

Multi-modal contrastive learning adapts to *intrinsic dimension* of **shared** latent variables

Yu Gui¹ Cong Ma² Zongming Ma³

¹ Department of Statistics and Data Science, University of Pennsylvania, ² Department of Statistics, University of Chicago, ³ Department of Statistics and Data Science, Yale University
yugui@wharton.upenn.edu, congma@uchicago.edu, zongming.ma@yale.edu

Growing availability of multi-modal data



How can one efficiently integrate data from multi-modalities?

Multi-modal Contrastive Learning

maximize $\text{sim}(f(X_i), g(Y_i))$, minimize $\text{sim}(f(X_i), g(Y_j))$, $i \neq j$

- Contrastive language-image pre-training (CLIP)[2] has been the SOTA pipeline for multi-modal learning
- infoNCE loss $\mathcal{L}^N(f, g, \tau)$ with **temperature optimization**

$$\mathcal{L}^N(f, g, \tau) = -\frac{1}{N} \sum_{i \in [N]} \log \frac{\exp\left(\frac{\sigma(f(X_i), g(Y_i))}{\tau}\right)}{N^{-1} \sum_{j \in [N]} \exp\left(\frac{\sigma(f(X_i), g(Y_j))}{\tau}\right)} + \text{symmetric term}$$

