

The Benefits of Balance

From Information Projections to Variance Reduction



NeurIPS 2024



Team



Lang Liu
University of
Washington



Ronak Mehta
University of
Washington



Zaid Harchaoui
University of
Washington



Learning Transferable Visual Models From Natural Language Supervision

Alec Radford^{*1} Jong Wook Kim^{*1} Chris Hallacy¹ Aditya Ramesh¹ Gabriel Goh¹ Sandhini Agarwal¹
Girish Sastry¹ Amanda Askell¹ Pamela Mishkin¹ Jack Clark¹ Gretchen Krueger¹ Ilva Sutskever¹

Unsupervised Learning of Visual Features by Contrasting Cluster Assignments

Mathilde Caron^{1,2} Ishan Misra² Julien Mairal¹
Priya Goyal² Piotr Bojanowski² Armand Joulin²

¹ Inria^{*}

² Facebook AI Research

SELF-LABELLING VIA SIMULTANEOUS CLUSTERING AND REPRESENTATION LEARNING

Yuki M. Asano Christian Rupprecht Andrea Vedaldi

Visual Geometry Group

cs.ox.ac.uk

DEMYSTIFYING CLIP DATA

Hu Xu¹ Saining Xie² Xiaoqing Ellen Tan¹ Po-Yao Huang¹ Russell Howes¹ Vasu Sharma¹
Shang-Wen Li¹ Gargi Ghosh¹ Luke Zettlemoyer^{1,3} Christoph Feichtenhofer¹
¹FAIR, Meta AI ²New York University ³University of Washington

DINOv2: Learning Robust Visual Features without Supervision

Maxime Oquab^{**}, Timothée Darcet^{**}, Théo Moutakanni^{**},
Hieu Vo^{*}, Vasil Khalidov^{*}, Pierre Fernandez, Daniel Haziza,
Vedant Matulkar, Nouby, Mahmoud Assran, Nicolas Ballas, Wojciech Galuba,
Michael Rabbat, Han Hu, Huo, Guang, Shang-Wen Li, Ishan Misra, Michael Rabbat,
Armand Joulin, Daniel Synnaeve, Hu Xu, Hervé Jegou, Julien Mairal¹,
Michael Neumann, Olivier Tschumper, Dávid Vincenty, Armand Joulin^{*}, Piotr Bojanowski^{*}

¹ Meta AI Research ¹ Inria

core team ^{**}equal contribution

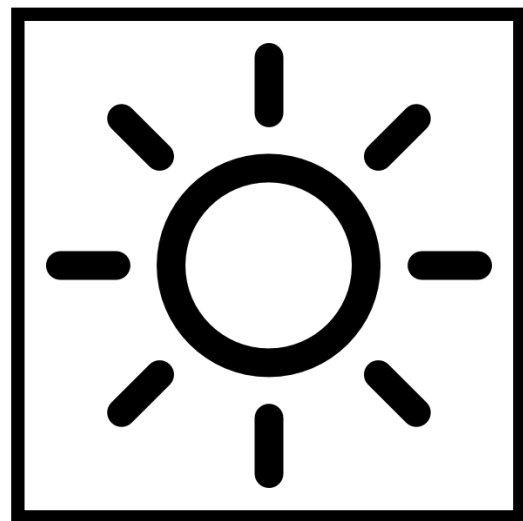
DATAComp: In search of the next generation of multimodal datasets

Aditya Khosla¹, Gabriel Ilharco^{*1}, Alex Fang^{*1}, Jonathan Hayase¹,
Nikhil Doherty⁵, Thao Nguyen¹, Ryan Marten^{7,9}, Mitchell Wortsman¹,
Sreyas Jayaram¹, Eyal Orgad³, Rahim Entezari¹⁰, Giannis Daras⁵,
Vivek Ramanujan¹, Yonatan Bitton¹¹, Kalyani Marathe¹,
Stephen Mussmann¹, Richard Vencu⁶, Mehdi Cherti^{6,8}, Ranjay Krishna¹,
Pang Wei Koh^{1,12}, Olga Saukh¹⁰, Alexander Ratner^{1,13}, Shuran Song²,
Hannaneh Hajishirzi^{1,7}, Ali Farhadi¹, Romain Beaumont⁶,
Sewoong Oh¹, Alex Dimakis⁵, Jenia Jitsev^{6,8},
Yair Carmon³, Vaishaal Shankar⁴, Ludwig Schmidt^{1,6,7}

Discriminative clustering with representation learning with any ratio of labeled to unlabeled data

