
Learning More Effective Cell Representations Efficiently

Jason Xiaotian Dou
University of Pittsburgh
jasondpku@gmail.com

Minxue Jia
University of Pittsburgh
minxue.jia@pitt.edu

Nika Zaslavsky
Carnegie Mellon University
nzaslavsky@cmu.edu

Mark Ebeid
University of Pittsburgh
mae117@pitt.edu

Runxue Bao
University of Pittsburgh
runxue.bao@pitt.edu

Shiyi Zhang
Carnegie Mellon University
shiyiz@andrew.cmu.edu

Ke Ni
University of Pittsburgh
ken67@pitt.edu

Paul Pu Liang
Carnegie Mellon University
paul.liangpu@gmail.com

Haiyi Mao
University of Pittsburgh
ham112@pitt.edu

Zhi-hong Mao
University of Pittsburgh
zhm4@pitt.edu

Abstract

1 Capturing similarity among cells is at the core of many tasks in single-cell tran-
2 scriptomics, such as the identification of cell types and cell states. This problem
3 can be formulated in a paradigm called metric learning. Metric learning aims to
4 learn data embeddings (feature vectors) in a way that reduces the distance between
5 similar feature vectors corresponding to cells of the same cell type, and increases
6 the distance between feature vectors corresponding to cells of different cell types.
7 As a variation of metric learning, deep metric learning uses neural networks to
8 automatically learn discriminative features from the cells and then compute the
9 distance. These (deep) metric learning approaches have been successfully applied
10 to computational biology tasks like similar cell identification, and synthesis of het-
11 erogeneous single-cell modalities. Here, we identify two computational challenges:
12 precise distance measurement between cells, and scalability over a large amount of
13 data in the applications of (deep) metric learning. We then propose our solutions:
14 optimal transport and coresset optimization. Optimal transport has the potential to
15 measure cell similarity more effectively, and coresset optimization is promising to
16 train representation learning models more efficiently. Empirical studies in image
17 retrieval and clustering tasks show the promise of the proposed approaches. We
18 propose to further explore the applicability of our methods to cell representation
19 learning.

20 1 Introduction

21 The success of machine learning algorithms largely depends on data representation. Metric learning
22 learns data embeddings and feature vectors in a way that reduces the distance between feature vectors
23 corresponding to objects belonging to the same class and increases the distance between the feature
24 vectors corresponding to different classes. Deep metric learning, on the other hand, uses neural

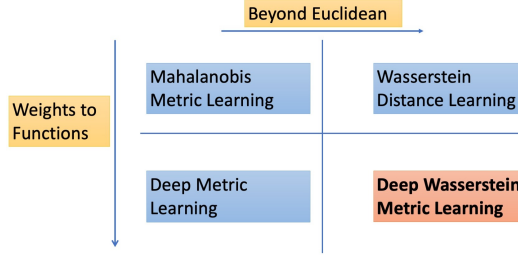


Figure 1: Overview of the Deep Wasserstein Metric Learning Framework

25 networks to automatically learn discriminative features from the objects and then compute the metric.
 26 (Deep) metric learning falls into the broad umbrella of representation learning, whose quest for
 27 representation learning is motivating the design of more powerful representation [3, 13].

28 These representation learning methods have achieved great successes in biology applications [40,
 29 43, 19, 23, 35, 56, 23, 50, 1, 11, 30, 14, 46, 47]. For example, Schema [40] uses a principled metric
 30 learning strategy to identify informative features in a modality to synthesize disparate modalities
 31 into a single coherent interpretation. It is used to infer cell types by integrating gene expression
 32 and chromatin accessibility data. Specifically, [53] presents an approach for integrating different
 33 modalities by learning a probabilistic coupling among them using autoencoders to map to a shared
 34 latent space. These methods can complement start of art single cell analytics tool such as Dynamo
 35 [38] to gain new insights into dynamic biological processes.

