

---

# Supplementary document for “STimage-1K4M: A histopathology image-gene expression dataset for spatial transcriptomics”

---

Jiawen Chen   Muqing Zhou   Wenrong Wu   Jinwei Zhang   Yun Li\*   Didong Li\*

University of North Carolina at Chapel Hill  
{jiawenn, muqingz, wenrong, jinweizh, yun\_li, didongli}@unc.edu

## 1 S1 DataSheet for STimage-1K4M

2 We present a DataSheet (Gebreu et al., 2021) for STimage-1K4M in this section.

### 3 1. Motivation

- 4 • **For what purpose was the dataset created?** STimage-1K4M was created for training  
5 histopathology multi-modal, self-supervised models, for understanding pathology  
6 images and gene expression data.
- 7 • **Who created the dataset (e.g., which team, research group) and on behalf of which  
8 entity (e.g., company, institution, organization)?** This dataset was curated by Dr.  
9 Yun Li and Dr. Didong Li’s group on behalf of University of North Carolina at Chapel  
10 Hill.
- 11 • **Who funded the creation of the dataset?** This work was supported by R01  
12 AG079291, R01 MH125236, U01 HG011720, P50 HD103573, R56 LM013784,  
13 R01 HL149683, and UM1 TR004406.

### 14 2. Composition

- 15 • **What do the instances that comprise the dataset represent (e.g., documents, photos,  
16 people, countries)?** This dataset includes histopathology images, spatial coordinates  
17 for spots and the paired gene expressions.
- 18 • **How many instances are there in total (of each type, if appropriate)?** STimage-  
19 1K4M includes 1,149 whole-slide images and 4,293,195 spots (sub-tiles) and the  
20 expression of 15,000-30,000 genes associated with each spot.
- 21 • **Does the dataset contain all possible instances or is it a sample (not necessarily  
22 random) of instances from a larger set?** The STimage-1K4M dataset is not a subset  
23 of a larger collection but rather an extensive compilation of all available instances we  
24 could identify and collect. We made a comprehensive effort to collect as much data  
25 as possible, specifically targeting datasets from the Gene Expression Omnibus (GEO)  
26 that contain both pathology images and gene expression data. While it may not cover  
27 every possible instance due to the inherent limitations of data availability and access,  
28 it represents the most exhaustive collection of such data currently available. We are  
29 committed to updating the dataset as more data or technologies become available.

---

\*Corresponding authors

- 30 • **What data does each instance consist of?** We consider each spot as an instance,  
31 which has high dimensional gene expression data, image data and spot coordinate data.
- 32 • **Is there a label or target associated with each instance?** Yes, the gene expression  
33 data could be treated as label for each image.
- 34 • **Is any information missing from individual instances?** We provide extra information  
35 like abstract, paper title. Such information is missed in datasets without a valid  
36 publication id.
- 37 • **Are relationships between individual instances made explicit (e.g., users' movie  
38 ratings, social network links)?** Yes, all instances of the same slide/dataset and their  
39 spatial relationship could be analyzed using the spatial coordinate file.
- 40 • **Are there recommended data splits (e.g., training, development/validation, test-  
41 ing)?** There are no recommended data splits, but potentially the data could be split by  
42 tissue type.
- 43 • **Are there any errors, sources of noise, or redundancies in the dataset?** While the  
44 STImage-1K4M dataset is carefully curated, the source ST data may contain inherent  
45 noise and errors due to the limitations of the technology used. We are not aware of any  
46 redundancies in the dataset.
- 47 • **Is the dataset self-contained, or does it link to or otherwise rely on external  
48 resources (e.g., websites, tweets, other datasets)?** The dataset is self-contained.
- 49 • **Does the dataset contain data that might be considered confidential (e.g., data  
50 that is protected by legal privilege or by doctor-patient confidentiality, data that  
51 includes the content of individuals' non-public communications)?** All the data in  
52 STImage-1K4M is from public available source, we are not aware of such confidential  
53 information.
- 54 • **Does the dataset contain data that, if viewed directly, might be offensive, insulting,  
55 threatening, or might otherwise cause anxiety?** No.
- 56 • **Does the dataset identify any subpopulations (e.g., by age, gender)?** Not explicitly.
- 57 • **Is it possible to identify individuals (i.e., one or more natural persons), either  
58 directly or indirectly (i.e., in combination with other data) from the dataset?** No.
- 59 • **Does the dataset contain data that might be considered sensitive in any way (e.g.,  
60 data that reveals race or ethnic origins, sexual orientations, religious beliefs,  
61 political opinions or union memberships, or locations; financial or health data;  
62 biometric or genetic data; forms of government identification, such as social  
63 security numbers; criminal history)?** No.

