Textomics could  
797 also be used to help scientific paper analysis tasks,  
798 such as paper recommendation (Bai et al., 2020),  
799 citation text generation (Luu et al., 2020), and cita-  
tion prediction (Suzen et al., 2021).



## References

- Xiaomei Bai, Mengyang Wang, Ivan Lee, Zhuo Yang, Xiangjie Kong, and Feng Xia. 2020. [Scientific paper recommendation: A survey](#). 850
- Yael Baran, Akhiad Bercovich, Arnau Sebe-Pedros, Yaniv Lubling, Amir Giladi, Elad Chomsky, Zohar Meir, Michael Hoichman, Aviezer Lifshitz, and Amos Tanay. 2019. Metacell: analysis of single-cell rna-seq data using k-nn graph partitions. *Genome biology*, 20(1):1–19. 851
- Iz Beltagy, Kyle Lo, and Arman Cohan. 2019. [Scibert: A pretrained language model for scientific text](#). 852
- Samuel R. Bowman, Luke Vilnis, Oriol Vinyals, Andrew M. Dai, Rafal Jozefowicz, and Samy Bengio. 2016. [Generating sentences from a continuous space](#). 853
- Liquan Chen, Guoyin Wang, Chenyang Tao, Dinghan Shen, Pengyu Cheng, Xinyuan Zhang, Wenlin Wang, Yizhe Zhang, and Lawrence Carin. 2019. [Improving textual network embedding with global attention via optimal transport](#). 854
- Wenhu Chen, Yu Su, Xifeng Yan, and William Yang Wang. 2020. [Kgpt: Knowledge-grounded pre-training for data-to-text generation](#). 855
- Arman Cohan, Sergey Feldman, Iz Beltagy, Doug Downey, and Daniel S. Weld. 2020. [Specter: Document-level representation learning using citation-informed transformers](#). 856
- Marcella Cornia, Matteo Stefanini, Lorenzo Baraldi, and Rita Cucchiara. 2020. [Meshed-memory transformer for image captioning](#). 857
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [Bert: Pre-training of deep bidirectional transformers for language understanding](#). 858
- Konstantinos Drossos, Samuel Lipping, and Tuomas Virtanen. 2019. [Clotho: An audio captioning dataset](#). 859
- Ondrej Dusek, Jekaterina Novikova, and Verena Rieser. 2019. [Evaluating the state-of-the-art of end-to-end natural language generation: The E2E NLG challenge](#). *CoRR*, abs/1901.07931. 860
- Ron Edgar, Michael Domrachev, and Alex E. Lash. 2002. Gene expression omnibus: Ncbi gene expression and hybridization array data repository. *Nucleic acids research*, 30 1. 861
- Claire Gardent, Anastasia Shimorina, Shashi Narayan, and Laura Perez-Beltrachini. 2017. [Creating training corpora for NLG micro-planners](#). In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 179–188, Vancouver, Canada. Association for Computational Linguistics. 862
- Simon Ging, Mohammadreza Zolfaghari, Hamed Pirsiavash, and Thomas Brox. 2020. [Coot: Cooperative hierarchical transformer for video-text representation learning](#). 863
- Sara Goodwin, John D McPherson, and W Richard McCombie. 2016. Coming of age: ten years of next-generation sequencing technologies. *Nature Reviews Genetics*, 17(6):333–351. 864
- Jiaqi Guan, Runzhe Li, Sheng Yu, and Xuegong Zhang. 2018. [Generation of synthetic electronic medical record text](#). 865
- Junxian He, Graham Neubig, and Taylor Berg-Kirkpatrick. 2021. [Efficient nearest neighbor language models](#). In *EMNLP*. 866
- Byungjin Hwang, Ji Hyun Lee, and Duhee Bang. 2018. [Single-cell rna sequencing technologies and bioinformatics pipelines](#). *Experimental & molecular medicine*, 50(8):1–14. 867
- Minoru Kanehisa and Peer Bork. 2003. [Bioinformatics in the post-sequence era](#). *Nature genetics*, 33(3):305–310. 868
- Lei Ke, Wenjie Pei, Ruiyu Li, Xiaoyong Shen, and Yu-Wing Tai. 2019. [Reflective decoding network for image captioning](#). 869
- Urvashi Khandelwal, Angela Fan, Dan Jurafsky, Luke Zettlemoyer, and Mike Lewis. 2021. [Nearest neighbor machine translation](#). *ArXiv*, abs/2010.00710. 870
- Urvashi Khandelwal, Omer Levy, Dan Jurafsky, Luke Zettlemoyer, and Mike Lewis. 2019. [Generalization through memorization: Nearest neighbor language models](#). *arXiv preprint arXiv:1911.00172*. 871
- Michael Kramer, Janusz Dutkowski, Michael Yu, Vineet Bafna, and Trey Ideker. 2014. [Inferring gene ontologies from pairwise similarity data](#). *Bioinformatics*, 30(12):i34–i42. 872
- Remi Lebret, David Grangier, and Michael Auli. 2016. [Neural text generation from structured data with application to the biography domain](#). 873
- Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang. 2020. [Biobert: a pre-trained biomedical language representation model for biomedical text mining](#). *Bioinformatics*, 36(4):1234–1240. 874
- Scott H. Lee. 2018. [Natural language generation for electronic health records](#). 875
- Jacob H. Levine, Erin F. Simonds, Sean C. Bendall, Kara L. Davis, El ad D. Amir, Michelle D. Tadmor, Oren Litvin, Harris G. Fienberg, Astraea Jager, Eli R. Zunder, Rachel Finck, Amanda L. Gedman, Ina Radtke, James R. Downing, Dana Pe’er, and Garry P. Nolan. 2015. [Data-driven phenotypic dissection of AML reveals progenitor-like cells that correlate with prognosis](#). *Cell*, 162(1):184–197. 876

