# The Benefits of Balance: From Information Projections to Variance Reduction

Lang Liu, Ronak Mehta, Soumik Pal, Zaid Harchaoui

## Data Balancing

**Motivation:** High-quality, large-scale datasets of paired observations (features + labels, images + captions) are scarce, while unpaired observations might be abundant.

$$(X_1, Y_1), \ldots, (X_n, Y_n) \overset{\text{i.i.d}}{\sim} P$$

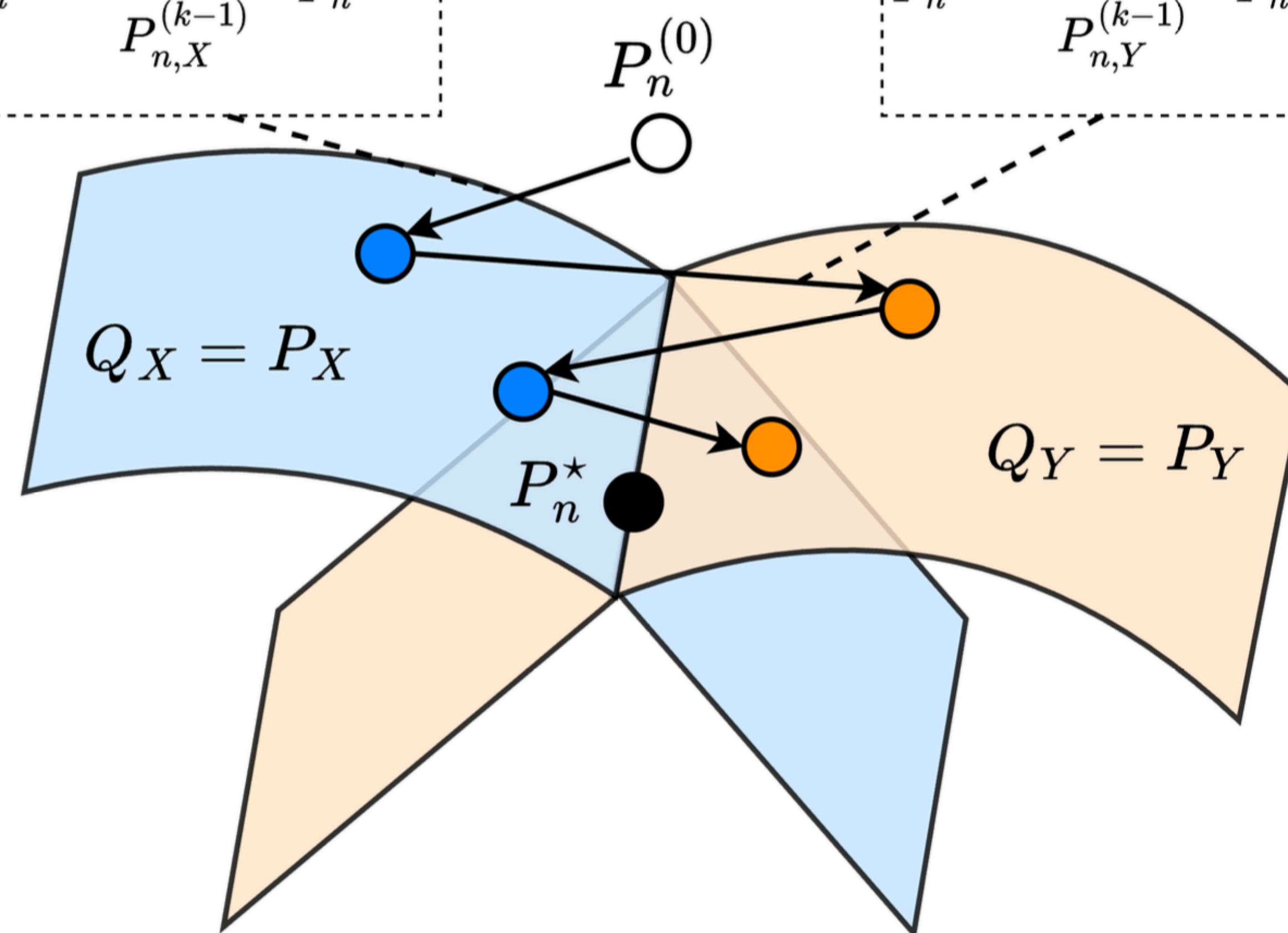marginal distributions (**known**) $(P_X, P_Y)$ joint distribution (**unknown**)