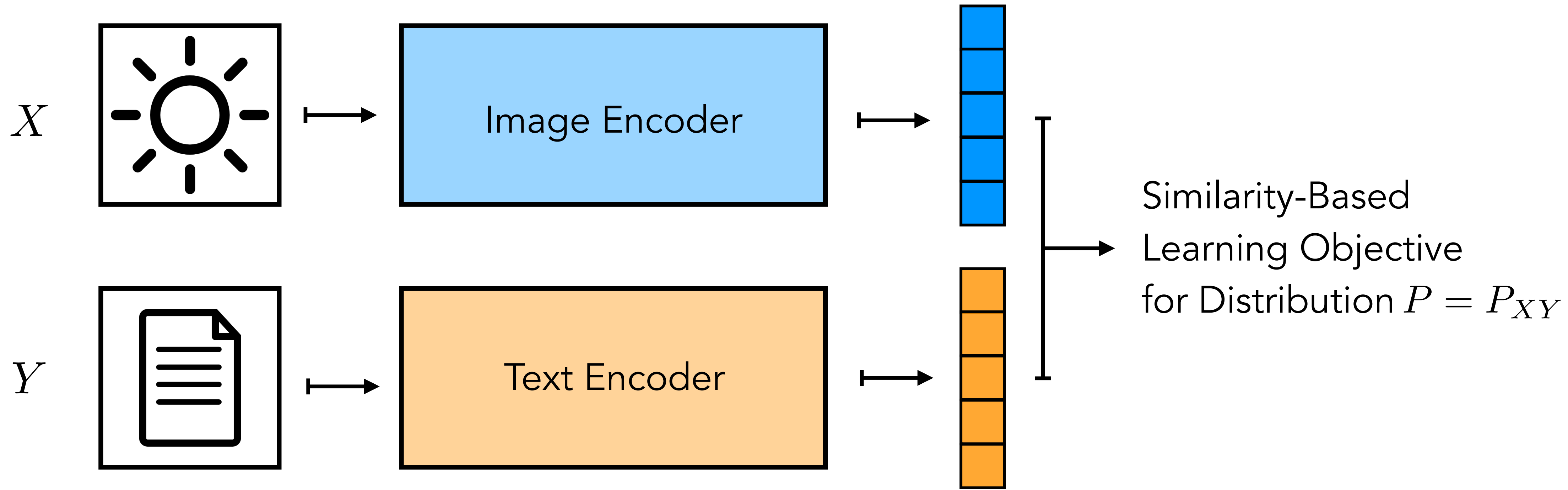
Corinne Jones¹ · Vincent Roulet² · Zaid Harchaoui²

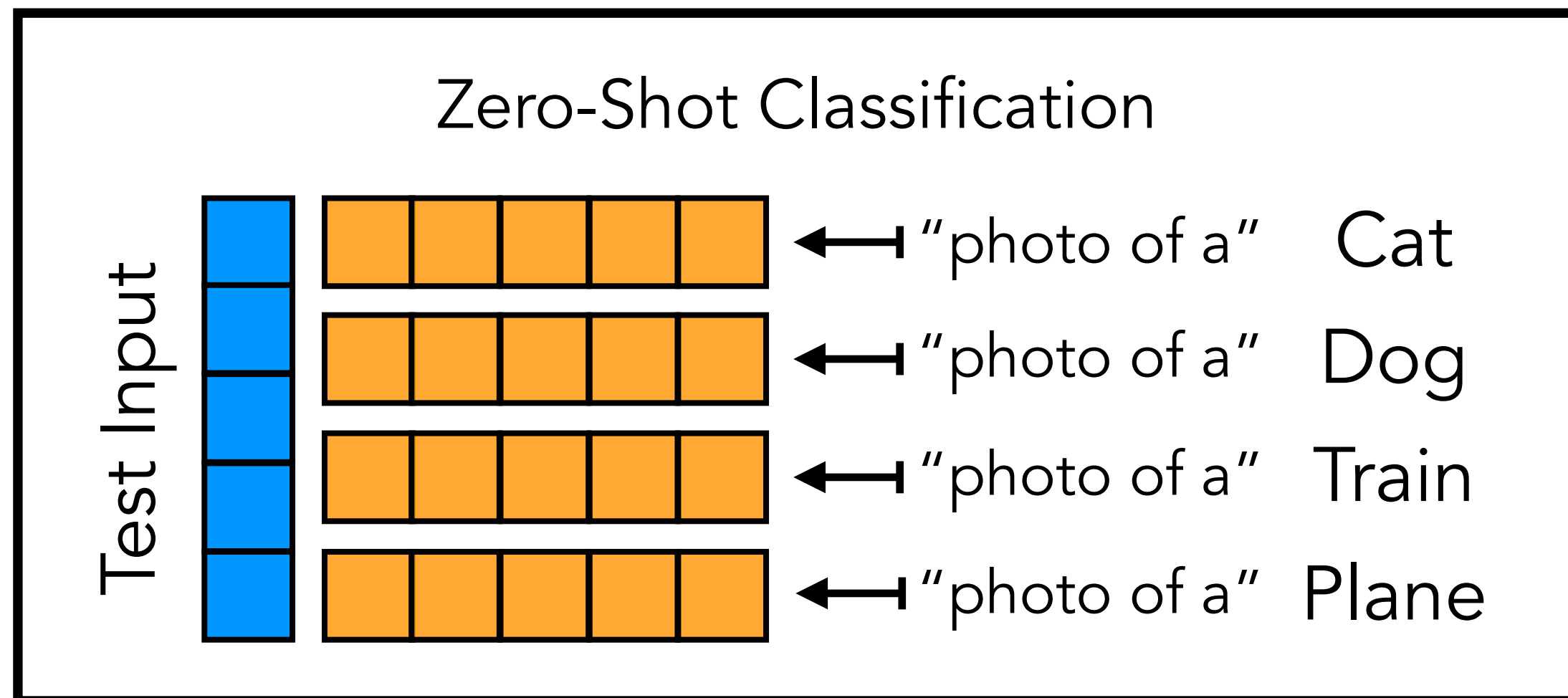
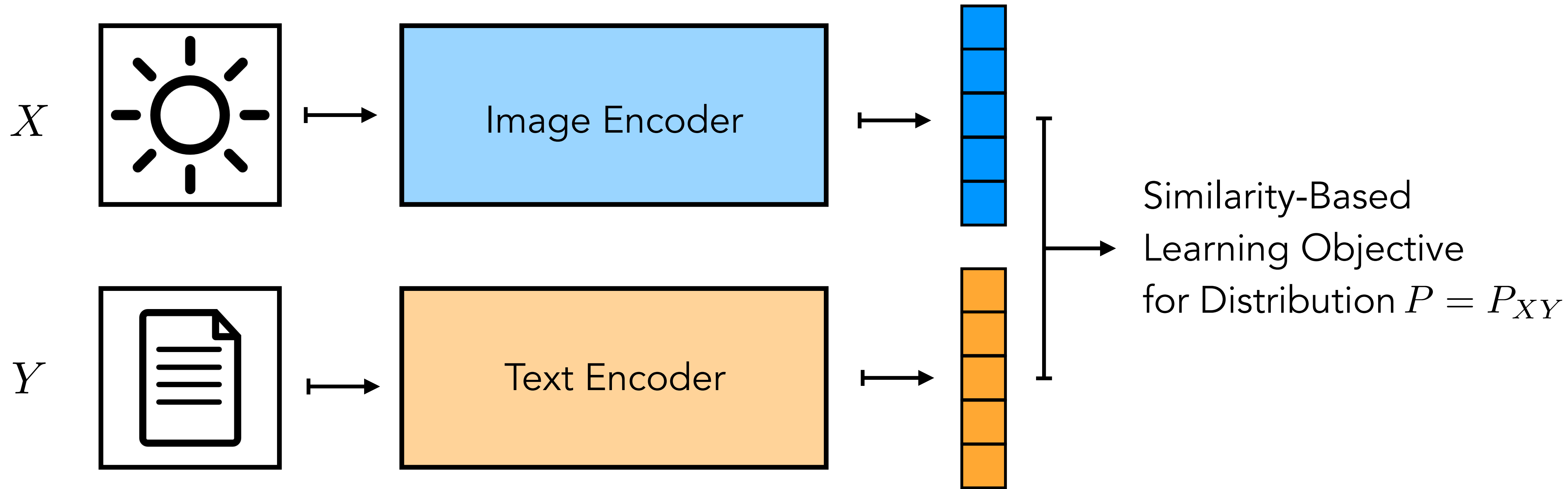
X