- 900 Linjie Li, Yen-Chun Chen, Yu Cheng, Zhe Gan, 950  
901 Licheng Yu, and Jingjing Liu. 2020. [Hero: 951](#)  
902 [Hierarchical encoder for video+language omni- 952](#)  
903 [representation pre-training.](#) 953  
904 Yehao Li, Yingwei Pan, Jingwen Chen, Ting Yao, and 954  
905 Tao Mei. 2021. [X-modaler: A versatile and high- 955](#)  
906 [performance codebase for cross-modal analytics.](#) 956  
907 Tsung-Yi Lin, Michael Maire, Serge Belongie, 957  
908 Lubomir Bourdev, Ross Girshick, James Hays, 958  
909 Pietro Perona, Deva Ramanan, C. Lawrence Zitnick, 959  
910 and Piotr Dollár. 2015. [Microsoft coco: Common 960](#)  
911 [objects in context.](#) 961  
912 Zequn Liu, Shukai Wang, Yiyang Gu, Ruiyi Zhang, 962  
913 Ming Zhang, and Sheng Wang. 2021. Graphine: A 963  
914 dataset for graph-aware terminology definition gen- 964  
915 eration. *arXiv preprint arXiv:2109.04018.* 965  
916 Kelvin Luu, Rik Koncel-Kedziorski, Kyle Lo, Isabel 966  
917 Cachola, and Noah A. Smith. 2020. Citation text 967  
918 generation. *ArXiv*, abs/2002.00317. 968  
919 Ilaria Manco, Emmanouil Benetos, Elio Quinton, and 969  
920 Gyorgy Fazekas. 2021. [Muscaps: Generating cap- 970](#)  
921 [tions for music audio.](#) 971  
922 Xinhao Mei, Qiushi Huang, Xubo Liu, Gengyun 972  
923 Chen, Jingqian Wu, Yusong Wu, Jinzheng Zhao, 973  
924 Shengchen Li, Tom Ko, H Lilian Tang, Xi Shao, 974  
925 Mark D. Plumbley, and Wenwu Wang. 2021. An 975  
926 encoder-decoder based audio captioning system 976  
927 with transfer and reinforcement learning. 977  
928 Oren Melamud and Chaitanya Shivade. 2019. [Towards 978](#)  
929 [automatic generation of shareable synthetic clinical 979](#)  
930 [notes using neural language models.](#) 980  
931 Yishu Miao and Phil Blunsom. 2016. Language as a 981  
932 latent variable: Discrete generative models for sen- 982  
933 tence compression. In *EMNLP*. 983  
934 Yasuhide Miura, Yuhao Zhang, Emily Tsai, Curtis Lan- 984  
935 glotz, and Dan Jurafsky. 2021. Improving factual 985  
936 completeness and consistency of image-to-text ra- 986  
937 diology report generation. In *Proceedings of the 2021 987*  
938 *Conference of the North American Chapter of the As- 988*  
939 *sociation for Computational Linguistics (NAACL)*. 989  
940 Ivan Montero, Nikolaos Pappas, and Noah A. Smith. 990  
941 2021. Sentence bottleneck autoencoders from trans- 991  
942 former language models. In *EMNLP*. 992  
943 Ratish Puduppully and Mirella Lapata. 2021. [Data-to- 993](#)  
944 [text generation with macro planning.](#) 994  
945 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya 995  
946 Ramesh, Gabriel Goh, Sandhini Agarwal, Girish 996  
947 Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, 997  
948 Gretchen Krueger, and Ilya Sutskever. 2021. [Learn- 998](#)  
949 [ing transferable visual models from natural language 999](#)  
950 [supervision.](#)  
951 Colin Raffel, Noam Shazeer, Adam Roberts, Katherine 952  