How can we incorporate marginal information?

empirical measure
$$P_n^{(0)} = P_n = \frac{1}{n} \sum_{i=1}^{n} \delta_{(X_i, Y_i)}$$
marginal likelihood ratio

$$P_n^{(k)} = \frac{P_X}{P_{n,X}^{(k-1)}} \cdot P_n^{(k-1)} \qquad P_n^{(k)} = \frac{P_Y}{P_{n,Y}^{(k-1)}} \cdot P_n^{(k-1)}$$

$P_n^{(0)}$

$Q_X = P_X$

$P_n^\star$

$Q_Y = P_Y$

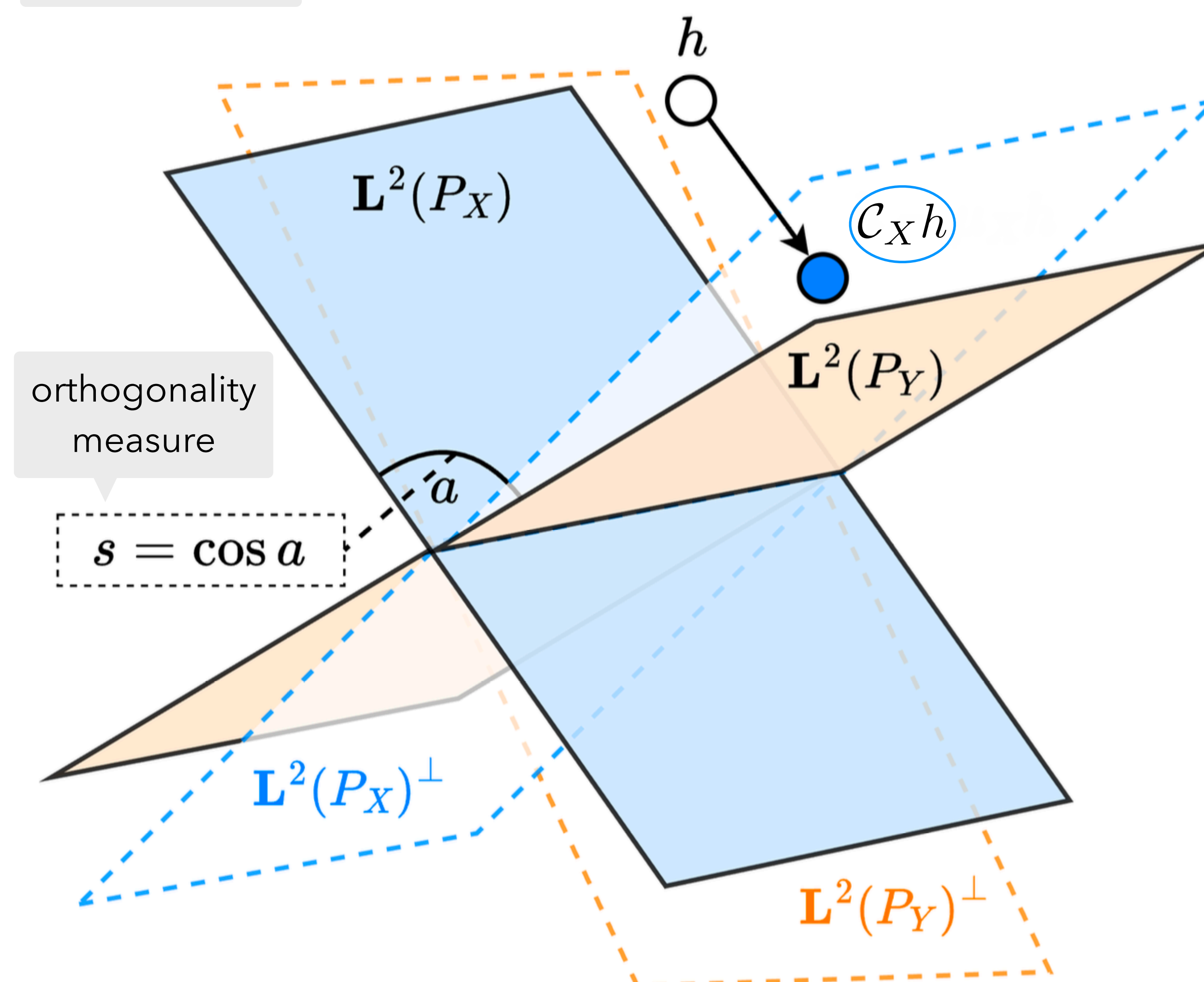| **Estimand** | $P(h) = \mathbb{E}_P[h(X,Y)]$ |
| --- | --- |
| **Estimator** | $P_n^{(k)}(h) = \mathbb{E}_{P_n^{(k)}}[h(X,Y)]$ |

How does balancing improve estimation and learning?