36 The deep metric learning framework SCimilarity (Invited Talk “Design for Inference in Drug Discov-
 37 ery and Development” by Aviv Regev, ICML 2022), which employs a standard triplet loss design, has
 38 achieved impressive performance in identifying similar cells in a massive collection of scRNA-Seq
 39 datasets. The results can help answer questions like in which tissues and diseases we find fibrotic
 40 macrophage-like cells.

41 [15] proposes the Deep Wasserstein Metric Learning Framework, as shown in Figure 1, which
 42 conducts multiple steps of adjustments over the original metric learning framework and achieves
 43 improved performance in image retrieval and clustering tasks. In the rest of the proposal, we introduce
 44 the adjustments accordingly and propose empirical studies for single-cell applications.

45 2 Methods

46 2.1 Representation Learning: Metric Learning and Deep Metric Learning

47 Representation learning is a class of machine learning approaches that allow a system to discover the
 48 representations required for feature detection or classification from raw data [38, 53, 18, 24, 17, 10,
 49 26, 12, 55, 52, 42] The requirement for manual feature engineering is reduced by allowing a machine
 50 to learn the features and apply them to a given activity. Metric learning and deep metric learning,
 51 specifically, focus on similarity-based approaches to learning the representations. Thus the similarity
 52 measurement becomes very important. Previous work [7] studies similarity measurement in gene
 53 expression. Metric learning has only limited capability to capture non-linearity in the data, while
 54 deep metric learning captures non-linear features better by learning the non-linear transformation.
 55 The most widely used loss functions for deep metric learning are the contrastive loss and the triplet
 56 loss, both use euclidean distance to measure the distance between objects. A more comprehensive
 57 illustration of so-called “ranking-based” loss functions are summarized in Figure 2. Given an image
 58 pair, the contrastive loss minimizes their distance in the embedding space if their classes are the same,
 59 and separates them a fixed margin away otherwise. The triplet loss takes triplets of anchor, positive,
 60 and negative images, and enforces the distance between the anchor and the positive to be smaller than
 61 that between the anchor and the negative. The formation of contrastive loss is as the following. We
 62 first have embedding pairs \mathcal{P} , which is sampled from a minibatch of size b . The pair contains an
 63 anchor ϕ_a from class y_a and either a positive ϕ_p with $y_a = y_p$ or a negative ϕ_n from a different class,
 64 $y_a \neq y_n$. The distance function we utilize is the standard Euclidean distance $d_e(x, y) = \|x - y\|_2$.

65 Then the network ϕ is trained to minimize:

$$\mathcal{L}_{\text{contrastive}} = \frac{1}{b} \sum_{(i,j) \in \mathcal{P}} \mathbb{I}_{y_i=y_j} d_e(\phi_i, \phi_j) + \mathbb{I}_{y_i \neq y_j} [\gamma - d_e(\phi_i, \phi_j)]_+ \quad (1)$$

66 Triplets extend the contrastive formulation by providing a triplets \mathcal{T} sampled from a mini-batch:

$$\mathcal{L}_{\text{triplet}} = \frac{1}{b} \sum_{\substack{(a,p,n) \in \mathcal{T} \\ y_a=y_p \neq y_n}} [d_e(\phi_a, \phi_p) - d_e(\phi_a, \phi_n) + \gamma]_+ \quad (2)$$

67 In the following, we present two adjustments that can improve (deep) metric learning’s performances,
68 one is optimal transport [9], and one is coresets optimization [31].

69 2.2 Optimal Transport

70 Optimal transport (OT) is the general problem of moving one distribution of mass to another as
71 efficiently as possible [41, 15, 34, 49]. Optimal transport has been used tremendously in computational
72 biology [39, 37, 44, 45]. For example, [39] uses scRNA-seq data collected across a time course to
73 infer how these probability distributions evolve over time, by using the mathematical approach of
74 optimal transport.