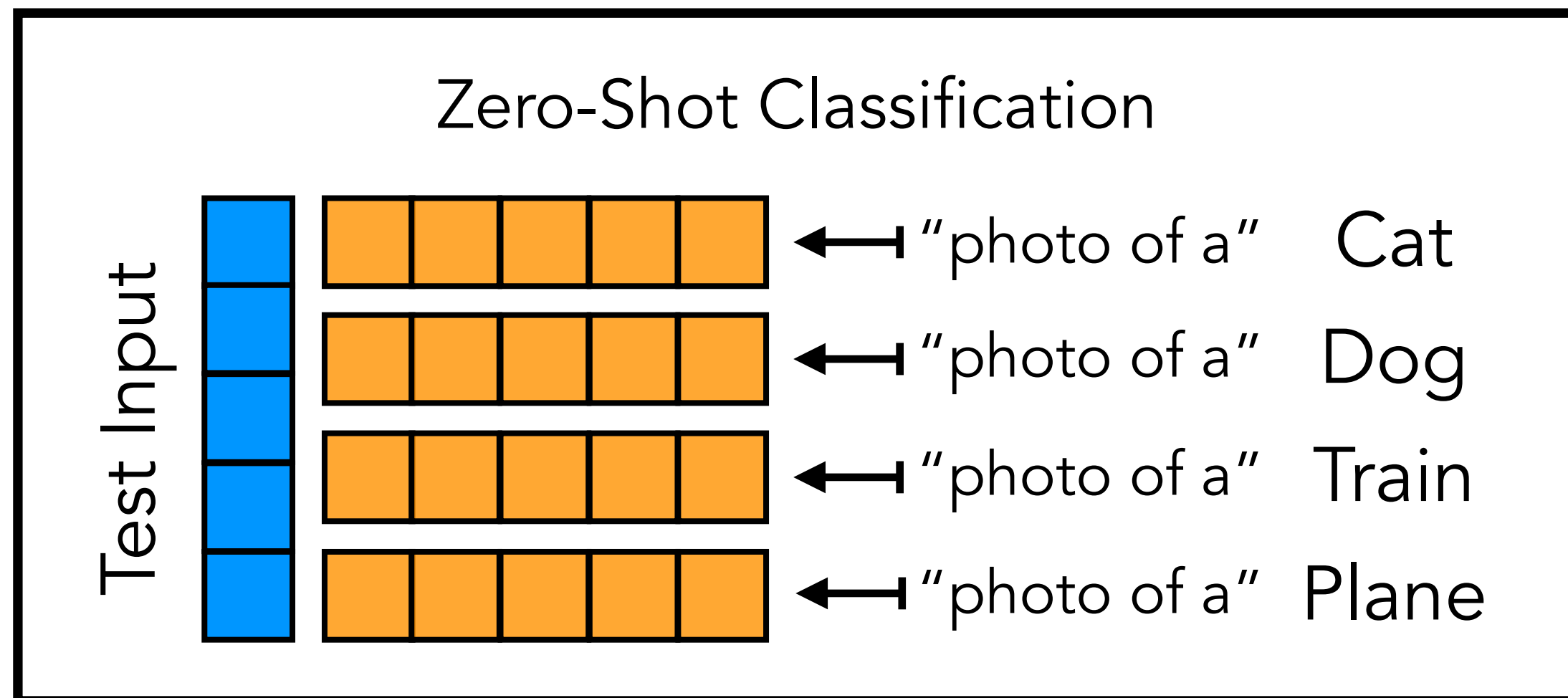
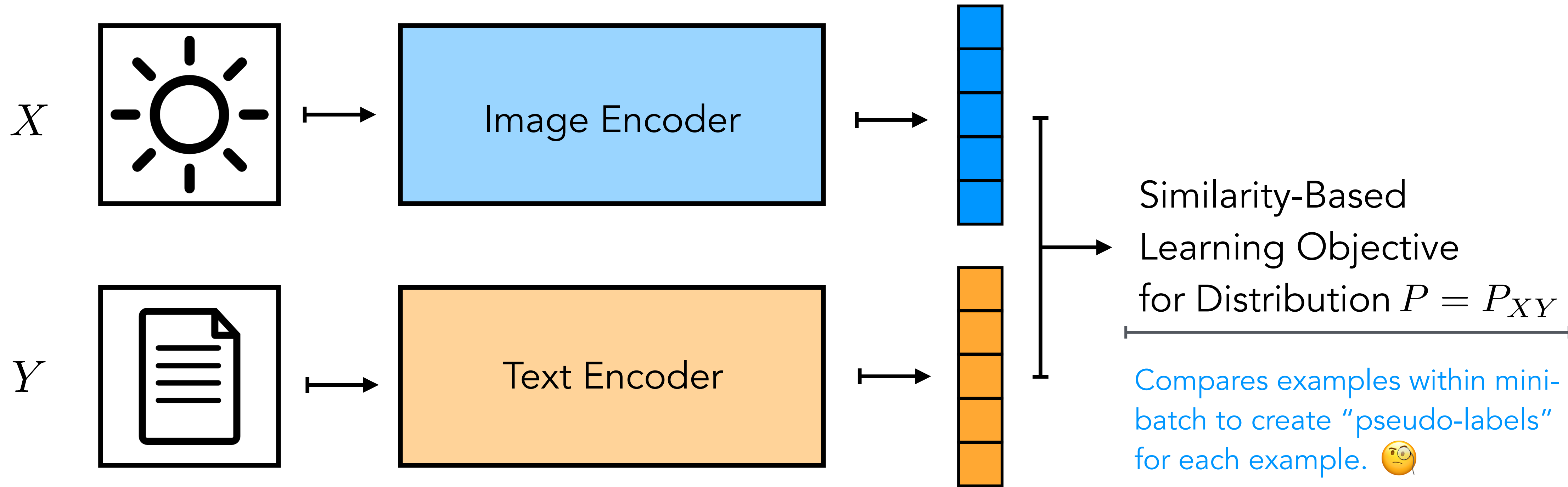