952 Lee, Sharan Narang, Michael Matena, Yanqi Zhou, 953  
953 Wei Li, and Peter J. Liu. 2020. Exploring the lim- 954  
954 its of transfer learning with a unified text-to-text 955  
955 transformer. *Journal of Machine Learning Research*, 956  
956 21(140):1–67. 957  
957 Clément Rebuffel, Marco Roberti, Laure Soulier, Geof- 958  
958 frey Scoutheeten, Rossella Cancelliere, and Patrick 959  
959 Gallinari. 2021. [Controlling hallucinations at word 960](#)  
960 [level in data-to-text generation.](#) 961  
961 Clément Rebuffel, Laure Soulier, Geoffrey 962  
962 Scoutheeten, and Patrick Gallinari. 2019. [A 963](#)  
963 [hierarchical model for data-to-text generation.](#) 964  
964 Leonardo F. R. Ribeiro, Yue Zhang, Claire Gardent, 965  
965 and Iryna Gurevych. 2020. [Modeling global and 966](#)  
966 [local node contexts for text generation from knowl- 967](#)  
967 [edge graphs.](#) 968  
968 Linfeng Song, Ante Wang, Jinsong Su, Yue Zhang, 969  
969 Kun Xu, Yubin Ge, and Dong Yu. 2021. [Structural 970](#)  
970 [information preserving for graph-to-text generation.](#) 971  
971 Tianbao Song, Jingbo Sun, Bo Chen, Weiming Peng, 972  
972 and Jihua Song. 2019. Latent space expanded vari- 973  
973 ational autoencoder for sentence generation. *IEEE 974*  
974 *Access*, 7:144618–144627. 975  
975 Neslihan Suzen, Alexander Gorban, Jeremy Levesley, 976  
976 and Evgeny Mirkes. 2021. [Semantic analysis for 977](#)  
977 [automated evaluation of the potential impact of re- 978](#)  
978 [search articles.](#) 979  
979 Jean-Francois Ton, Walter A. Talbott, Shuangfei Zhai, 980  
980 and Joshua M. Susskind. 2021. Regularized train- 981  
981 ing of nearest neighbor language models. *ArXiv*, 982  
982 abs/2109.08249. 983  
983 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob 984  
984 Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz 985  
985 Kaiser, and Illia Polosukhin. 2017. [Attention is all 986](#)  
986 [you need.](#) 987  
987 Bin Wang and C. C. Jay Kuo. 2020. [Sbert-wk: A sen- 988](#)  
988 [tence embedding method by dissecting bert-based 989](#)  
989 [word models.](#) 990  
990 Sheng Wang, Jianzhu Ma, Michael Ku Yu, Fan Zheng, 991  
991 Edward W Huang, Jiawei Han, Jian Peng, and Trey 992  
992 Ideker. 2018. Annotating gene sets by mining large 993  
993 literature collections with protein networks. In *PA- 994*  
994 *CIFIC SYMPOSIUM ON BIOCOMPUTING 2018: 995*  
995 *Proceedings of the Pacific Symposium*, pages 602– 996  
996 613. World Scientific. 997  
997 Wenlin Wang, Chenyang Tao, Zhe Gan, Guoyin Wang, 998  
998 Liqun Chen, Xinyuan Zhang, Ruiyi Zhang, Qian 999  
999 Yang, Ricardo Henao, and Lawrence Carin. 2019. [Im- 999](#)  
1000 [proving textual network learning with variational 999](#)  
1000 [homophilic embeddings.](#)  
1000 Yixin Wang, Zihao Lin, Jiang Tian, Zhongchao Shi, 997  
998 Yang Zhang, Jianping Fan, and Zhiqiang He. 2021. [Confidence- 998](#)  
999 [guided radiology report generation.](#) 999