75 Wasserstein distance provides the mathematical tool to measure distances between functions, his-
76 tograms, or more general objects in the optimal transport problem. Wasserstein distance is also called
77 Earth Mover’s Distance, which is employed to develop PhEMD (Phenotypic Earth Mover’s Distance)
78 [6], which is used to embed the space of drug perturbations on the basis of the drugs’ effects on cell
79 populations. Wasserstein distance-based loss functions have shown superior performance in learning
80 tasks [20]. Thus a new set of loss functions are proposed to replace the Euclidean distance with
81 Wasserstein distance in original contrastive loss and triplet loss by defining

$$d_w(x, y) = W_1(x, y) \quad (3)$$

82 The new Wasserstein-contrastive (wcontrastive) loss and Wasserstein-triplet (wtriplet) loss can be
83 formulated as [15]:

$$\mathcal{L}_{\text{wcontrastive}} = \frac{1}{b} \sum_{(i,j) \in \mathcal{P}} \mathbb{I}_{y_i=y_j} d_w(\phi_i, \phi_j) + \mathbb{I}_{y_i \neq y_j} [\gamma - d_w(\phi_i, \phi_j)]_+, \quad (4)$$

84

$$\mathcal{L}_{\text{wtriplet}} = \frac{1}{b} \sum_{\substack{(a,p,n) \in \mathcal{T} \\ y_a=y_p \neq y_n}} [d_w(\phi_a, \phi_p) - d_w(\phi_a, \phi_n) + \gamma]_+. \quad (5)$$

85 We propose to apply the two new loss functions to similar cell identification and synthesis of
86 heterogeneous modalities applications. In the similar cell identification task, the loss functions
87 above can impose a discriminative constraint on the feature embedding to improve the similarity
88 measurement [51].

89 2.3 Coreset Optimization

90 Coreset optimization is about data-efficient methods to find subsets of massive data that can generalize
91 to the full data when trained on. In other words, a coreset is a subset of the original training set that
92 is representative to train machine learning models [15, 54, 32, 28]. More specifically, Wasserstein
93 measure coreset [8], is an extension of coresets that takes into account continuous data distribution
94 and generalization. Recently coreset has been successfully applied to the purification of single-cell
95 transcriptomics data [33]. It focuses on alleviating potential replicate-specific biases within single-cell
96 datasets. The key is to select a “representative” subset (coresets) of cells from areas of the single-cell
97 landscape where multiple replicates are represented. The approach [33] takes is solving the exemplar
98 clustering problem, which minimizes the sum of pairwise dissimilarities between cells in the coreset
99 and the rest of cells. We follow the idea in [33] but use Wasserstein distance to replace the Gaussian

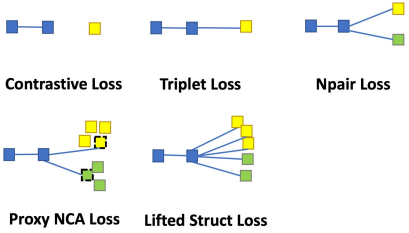


Figure 2: Illustration of different ranking-based loss functions. Different colors (blue, yellow, green) represent different classes. For simplicity, only 3 classes are shown. The left-most blue square is an anchor (query). In Contrastive loss, the anchor is compared with only one positive example. In Triplet loss, the anchor is compared with only one negative example and one positive example. In Npair, ProxyNCA, and Lifted Struct losses, one positive example and multiple negative classes are incorporated. Npair loss randomly selects one example per negative class. ProxyNCA loss pushes the anchor away from negative proxies instead of negative examples. Lifted Struct loss uses all examples from all negative classes [16, 27].

100 kernel to define similarity between cells. Here $r|V|$ is the exemplar cells from the groundset V and S
 101 is the targeted coreset.

$$S^* \in \arg \max_{|S| \leq r|V|} S \subseteq V \sum_{x \in V} \max_{y \in S} d_w(x, y) \quad (6)$$

102 [40] processes data from a Slide-seq replicate (three modalities with 20823 transcriptomes * 17607
 103 genes) in 34 mins. [31] demonstrates a specific coreset optimization algorithm CRAIG can achieve
 104 the average speedup of 3x for similar loss residual and error rate. So we expect for the single cell
 105 synthesis task in [40], we can reduce data processing time from 34 mins to 11 mins. Furthermore,
 106 feature-efficient methods [1, 2, 48] can be applied to remove irrelevant variables during the training
 107 of coreset optimization algorithms.