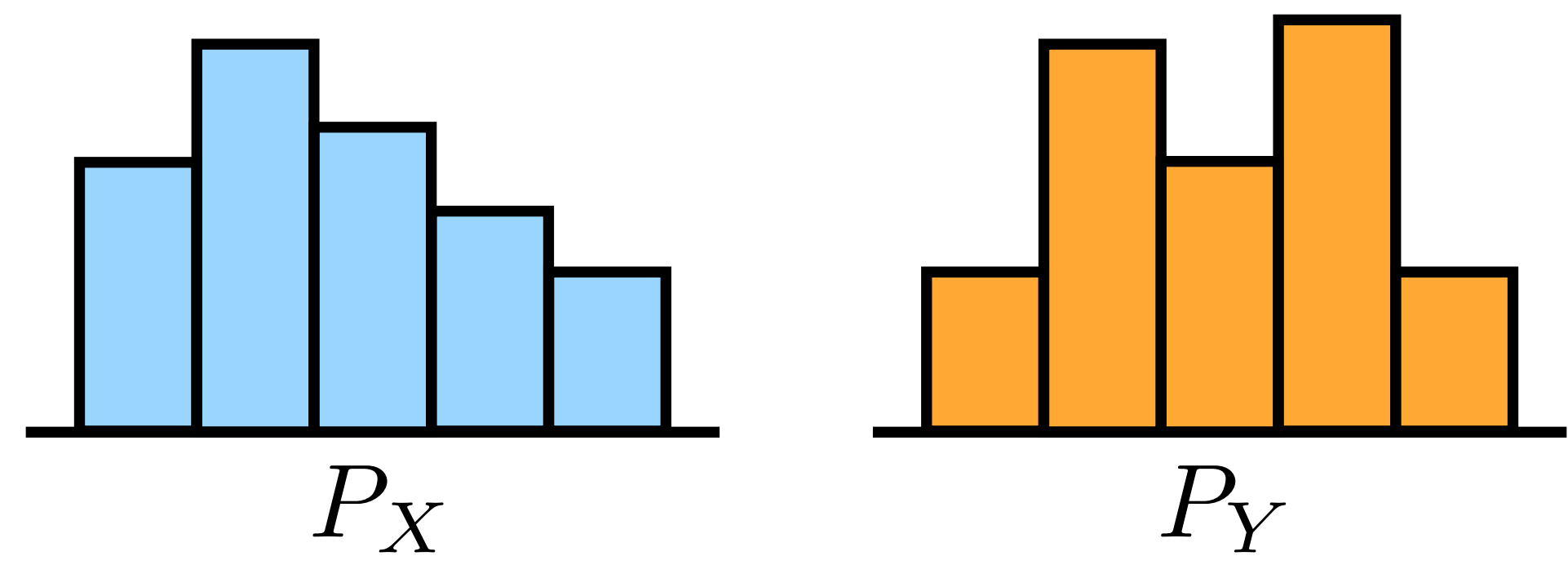
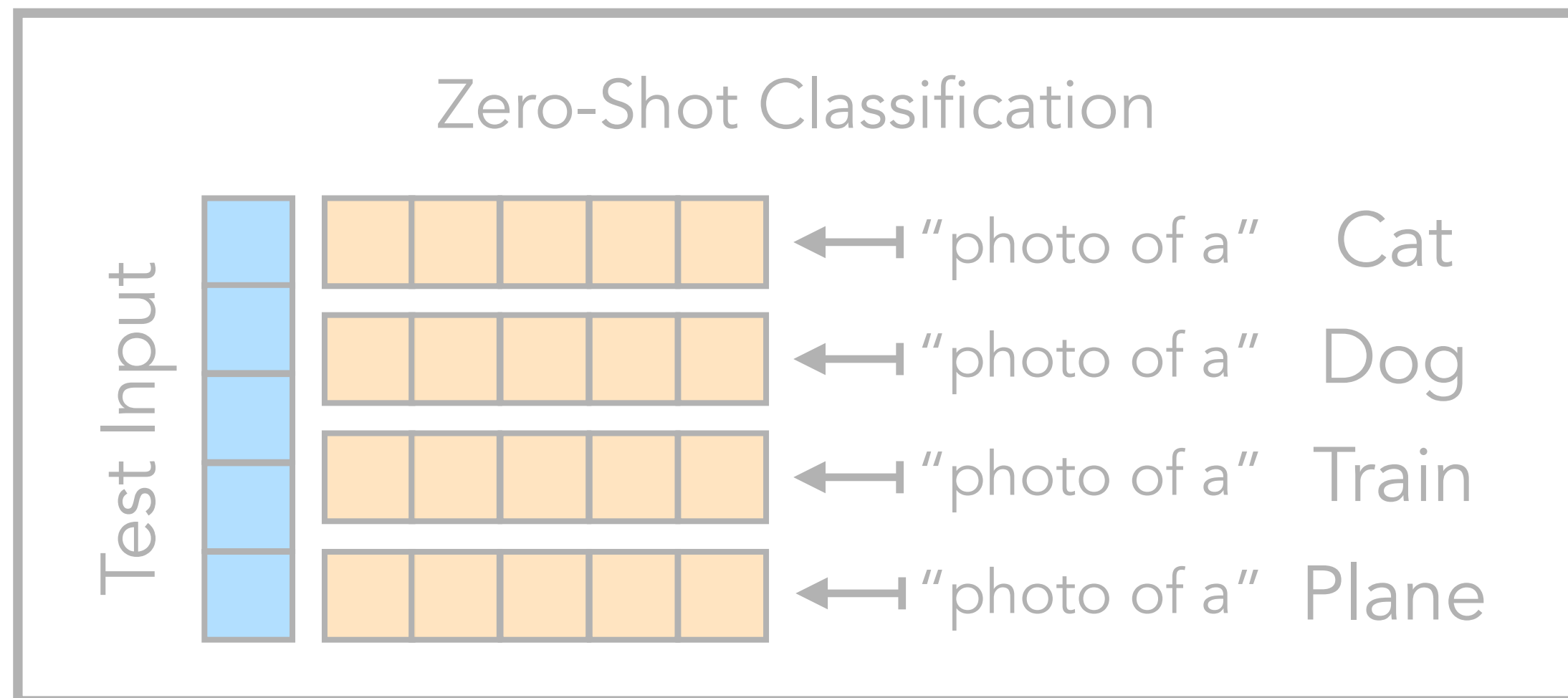