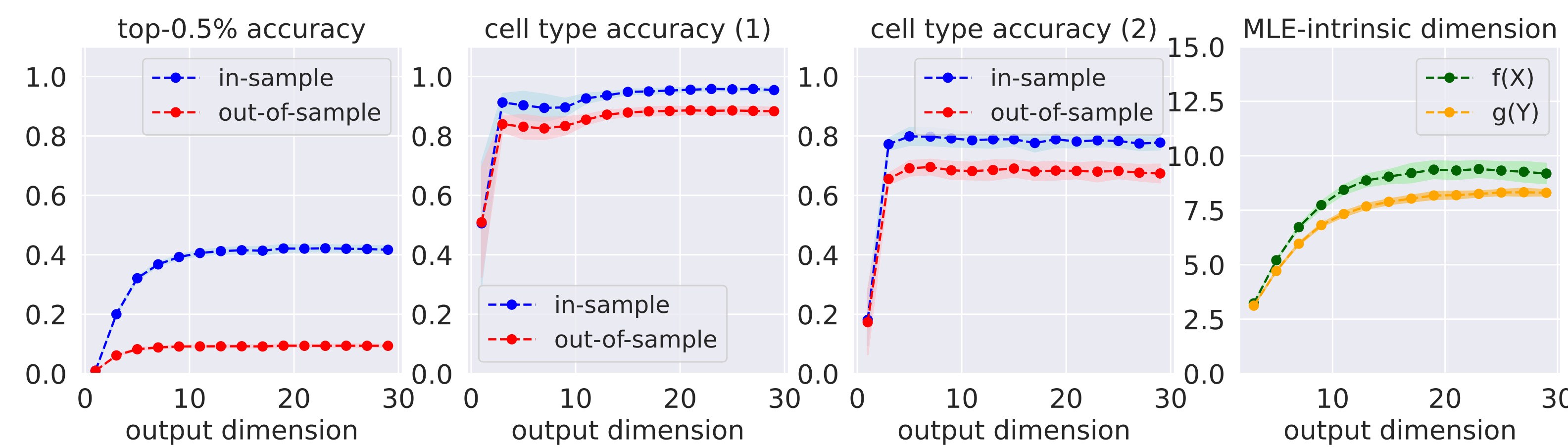
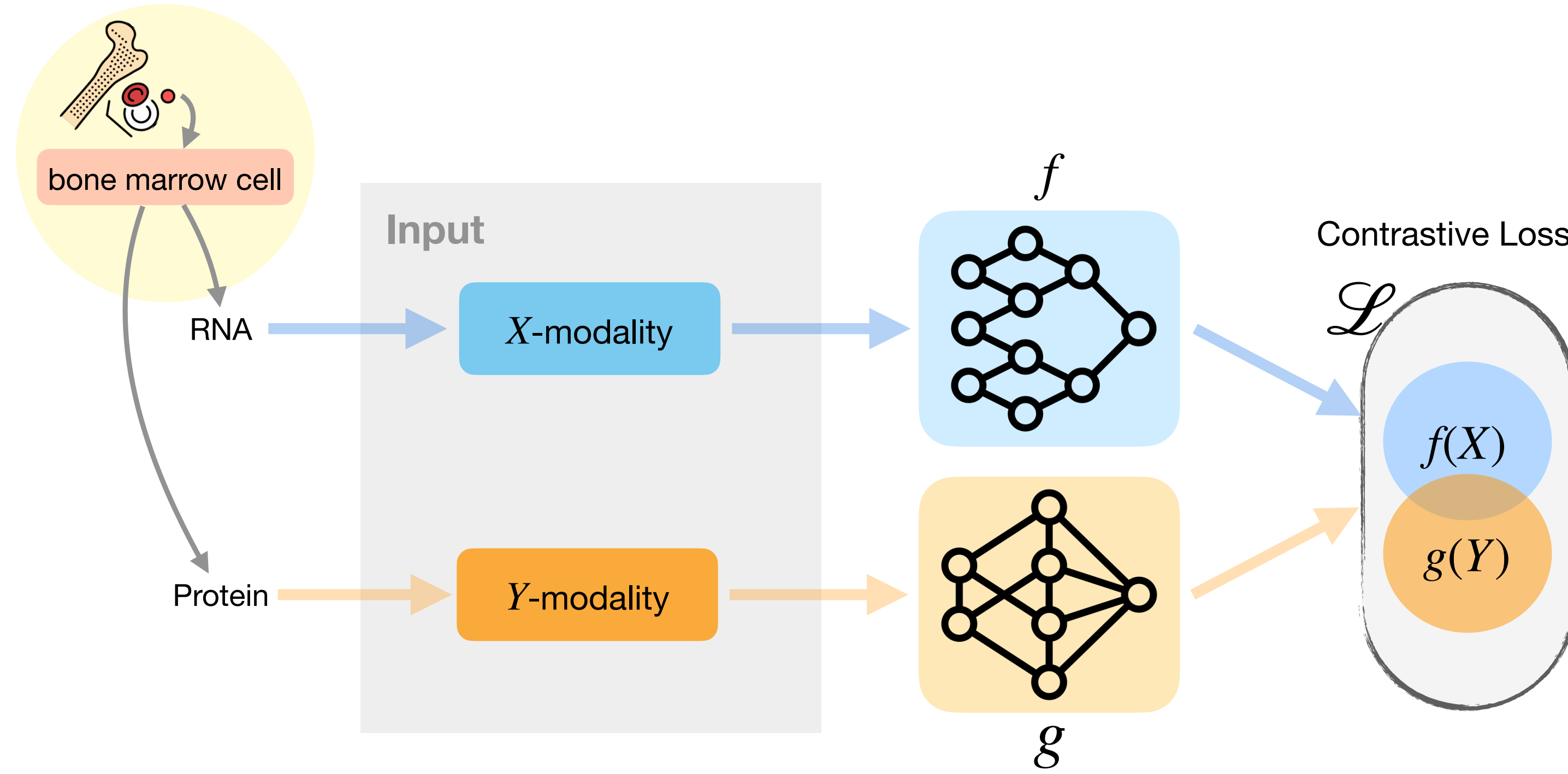