108 3 Experiments Design

109 The experiments conducted on various datasets have demonstrated that optimal transport and coreset
 110 optimization can achieve superior performance on image retrieval and clustering tasks [15, 16]. In
 111 the following we lay out empirical studies to explore the applicability in building cell representations.
 112 The specific detail of computational studies is under investigation.

113 3.1 Cell Similarity Identification

114 For the cell similarity identification task, we plan to build on the datasets PBMC [25] and SLN-all
 115 [22] which are included in the phenomenal scvi-tools [21]. The PBMC dataset is measured with
 116 CITE-seq. The SLN-all dataset contains Immune cells from the murine spleen and lymph nodes.

117 3.2 Multimodal Integration

118 Regarding the multimodal integration challenge, we plan to follow the setup in [29] to apply the
 119 methods on multiple single-cell datasets including sci-CAR cell line [4], SNARE-seq cell line [5],
 120 and 10X Multiome T-cell depleted bone marrow [36] to validate the methods' effectiveness and
 121 develop new computational and biological insights from the downstream tasks.

122 4 Discussion

123 In this essay, we propose to apply two computational methods: optimal transport and coreset
 124 optimization, which are successfully demonstrated usability in image representation learning [15, 16],
 125 to cell representation learning with applications in cell similarity identification and multimodal
 126 integration. We will report results and insights in empirical studies in follow-up research.

References

- 127
- 128 [1] R. Bao, B. Gu, and H. Huang. Fast oscar and owl regression via safe screening rules. In
129 *International Conference on Machine Learning*, pages 653–663. PMLR, 2020.
- 130 [2] R. Bao, B. Gu, and H. Huang. An accelerated doubly stochastic gradient method with faster
131 explicit model identification. In *Proceedings of the 31st ACM International Conference on*
132 *Information & Knowledge Management*, pages 57–66, 2022.
- 133 [3] Y. Bengio, A. Courville, and P. Vincent. Representation learning: A review and new perspectives,
134 2012.
- 135 [4] J. Cao, D. A. Cusanovich, V. Ramani, D. Aghamirzaie, H. A. Pliner, A. J. Hill, R. M. Daza,
136 J. L. McFaline-Figueroa, J. S. Packer, L. Christiansen, F. J. Steemers, A. C. Adey, C. Trapnell,
137 and J. Shendure. Joint profiling of chromatin accessibility and gene expression in thousands of
138 single cells. *Science*, 361(6409):1380–1385, 2018.
- 139 [5] S. Chen, B. B. Lake, and K. Zhang. High-throughput sequencing of the transcriptome and