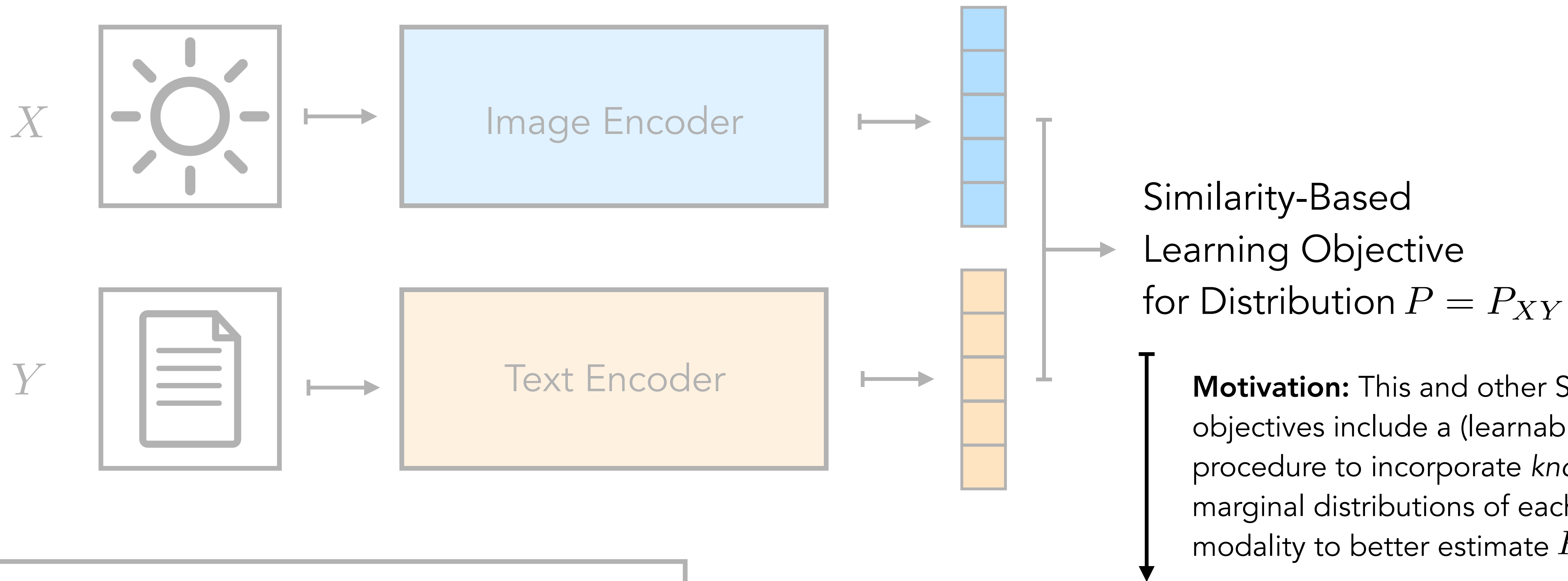
Y