### 64 3. Collection Process

- 65 • **How was the data associated with each instance acquired?** We queried the GEO  
66 website using keywords "spatial transcriptomics", specifically targeting supplementary  
67 files including files in JPG, PNG, or TIFF formats. This search resulted in 856 datasets  
68 from 121 unique GEO studies. Additionally, we gathered 58 Visium and 4 VisiumHD  
69 datasets from 10X Genomics, complementing these with 233 slides manually collected  
70 from 10 additional studies.
- 71 • **What mechanisms or procedures were used to collect the data (e.g., hardware  
72 apparatuses or sensors, manual human curation, software programs, software  
73 APIs)?** We used rvest R package to gather download links for datasets in GEO. For  
74 other datasets, data were collected by human manual curation.
- 75 • **If the dataset is a sample from a larger set, what was the sampling strategy (e.g.,  
76 deterministic, probabilistic with specific sampling probabilities)?** Not applicable,  
77 this dataset is not a sample from a larger set.
- 78 • **Who was involved in the data collection process (e.g., students, crowdworkers,  
79 contractors) and how were they compensated (e.g., how much were crowdworkers  
80 paid)?** Jiawen Chen, Wenrong Wu and Jinwei Zhang are involved in the data collection  
81 process. All of them are graduate students.

- 82
- **Over what timeframe was the data collected?** STimage-1K4M includes data generated from 2016-2024.
  - 83
  - 84 • **Were any ethical review processes conducted (e.g., by an institutional review board)?** No official ethical review processes were conducted.
  - 85
  - 86 • **Did you collect the data from the individuals in question directly, or obtain it via third parties or other sources (e.g., websites)?** The data were collected from
  - 87 websites.
  - 88
  - 89 • **Were the individuals in question notified about the data collection?** Not applicable,
  - 90 we collected the data from publicly available sources with no contact with individuals
  - 91 involved in the study.
  - 92 • **Did the individuals in question consent to the collection and use of their data?** Not
  - 93 applicable, we collected the data from publicly available sources without contact with
  - 94 individuals involved in the study. We cite all the studies included in the dataset.
  - 95 • **If consent was obtained, were the consenting individuals provided with a mechanism to revoke their consent in the future or for certain uses?** Not applicable, we
  - 96 collected the data from publicly available sources.
  - 97
  - 98 • **Has an analysis of the potential impact of the dataset and its use on data subjects (e.g., a data protection impact analysis) been conducted?** Not applicable, we
  - 99 collected the data from publicly available sources.
  - 100

#### 101 4. Preprocessing/cleaning/labeling

- 102 • **Was any preprocessing/cleaning/labeling of the data done (e.g., discretization or bucketing, tokenization, part-of-speech tagging, SIFT feature extraction, removal of instances, processing of missing values)?** Yes, the meta data like tissue type
- 103 were cleaned manually. All code related to label cleaning is available in the GitHub
- 104 repository <https://github.com/JiawenChenn/STimage-1K4M>.
- 105
- 106 • **Was the “raw” data saved in addition to the preprocessed/cleaned/labeled data (e.g., to support unanticipated future uses)?** Yes, all raw labels were saved.
- 107
- 108 • **Is the software that was used to preprocess/clean/label the data available?** Yes, all
- 109 code related to label cleaning is available in the same GitHub repository.
- 110