140 chromatin accessibility in the same cell. *Nature Biotechnology*, 37(12):1452–1457.
- 141 [6] W. S. Chen, N. Zivanovic, D. van Dijk, G. Wolf, B. Bodenmiller, and S. Krishnaswamy.
142 Embedding the single-cell experimental variable state space to reveal manifold structure of drug
143 perturbation effects in breast cancer. *bioRxiv*, page 455436, 2018.
- 144 [7] M. Chikina. *Devising effective similarity measures and learning algorithms for the study of*
145 *metazoan gene expression*. Princeton University, 2011.
- 146 [8] S. Clatici, A. Genevay, and J. Solomon. Wasserstein measure coresets, 2020.
- 147 [9] M. Cuturi. Sinkhorn distances: Lightspeed computation of optimal transport. In *Advances in*
148 *neural information processing systems*, pages 2292–2300, 2013.
- 149 [10] X. Ding, L. Zhao, and L. Akoglu. Hyperparameter sensitivity in deep outlier detection: Analysis
150 and a scalable hyper-ensemble solution. In A. H. Oh, A. Agarwal, D. Belgrave, and K. Cho,
151 editors, *Advances in Neural Information Processing Systems*, 2022.
- 152 [11] J. X. Dou. Impartial redistricting: a markov chain approach to the "gerrymandering problem".
153 *arXiv preprint arXiv:1711.04618*, 2017.
- 154 [12] J. X. Dou and R. Bao. Clinical decision system using machine learning and deep learning: a
155 survey.
- 156 [13] J. X. Dou, M. Jia, R. Bao, and H. H. Mao. Enhance ‘similar’ cell identification through optimal
157 transport.
- 158 [14] J. X. Dou, M. Liu, H. Muneer, and A. Schluskel. What words do we use to lie?: Word choice in
159 deceptive messages. *ArXiv*, abs/1710.00273, 2017.
- 160 [15] J. X. Dou, L. Luo, and R. M. Yang. An optimal transport approach to deep metric learning
161 (student abstract). In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36,
162 pages 12935–12936, 2022.
- 163 [16] J. X. Dou, A. Q. Pan, R. Bao, H. H. Mao, L. Luo, and Z. Mao. Sampling through the lens of
164 sequential decision making. *arXiv preprint arXiv:2208.08056*, 2022.
- 165 [17] J. X. Dou, N. Sun, and X. Zou. “draw my topics”: Find desired topics fast from large scale of
166 corpus. *ArXiv*, abs/1602.01428, 2016.
- 167 [18] Y. Ektefaie, G. Dasoulas, A. Noori, M. Farhat, and M. Zitnik. Geometric multimodal represen-
168 tation learning, 2022.
- 169 [19] M. Flores, Z. Liu, T. Zhang, M. M. Hasib, Y.-C. Chiu, Z. Ye, K. Paniagua, S. Jo, J. Zhang, S.-J.
170 Gao, et al. Deep learning tackles single-cell analysis—a survey of deep learning for scrna-seq
171 analysis. *Briefings in Bioinformatics*, 23(1):bbab531, 2022.