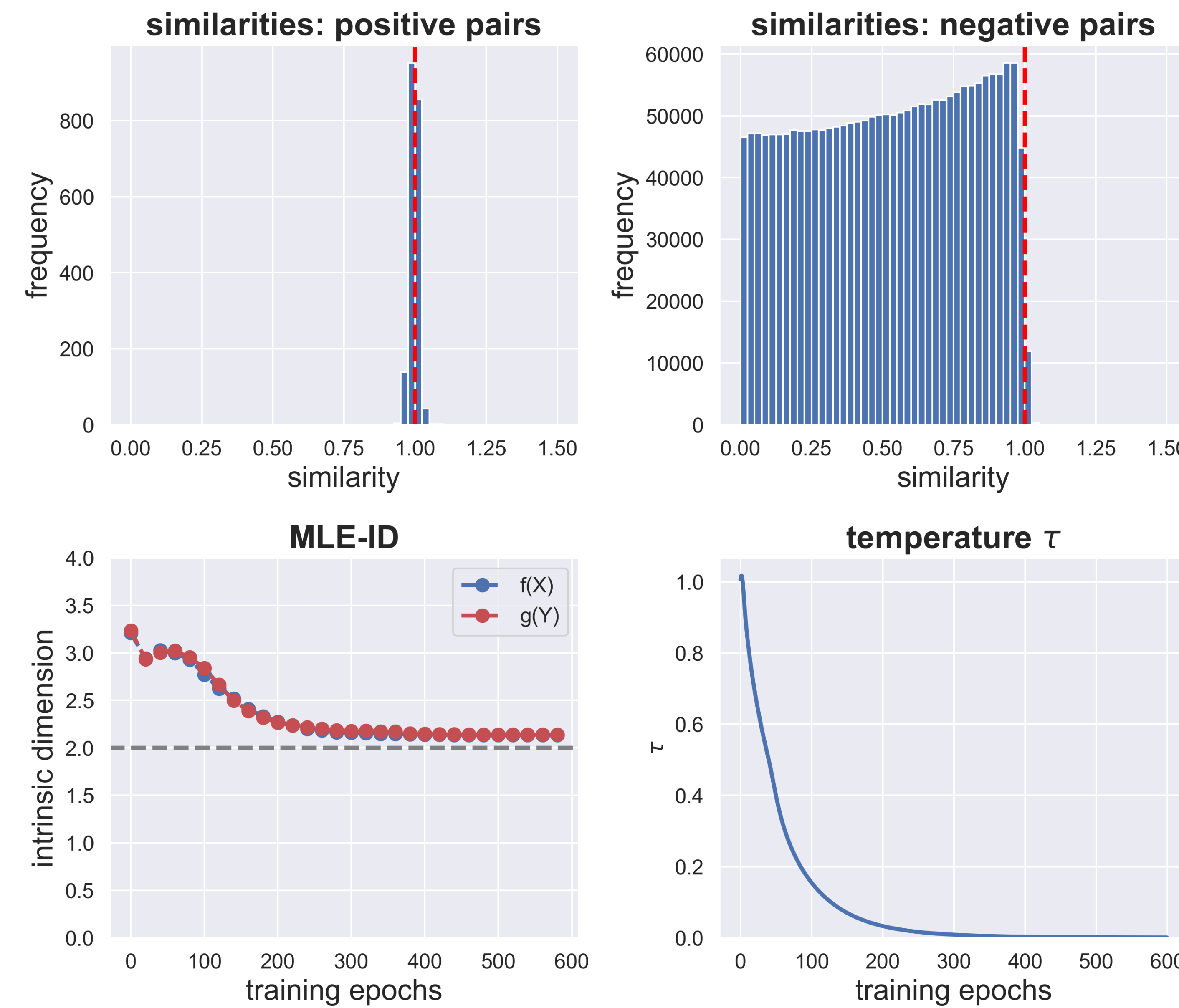
Figure 1. Results with CITE-seq dataset.

Question

- What representations does CLIP learn?
- What is the role of temperature τ ?