## Information Projections → Orthogonal Projections

$$[P_n^{(k)} - P](h) = [P_n^{(k-1)} - P](\mathcal{C}_X h) + O_p(n^{-1})$$
$$= [P_n^{(k-2)} - P](\mathcal{C}_Y \mathcal{C}_X h) + O_p(n^{-1})$$

novel recursion formula
$$= [P_n - P](\mathcal{C}_Y \mathcal{C}_X \ldots \mathcal{C}_Y \mathcal{C}_X h) + O_p(n^{-1})$$

$h$

$\mathbf{L}^2(P_X)$

$\mathcal{C}_X h$

$\mathbf{L}^2(P_Y)$

orthogonality measure

$a$

$$s = \cos a$$

$\mathbf{L}^2(P_X)^\perp$

$\mathbf{L}^2(P_Y)^\perp$

## Orthogonal Projections → Variance Reduction

We compare the mean squared errors of the empirical versus balanced mean.

$$\sigma^2 = \mathrm{Var}[h(X,Y)] \implies \mathrm{Var}[P_n(h)] = \frac{\sigma^2}{n}$$
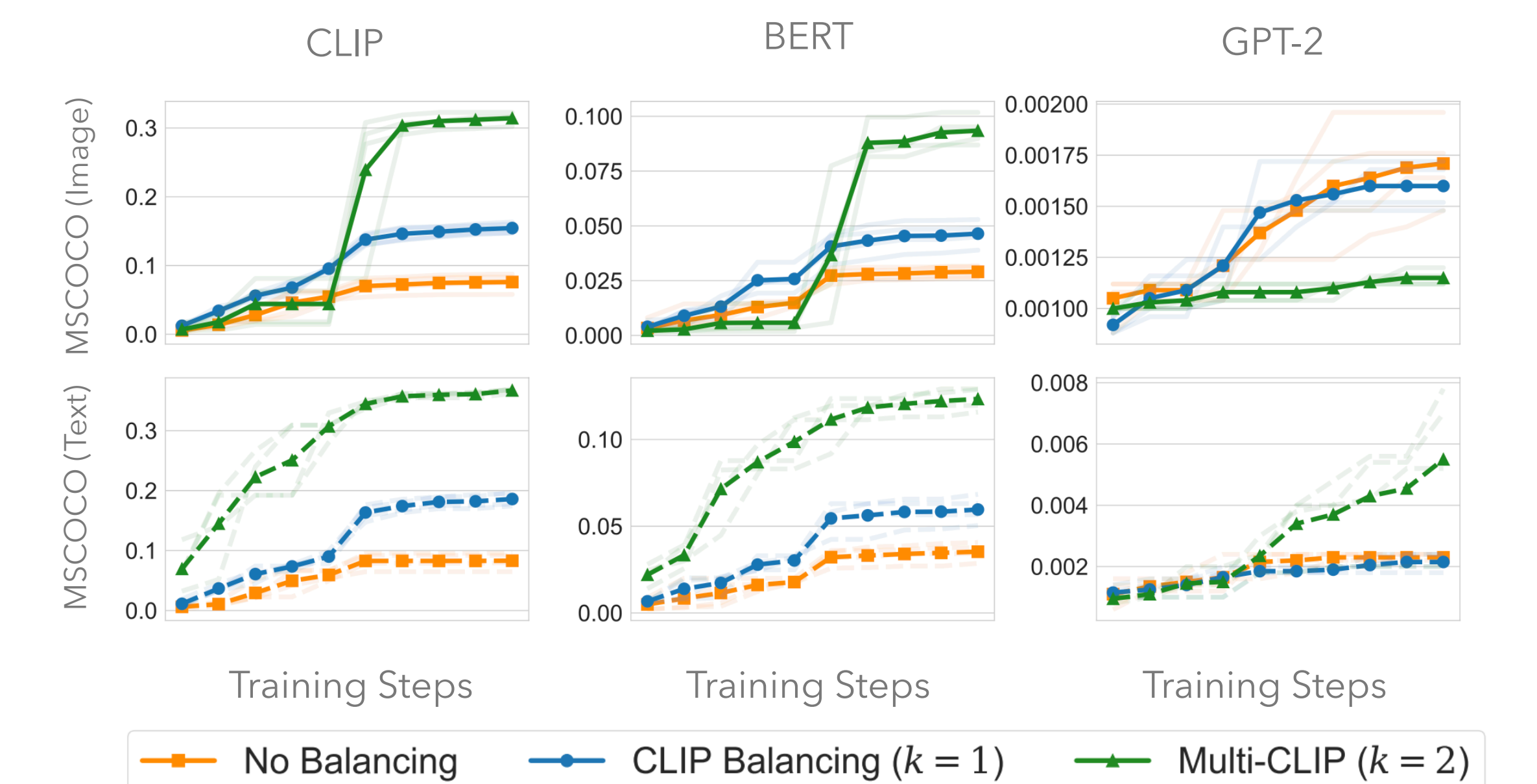
**Theorem.** The iterates of balancing satisfy