- 172 [20] C. Frogner, C. Zhang, H. Mobahi, M. Araya, and T. A. Poggio. Learning with a wasserstein
173 loss. *Advances in neural information processing systems*, 28, 2015.
- 174 [21] A. Gayoso, R. Lopez, G. Xing, P. Boyeau, V. Valiollah Pour Amiri, J. Hong, K. Wu, M. Jaya-
175 suriya, E. Mehlman, M. Langevin, Y. Liu, J. Samaran, G. Misrachi, A. Nazaret, O. Clivio, C. Xu,
176 T. Ashuach, M. Gabitto, M. Lotfollahi, V. Svensson, E. da Veiga Beltrame, V. Kleshchevnikov,
177 C. Talavera-López, L. Pachter, F. J. Theis, A. Streets, M. I. Jordan, J. Regier, and N. Yosef. A
178 python library for probabilistic analysis of single-cell omics data. *Nature Biotechnology*, Feb
179 2022.
- 180 [22] A. Gayoso, Z. Steier, R. Lopez, J. Regier, K. L. Nazor, A. Streets, and N. Yosef. Joint
181 probabilistic modeling of single-cell multi-omic data with totalvi. *Nature methods*, 18(3):272–
182 282, 2021.
- 183 [23] S. Grebinoski, Q. Zhang, A. R. Cillo, S. Manne, H. Xiao, E. A. Brunazzi, T. Tabib, C. Cardello,
184 C. G. Lian, G. F. Murphy, et al. Autoreactive cd8+ t cells are restrained by an exhaustion-like
185 program that is maintained by lag3. *Nature Immunology*, 23(6):868–877, 2022.
- 186 [24] W. L. Hamilton, R. Ying, and J. Leskovec. Representation learning on graphs: Methods and
187 applications. *arXiv preprint arXiv:1709.05584*, 2017.
- 188 [25] Y. Hao, S. Hao, E. Andersen-Nissen, W. M. Mauck III, S. Zheng, A. Butler, M. J. Lee, A. J.
189 Wilk, C. Darby, M. Zager, et al. Integrated analysis of multimodal single-cell data. *Cell*,
190 184(13):3573–3587, 2021.
- 191 [26] S. He, Y. Wang, S. Han, S. Zou, and F. Miao. A robust and constrained multi-agent reinforcement
192 learning framework for electric vehicle amod systems. *arXiv preprint arXiv:2209.08230*, 2022.
- 193 [27] S. Kim, D. Kim, M. Cho, and S. Kwak. Proxy anchor loss for deep metric learning. In
194 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages
195 3238–3247, 2020.
- 196 [28] T. Lin, Z. Zheng, E. Y. Chen, M. Cuturi, and M. I. Jordan. On projection robust optimal
197 transport: Sample complexity and model misspecification. *arXiv preprint arXiv:2006.12301*,
198 2020.
- 199 [29] H. Mao, M. Jia, J. X. Dou, H. Zhang, and P. V. Benos. Coem: Cross-modal embedding for
200 metacell identification. *arXiv preprint arXiv:2207.07734*, 2022.
- 201 [30] H. Mao, H. Liu, J. X. Dou, and P. V. Benos. Towards cross-modal causal structure and
202 representation learning. In *Machine Learning for Health*, 2022.
- 203 [31] B. Mirzasoleiman, J. Bilmes, and J. Leskovec. Coresets for data-efficient training of machine
204 learning models. In *International Conference on Machine Learning (ICML)*, July 2020.
- 205 [32] B. Mirzasoleiman, K. Cao, and J. Leskovec. Coresets for robust training of deep neural networks
206 against noisy labels. *Advances in Neural Information Processing Systems*, 33, 2020.
- 207 [33] R. Pálovics, T. Wyss-Coray, and B. Mirzasoleiman. Purification of single-cell transcriptomics
208 data with coreset selection.
- 209 [34] N. Papadakis. *Optimal transport for image processing*. PhD thesis, 2015.
- 210 [35] S. Paudel, B. E. Warner, R. Wang, J. Adams-Haduch, A. S. Reznik, J. X. Dou, Y. Huang, Y.-T.
211 Gao, W.-P. Koh, A. Bäckerholm, J.-M. Yuan, and K. H. Y. Shair. Serological profiling using an
212 Epstein-Barr virus mammalian expression library identifies EBNA1 IgA as a pre-diagnostic
213 marker for nasopharyngeal carcinoma. *Clinical Cancer Research*, 09 2022. CCR-22-1600.
- 214 [36] S. Persad, Z.-N. Choo, C. Dien, I. Masilionis, R. Chaligné, T. Nawy, C. C. Brown, I. Pe’er,
215 M. Setty, and D. Pe’er. Seacells: Inference of transcriptional and epigenomic cellular states
216 from single-cell genomics data. *bioRxiv*, 2022.
- 217 [37] N. Prasad, K. Yang, and C. Uhler. Optimal transport using gans for lineage tracing, 2020.