Problem Setting

Data from **unknown** joint probability distribution.

$$(X_1, Y_1), \dots, (X_n, Y_n) \sim P$$

Access to **known** marginal distributions.

$$(P_X, P_Y)$$

Goal: estimate the parameter:

$$P(h) = \mathbb{E}_{(X,Y) \sim P} [h(X, Y)]$$

and characterize how the marginals improve upon

$$P_n(h) = \frac{1}{n} \sum_{i=1}^n h(X_i, Y_i)$$

Problem Setting

Data from **unknown** joint probability distribution.

$$(X_1, Y_1), \dots, (X_n, Y_n) \sim P$$

Access to **known** marginal distributions.

$$(P_X, P_Y)$$

Goal: estimate the parameter:

$$P(h) = \mathbb{E}_{(X,Y) \sim P} [h(X, Y)]$$

and characterize how the marginals improve upon

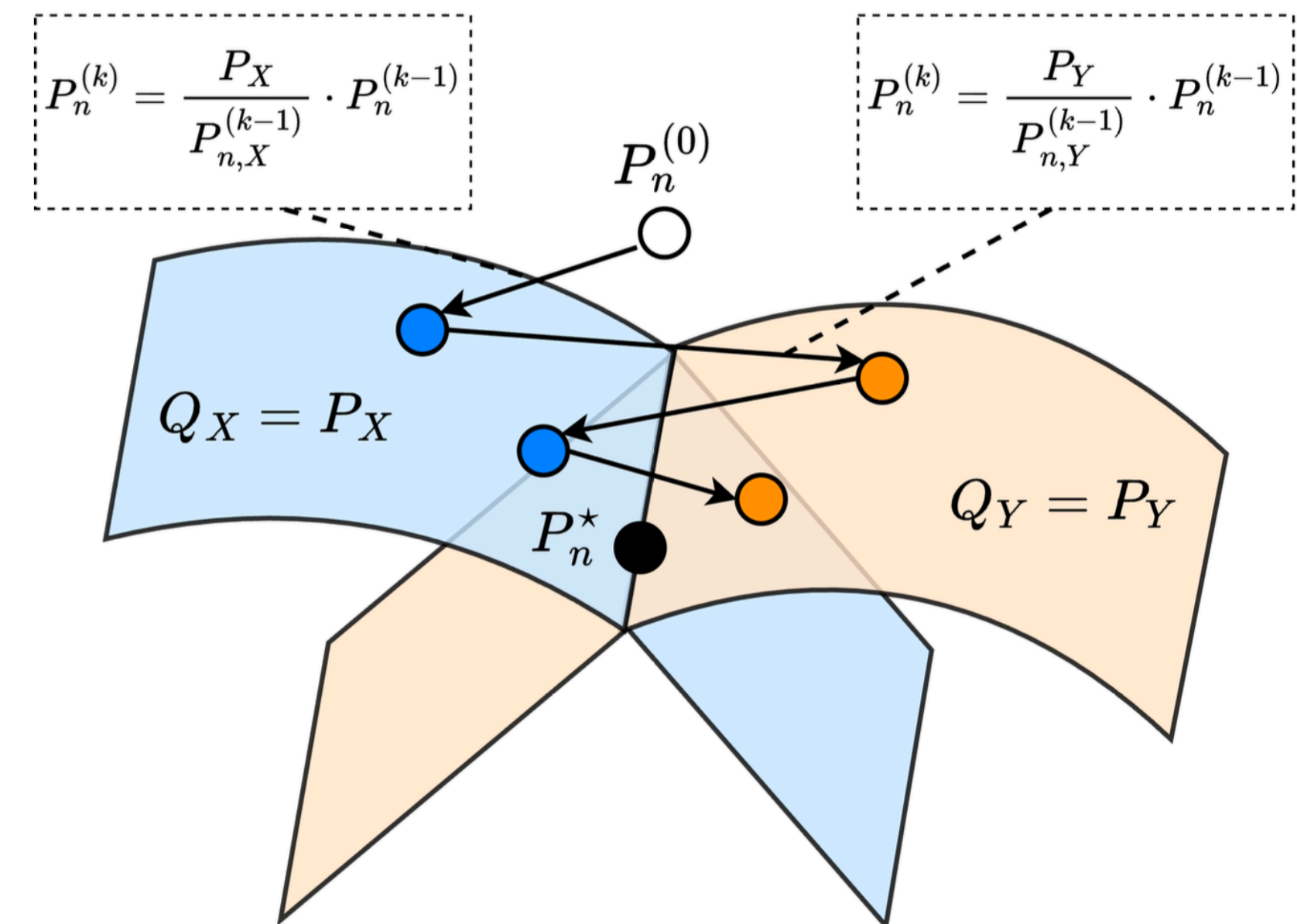
$$P_n(h) = \frac{1}{n} \sum_{i=1}^n h(X_i, Y_i)$$

Data Balancing

Rescale rows and columns by the desired marginals.