#### 111 5. Uses

- 112 • **Has the dataset been used for any tasks already?** No.
- 113 • **Is there a repository that links to any or all papers or systems that use the dataset?** Yes, all sources are available at <https://github.com/JiawenChenn/STimage-1K4M>.
- 114
- 115
- 116 • **What (other) tasks could the dataset be used for?** The dataset could be used for
- 117 training self-supervised models to better understand histopathology and gene expres-
- 118 sion.
- 119 • **Is there anything about the composition of the dataset or the way it was collected and preprocessed/cleaned/labeled that might impact future uses?** Yes, we would
- 120 recommend our format as the future data format release for ST data.
- 121
- 122 • **Are there tasks for which the dataset should not be used?** We are not aware of such
- 123 task.

#### 124 6. Distribution

- 125 • **Will the dataset be distributed to third parties outside of the entity (e.g., company, institution, organization) on behalf of which the dataset was created?** Yes, all the
- 126 data are distributed under a permissible license for research-based use.
- 127
- 128 • **How will the dataset will be distributed (e.g., tarball on website, API, GitHub)?**
- 129 The dataset is distributed on Huggingface: [https://huggingface.co/datasets/](https://huggingface.co/datasets/jiawennnn/STimage-1K4M)
- 130 [jiawennnn/STimage-1K4M](https://huggingface.co/datasets/jiawennnn/STimage-1K4M).
- 131 • **When will the dataset be distributed?** The dataset is released with a permissible
- 132 license for research-based use.

- 133 • **Will the dataset be distributed under a copyright or other intellectual property**  
134 **(IP) license, and/or under applicable terms of use(ToU)?** The use of data will be  
135 under a permissible license for research-based use.
- 136 • **Have any third parties imposed IP-based or other restrictions on the data associ-**  
137 **ated with the instances?** We are not aware of such restrictions.
- 138 • **Do any export controls or other regulatory restrictions apply to the dataset or to**  
139 **individual instances?** We are not aware of such restrictions.

## 140 7. Maintenance

- 141 • **Who will be supporting/hosting/maintaining the dataset?** The first author of the  
142 paper.
- 143 • **How can the owner/curator/manager of the dataset be contacted (e.g., email**  
144 **address)?** The first author and corresponding authors could be contacted using email  
145 listed in the paper or through GitHub.
- 146 • **Is there an erratum?** No.
- 147 • **Will the dataset be updated (e.g., to correct labeling errors, add new instances,**  
148 **delete instances)?** Yes, the dataset will be updated periodically to ensure data quality.  
149 We are also committed to continually expanding the dataset by adding new samples as  
150 they become available.
- 151 • **If the dataset relates to people, are there applicable limits on the retention of the**  
152 **data associated with the instances (e.g., were the individuals in question told that**  
153 **their data would be retained for a fixed period of time and then deleted)?** We are  
154 not aware of such limits.
- 155 • **Will older versions of the dataset continue to be supported/hosted/maintained?**  
156 All versions of the dataset will be available.
- 157 • **If others want to extend/augment/build on/contribute to the dataset, is there a**  
158 **mechanism for them to do so?** Not at this time.

## 159 S2 Additional dataset information

160 STimage-1K4M is released at <https://github.com/JiawenChenn/STimage-1K4M> with meta-  
161 data record also contained in this repository. The license for the data use is a permissible license  
162 for research-based use, which is described detailly on Huggingface: [https://huggingface.co/](https://huggingface.co/datasets/jiawennnn/STimage-1K4M)  
163 [datasets/jiawennnn/STimage-1K4M](https://huggingface.co/datasets/jiawennnn/STimage-1K4M). All code related to this project is under MIT license.

## 164 S3 Author Statement

165 The authors of the STimage-1K4M dataset bear full responsibility for the content and compliance of  
166 this project. All authors of this paper have confirmed the data license. The data is now hosted on ftp  
167 server and will be maintained by the first author of this paper.

## 168 References

169 Gebru, T., J. Morgenstern, B. Vecchione, J. W. Vaughan, H. Wallach, H. D. Iii, and K. Crawford  
170 (2021). Datasheets for datasets. *Communications of the ACM* 64(12), 86–92.