- 218 [38] X. Qiu, Y. Zhang, J. D. Martin-Rufino, C. Weng, S. Hosseinzadeh, D. Yang, A. N. Pogson, M. Y.
219 Hein, K. H. J. Min, L. Wang, et al. Mapping transcriptomic vector fields of single cells. *Cell*,
220 185(4):690–711, 2022.
- 221 [39] G. Schiebinger, J. Shu, M. Tabaka, B. Cleary, V. Subramanian, A. Solomon, J. Gould, S. Liu,
222 S. Lin, P. Berube, et al. Optimal-transport analysis of single-cell gene expression identifies
223 developmental trajectories in reprogramming. *Cell*, 176(4):928–943, 2019.
- 224 [40] R. Singh, B. L. Hie, A. Narayan, and B. Berger. Schema: metric learning enables interpretable
225 synthesis of heterogeneous single-cell modalities. *Genome biology*, 22(1):1–24, 2021.
- 226 [41] R. Singh, J. S. S. Li, S. G. Tattikota, Y. Liu, J. Xu, Y. Hu, N. Perrimon, and B. Berger. Optimal
227 transport analysis of single-cell transcriptomics directs hypotheses prioritization and validation.
228 *bioRxiv*, 2022.
- 229 [42] J. X. Tan and T. Finkel. A phosphoinositide signalling pathway mediates rapid lysosomal repair.
230 *Nature*, 609(7928):815–821, 2022.
- 231 [43] A. Tanay and A. Regev. Scaling single-cell genomics from phenomenology to mechanism.
232 *Nature*, 541(7637):331–338, 2017.
- 233 [44] A. Tong, J. Huang, G. Wolf, D. Van Dijk, and S. Krishnaswamy. Trajectorynet: A dynamic
234 optimal transport network for modeling cellular dynamics. In *International conference on*
235 *machine learning*, pages 9526–9536. PMLR, 2020.
- 236 [45] A. Y. Tong, G. Huguet, A. Natic, K. MacDonald, M. Kuchroo, R. Coifman, G. Wolf, and S. Kr-
237 ishnaswamy. Diffusion earth mover’s distance and distribution embeddings. In *International*
238 *Conference on Machine Learning*, pages 10336–10346. PMLR, 2021.
- 239 [46] L. C. Wu, J. X. Dou, D. Sleator, A. M. Frieze, and D. Miller. Impartial redistricting: A markov
240 chain approach. *ArXiv*, abs/1510.03247, 2015.
- 241 [47] L. Xiong, K. Tian, Y. Li, W. Ning, X. Gao, and Q. C. Zhang. Online single-cell data integration
242 through projecting heterogeneous datasets into a common cell-embedding space. *Nature*
243 *communications*, 13(1):1–17, 2022.
- 244 [48] A. Xu and H. Huang. Detached error feedback for distributed SGD with random sparsification.
245 In K. Chaudhuri, S. Jegelka, L. Song, C. Szepesvari, G. Niu, and S. Sabato, editors, *Proceedings*
246 *of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of*
247 *Machine Learning Research*, pages 24550–24575. PMLR, 17–23 Jul 2022.
- 248 [49] H. Xu, J. Liu, D. Luo, and L. Carin. Representing graphs via gromov-wasserstein factorization.
249 *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022.
- 250 [50] Q. Xu, Y. Yang, X. Zhang, and J. J. Cai. Association of pyroptosis and severeness of covid-19
251 as revealed by integrated single-cell transcriptome data analysis. *ImmunoInformatics*, 6:100013,
252 2022.
- 253 [51] J. Yan, E. Yang, C. Deng, and H. Huang. Metricformer: A unified perspective of correlation
254 exploring in similarity learning. In A. H. Oh, A. Agarwal, D. Belgrave, and K. Cho, editors,
255 *Advances in Neural Information Processing Systems*, 2022.
- 256 [52] H. Yang and J. X. Tan. The pitt pathway: Keeping lysosomes young. *Clinical and Translational*
257 *Medicine*, 12(10), 2022.
- 258 [53] K. D. Yang, A. Belyaeva, S. Venkatachalapathy, K. Damodaran, A. Katcoff, A. Radhakrishnan,
259 G. Shivashankar, and C. Uhler. Multi-domain translation between single-cell imaging and
260 sequencing data using autoencoders. *Nature communications*, 12(1):1–10, 2021.
- 261 [54] Y. Yang, T. Y. Liu, and B. Mirzasoleiman. Not all poisons are created equal: Robust training
262 against data poisoning. In *Proceedings of the 39th International Conference on Machine*
263 *Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 25154–25165.
264 PMLR, 17–23 Jul 2022.

- 265 [55] T.-H. Zhang, M. M. Hasib, Y.-C. Chiu, Z.-F. Han, Y.-F. Jin, M. Flores, Y. Chen, and Y. Huang.
266 Transformer for gene expression modeling (t-gem): An interpretable deep learning model for
267 gene expression-based phenotype predictions. *Cancers*, 14(19):4763, 2022.
- 268 [56] H. Zhong, S. Liu, F. Cao, Y. Zhao, J. Zhou, F. Tang, Z. Peng, Y. Li, S. Xu, C. Wang, et al.
269 Dissecting tumor antigens and immune subtypes of glioma to develop mrna vaccine. *Frontiers*
270 *in immunology*, 12, 2021.