$$P_n^{(0)} = P_n = \frac{1}{n} \sum_{i=1}^n \delta_{(X_i, Y_i)}$$

$$P_n^{(k)} = \begin{cases} \frac{P_X}{P_{n,X}^{(k-1)}} \cdot P_n^{(k-1)} & k \text{ odd} \\ \frac{P_Y}{P_{n,Y}^{(k-1)}} \cdot P_n^{(k-1)} & k \text{ even} \end{cases}$$



Problem Setting

Data from **unknown** joint probability distribution.

$$(X_1, Y_1), \dots, (X_n, Y_n) \sim P$$

Access to **known** marginal distributions.

$$(P_X, P_Y)$$

Goal: estimate the parameter:

$$P(h) = \mathbb{E}_{(X,Y) \sim P} [h(X, Y)]$$

and characterize how the marginals improve upon

$$P_n(h) = \frac{1}{n} \sum_{i=1}^n h(X_i, Y_i)$$

Data Balancing

Rescale rows and columns by the desired marginals.

$$P_n^{(0)} = P_n = \frac{1}{n} \sum_{i=1}^n \delta_{(X_i, Y_i)}$$
$$P_n^{(k)} = \begin{cases} \frac{P_X}{P_{n,X}^{(k-1)}} \cdot P_n^{(k-1)} & k \text{ odd} \\ \frac{P_Y}{P_{n,Y}^{(k-1)}} \cdot P_n^{(k-1)} & k \text{ even} \end{cases}$$

1. How does the balanced distribution improve upon the empirical measure **theoretically**?
2. What are the **practical** implications for SSL objectives such as CLIP?

Problem Setting

Data from **unknown** joint probability distribution.

$$(X_1, Y_1), \dots, (X_n, Y_n) \sim P$$

Access to **known** marginal distributions.

$$(P_X, P_Y)$$

Goal: estimate the parameter:

$$P(h) = \mathbb{E}_{(X,Y) \sim P} [h(X, Y)]$$

and characterize how the marginals improve upon

$$P_n(h) = \frac{1}{n} \sum_{i=1}^n h(X_i, Y_i)$$

Data Balancing

Rescale rows and columns by the desired marginals.

$$P_n^{(0)} = P_n = \frac{1}{n} \sum_{i=1}^n \delta_{(X_i, Y_i)}$$
$$P_n^{(k)} = \begin{cases} \frac{P_X}{P_{n,X}^{(k-1)}} \cdot P_n^{(k-1)} & k \text{ odd} \\ \frac{P_Y}{P_{n,Y}^{(k-1)}} \cdot P_n^{(k-1)} & k \text{ even} \end{cases}$$

Theorem. The iterates of balancing satisfy

$$\mathbb{E}_P \left| P_n^{(k)}(h) - P(h) \right|^2 = \frac{\sigma_k^2}{n} + \tilde{O} \left(\frac{k^6}{n^{3/2}} \right) \rightarrow \frac{\sigma_0^2 - \sigma_{\text{gap}}^2}{n}$$

Novel recursion formula for estimation error that is of independent interest (OT, data-centric ML, etc.)

Problem Setting

Data from **unknown** joint probability distribution.

$$(X_1, Y_1), \dots, (X_n, Y_n) \sim P$$

Access to **known** marginal distributions.

$$(P_X, P_Y)$$

Goal: estimate the parameter:

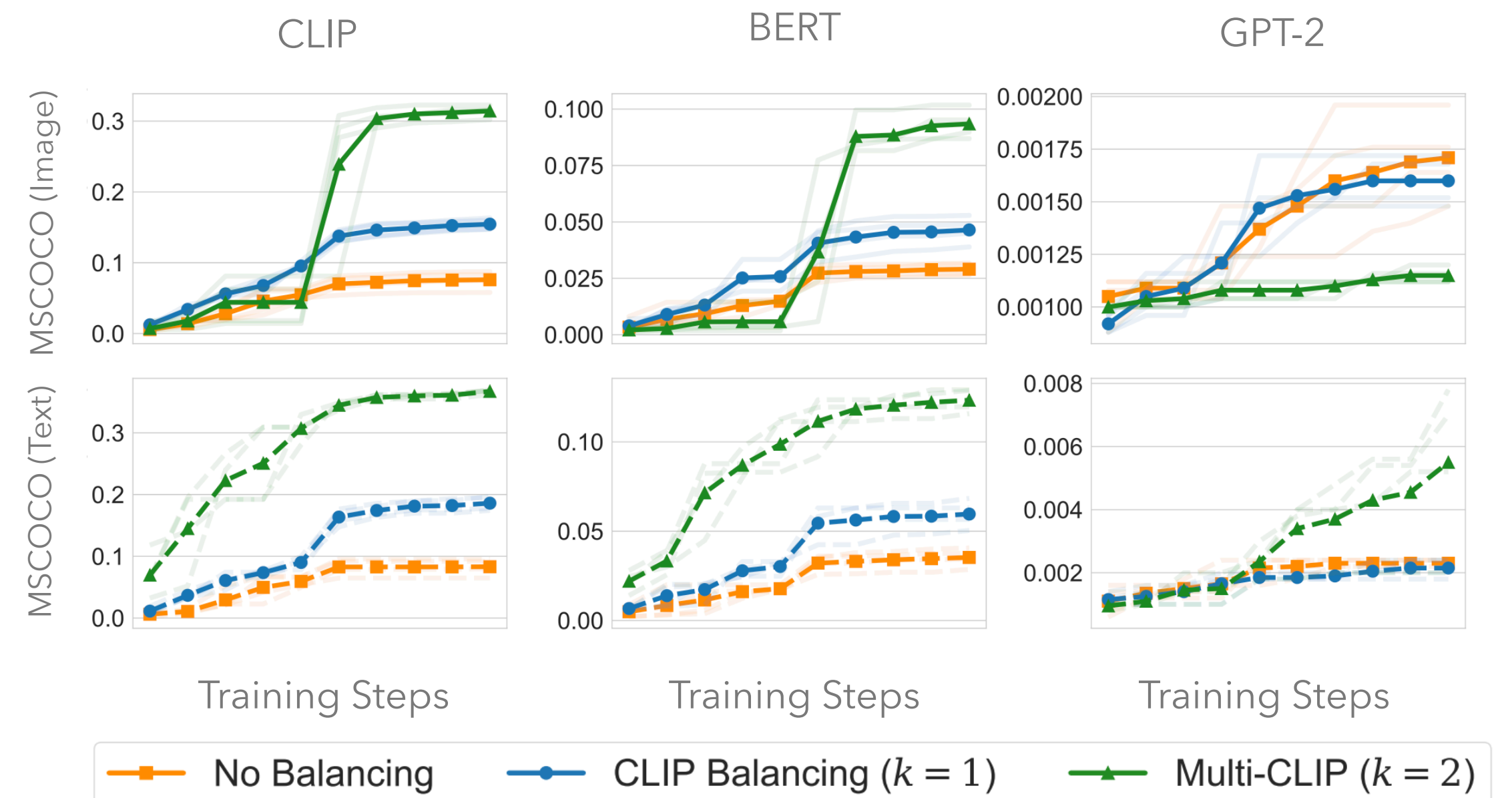
$$P(h) = \mathbb{E}_{(X,Y) \sim P} [h(X, Y)]$$