1000 Ho-Hsiang Wu, Prem Seetharaman, Kundan Kumar, 1050  
 1001 and Juan Pablo Bello. 2021. [Wav2clip: Learning](#) 1051  
 1002 [robust audio representations from clip.](#) 1052

1003 Jianbo Yuan, Haofu Liao, Rui Luo, and Jiebo Luo. 1053  
 1004 2019. [Automatic radiology report generation based](#) 1054  
 1005 [on multi-view image fusion and medical concept en-](#) 1055  
 1006 [richment.](#) 1056

1007 Xinyuan Zhang, Yi Yang, Siyang Yuan, Dinghan Shen, 1057  
 1008 and Lawrence Carin. 2019. [Syntax-infused vari-](#) 1058  
 1009 [ational autoencoder for text generation.](#) *ArXiv*, 1059  
 1010 [abs/1906.02181.](#) 1060

1011 Yanjian Zhang, Qin Chen, Yiteng Zhang, Zhongyu 1061  
 1012 Wei, Yixu Gao, Jiajie Peng, Zengfeng Huang, Wei- 1062  
 1013 jian Sun, and Xuan-Jing Huang. 2020. [Automatic](#) 1063  
 1014 [term name generation for gene ontology: Task and](#) 1064  
 1015 [dataset.](#) In *Proceedings of the 2020 Conference on* 1065  
 1016 *Empirical Methods in Natural Language Processing:* 1066  
 1017 *Findings*, pages 4705–4710. 1067

1018 Luowei Zhou, Yingbo Zhou, Jason J. Corso, Richard 1068  
 1019 Socher, and Caiming Xiong. 2018. [End-to-end](#) 1069  
 1020 [dense video captioning with masked transformer.](#) 1070

## 1020 A Appendices 1070

1021 We provided more details here about our dataset 1071  
 1022 and related experimental results here. In Table S1, 1072  
 1023 we summarized the statistics information of 9 Tex- 1073  
 1024 tomics platforms. There are 3 different 3 species 1074  
 1025 across 9 platforms, including Homo sapiens, Ara- 1075  
 1026 bidopsis thailiana, and Mus musculus. #Sample 1076  
 1027 (All) represents the entire number of samples for 9 1077  
 1028 platforms, #Sample (Vec2Text) represents the num- 1078  
 1029 ber of samples in the subset after BLEU filtering, 1079  
 1030 and #Sample (PMC) represents the number of sam- 1080  
 1031 ples in the subset with full scientific articles. 1081

1032 We also represented the results of Spearman cor- 1082  
 1033 relations between text-based similarity and vector- 1083  
 1034 based similarity across 9 platforms in Table S2. The 1084  
 1035 Spearman correlations are all higher than 0.2 in ev- 1085  
 1036 ery platform, which shows a substantial agreement 1086  
 1037 between text-based similarity and vector-based sim- 1087  
 1038 ilarity. 1088

1039 In Table S3, We represented the automatic eval- 1089  
 1040 uation metric scores for vec2text task, which in- 1090  
 1041 cluded BLEU-1, BLEU-2, ROUGE-1, ROUGE-L, 1091  
 1042 METEOR and NIST, which indicated consistent 1092  
 1043 improvement of our method over comparison ap- 1093  
 1044 proaches on different automatic metrics. 1094  
 1045 1095  
 1046 1096  
 1047 1097  
 1048 1098  
 1049 1099

Platform	Species	#Sample (All)	#Sample (PMC)	# Sample (Vec2Text)	#Feature	M. R.
GPL96	H. S.	1,371	353	240	100K	0.19
GPL198	A. T.	1,081	194	250	100K	0.03
GPL570	H. S.	5,822	1,879	1,004	100K	0.12
GPL1261	M. M.	4,563	1,326	1,059	100K	0.09
GPL6244	H. S.	1,831	659	307	100K	0.10
GPL6246	H. S.	2,366	850	388	100K	0.08
GPL6887	M. M.	1,150	407	240	100K	0.09
GPL10558	H. S.	2,580	1,261	519	100K	0.11
GPL13534	H. S.	1,509	762	234	100K	0.26

Table S1: Statistics of the Textomics data. Each row is a sequencing platform in Textomics. H. S. denotes Homo Sapiens. A. T. denotes Arabidopsis Thaliana. M. M. denotes Mus Musculus. M. R. denotes missing rate. All, PMC, Vec2Text represent number of samples without filtering, with associated PMC full text article, and after using automated filtering, respectively.

Textomics platform	GPL 96	GPL 198	GPL 570	GPL 1261	GPL 6244	GPL 6246	GPL 6887	GPL 10558	GPL 13534
Spearman correlation	0.36	0.20	0.24	0.34	0.44	0.45	0.22	0.38	0.30

Table S2: The result for spearman correlation

Platform	BLEU-1	ROUGE-1	ROUGE-L	METEOR	NIST
GPL96	0.179	0.233	0.166	0.143	0.817
GPL198	0.198	0.257	0.192	0.168	0.889
GPL570	0.212	0.269	0.205	0.182	0.936
GPL1261	0.229	0.283	0.226	0.202	0.980
GPL6244	0.183	0.250	0.179	0.156	0.750
GPL6246	0.219	0.269	0.210	0.187	0.950
GPL6887	0.198	0.260	0.196	0.171	0.847
GPL10558	0.191	0.257	0.177	0.165	0.842
GPL13534	0.242	0.332	0.279	0.260	1.124

Table S3: The first result for evaluating paper embedding using textomics