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

The quantity $s \in [0, 1)$ can be computed via the spectral properties of the two conditional mean operators.
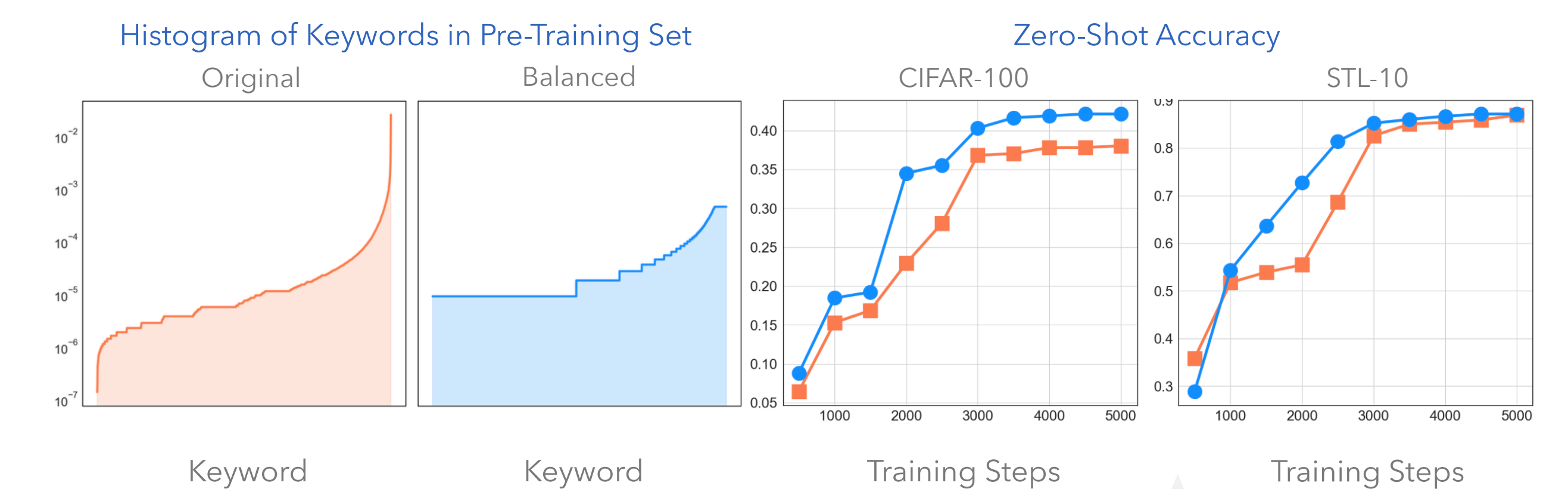
## Experiments

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.



No Balancing — CLIP Balancing ($k = 1$) — Multi-CLIP ($k = 2$)

Comparing CLIP models when balancing the entire pre-training set.



Histogram of Keywords in Pre-Training Set — Zero-Shot Accuracy

Balancing at scale improves performance on zero-shot classification.

Understanding performance under marginal misspecification.

Performance is resistant to marginal corruption.

Code



$\varepsilon = 0.25$ (IPWI)
$\varepsilon = 0.125$ (IPWI)
$\varepsilon = 0.0$ (IPWI)
$\varepsilon = 0.25$ (Bal.)
$\varepsilon = 0.125$ (Bal.)
$\varepsilon = 0.0$ (Bal.)
Emp. Measure