and characterize how the marginals improve upon

$$P_n(h) = \frac{1}{n} \sum_{i=1}^n h(X_i, Y_i)$$

Balancing mini-batches to improve the stability of the CLIP training objective.

Using a balanced objective increases zero-shot retrieval (recall) across datasets and embedding architectures.



Problem Setting

Data from **unknown** joint probability distribution.

$$(X_1, Y_1), \dots, (X_n, Y_n) \sim P$$

Access to **known** marginal distributions.

$$(P_X, P_Y)$$

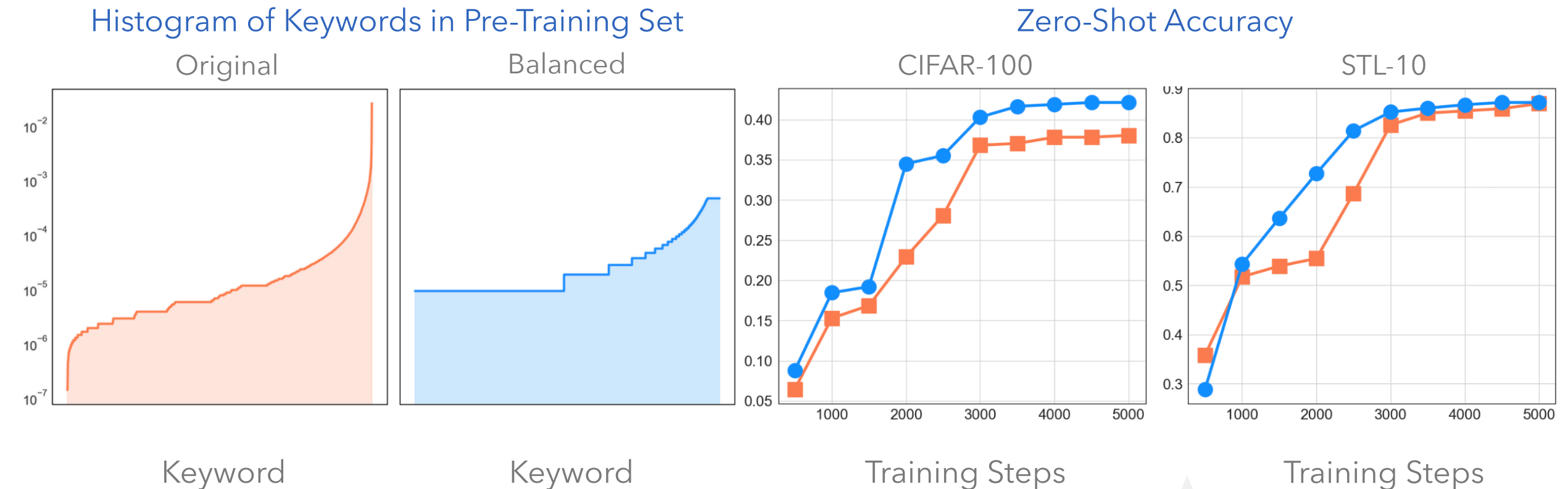
Goal: estimate the parameter:

$$P(h) = \mathbb{E}_{(X,Y) \sim P} [h(X, Y)]$$

and characterize how the marginals improve upon

$$P_n(h) = \frac{1}{n} \sum_{i=1}^n h(X_i, Y_i)$$

CLIP models trained on the balanced pre-training set improve over those trained on the original.



Balancing at scale improves performance on zero-shot classification.

Problem Setting

Data from **unknown** joint probability distribution.

$$(X_1, Y_1), \dots, (X_n, Y_n) \sim P$$

Access to **known** marginal distributions.

$$(P_X, P_Y)$$

Goal: estimate the parameter:

$$P(h) = \mathbb{E}_{(X,Y) \sim P} [h(X, Y)]$$

and characterize how the marginals improve upon

$$P_n(h) = \frac{1}{n} \sum_{i=1}^n h(X_i, Y_i)$$

Thank you!