A toy example

$$Y_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \mathbf{I}_{20}), \quad \xi_i \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \mathbf{I}_{20-k^*}), \quad X_i = (Y_{i1}, \dots, Y_{ik^*}, \xi_i^\top)^\top$$



- Similarity concentration:** For positive pairs, similarities concentrates around 1, while negative pairs are capped by 1
- Intrinsic dimension adaptation:** Although output $d = 3$, representations with intrinsic dimension $k^* = 2$ are preferred
- Temperature convergence:** The optimized temperature $\tau \rightarrow 0$

Intrinsic dimension:

$$\text{ID}(f) = \min \left\{ k \in \mathbb{N} : f = \phi \circ h, h : \mathbb{R}^{d_x} \rightarrow \mathbb{R}^d, \phi : \mathbb{R}^d \rightarrow \mathbb{R}^d \text{ injective} \right\}$$

Ideal Representations

- Alignment:** with $m_\sigma(f, g) = \text{ess sup}_{X \perp \tilde{Y}} \sigma(f(X), g(\tilde{Y}))$

$$\mathcal{A}(\mathcal{H}) = \left\{ (f, g) \in \mathcal{H} : \frac{f(X)}{\mathbb{E}\|f(X)\|} = \frac{g(Y)}{\mathbb{E}\|g(Y)\|}, \quad \sigma(f(X), g(Y)) = m_\sigma(f, g) \text{ a.s.} \right\}$$

- Mutual information maximization:** $I_M^*(\mathcal{H}) = \sup_{\mathcal{H}} I(f_M(X); g_M(Y))$

$$\mathcal{W}(\mathcal{H}) = \left\{ (f, g) \in \mathcal{H} : \liminf_{M \rightarrow +\infty} (I(f_M(X); g_M(Y)) - I_M^*(\mathcal{H})) \geq 0 \right\}$$

$$\mathcal{V}(\mathcal{H}) = \mathcal{A}(\mathcal{H}) \cap \mathcal{W}(\mathcal{H})$$

- Intrinsic dimension adaptation:**

Lemma. Suppose $\mathcal{V}(\mathcal{H}) \neq \emptyset$. Then, for all $(f, g) \in \mathcal{V}(\mathcal{H})$, we have $\text{ID}(f) = \text{ID}(g) = k^*$, i.e., maps in $\mathcal{V}(\mathcal{H})$ have the same intrinsic dimension k^*

Is any (approximate) minimizer of CLIP ideal?

$$\mathcal{O}_{\mathcal{L}, \eta}(\mathcal{H}) = \left\{ (f, g) \in \mathcal{H} : \exists \tau \geq \varepsilon(\eta), \limsup_{M \rightarrow +\infty} (\mathcal{L}(f_M, g_M, \tau) + 2I_M^*(\mathcal{H})) \leq 2\eta \right\}$$

Main results [1]

$$\mathcal{V}(\mathcal{H}) \neq \emptyset \implies \bigcap_{\eta \geq 0} \mathcal{O}_{\mathcal{L}, \eta}(\mathcal{H}) \neq \emptyset.$$

In addition, for any $(f, g) \in \bigcap_{\eta \geq 0} \mathcal{O}_{\mathcal{L}, \eta}(\mathcal{H})$, with $\sigma(f(X), g(Y)) = \frac{\langle f(X), g(Y) \rangle}{\mathbb{E}\|f(X)\| \cdot \mathbb{E}\|g(Y)\|}$,

- (similarity maximization) $\sigma(f(X), g(Y)) = m_\sigma(f, g)$ almost surely
- (intrinsic dimension adaptation) $\text{ID}(f) = \text{ID}(g) = k^*$
- (monotonicity in temperature) $\mathcal{L}(f, g, \tau)$ is increasing in τ
- (mutual information maximization) $(f, g) \in \mathcal{W}(\mathcal{H})$

References

- Gui, Y., Ma, C., and Ma, Z. (2025). Multi-modal contrastive learning adapts to intrinsic dimensions of shared latent variables. *NeurIPS*.
- Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al. (2021). Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR.