

Statistical Foundations of Foundation Modeling

Final Examination
June 02, 2025

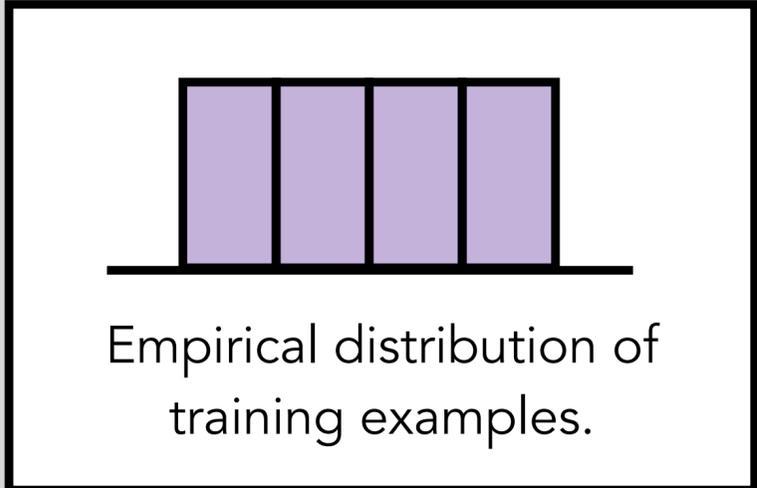


Ronak Mehta





Data Generating
Distribution
 P

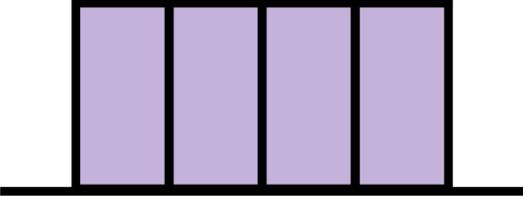




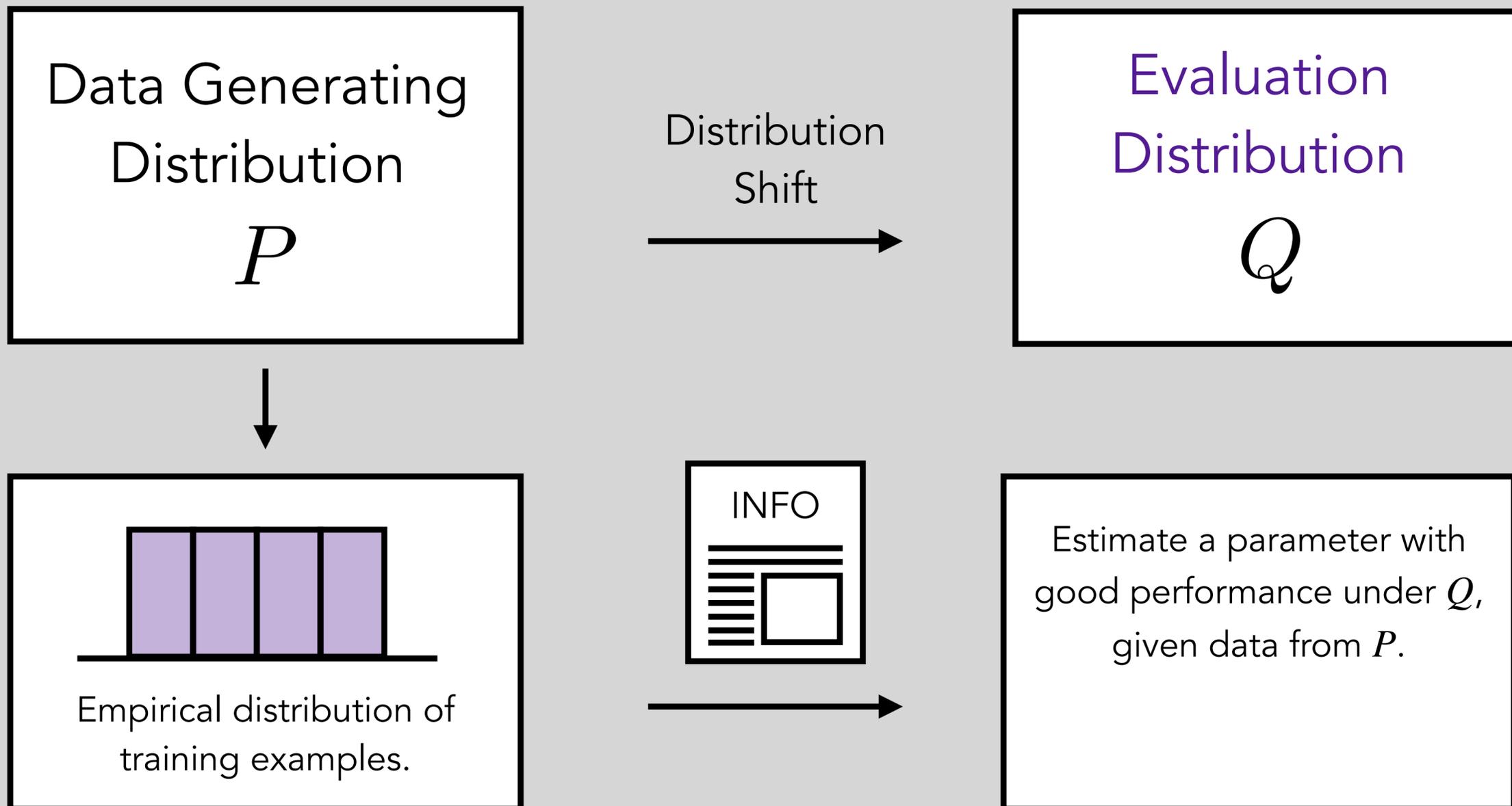
Data Generating
Distribution
 P

Distribution
Shift
→

Evaluation
Distribution
 Q



Empirical distribution of
training examples.



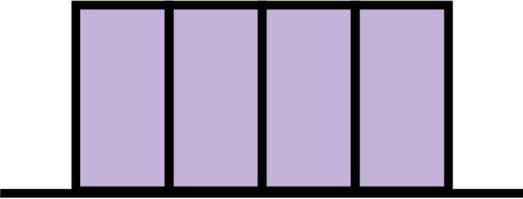


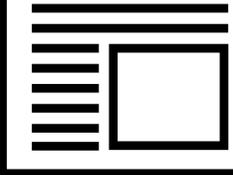
Data Generating
Distribution
 P

Distribution
Shift
→

Evaluation
Distribution
 Q

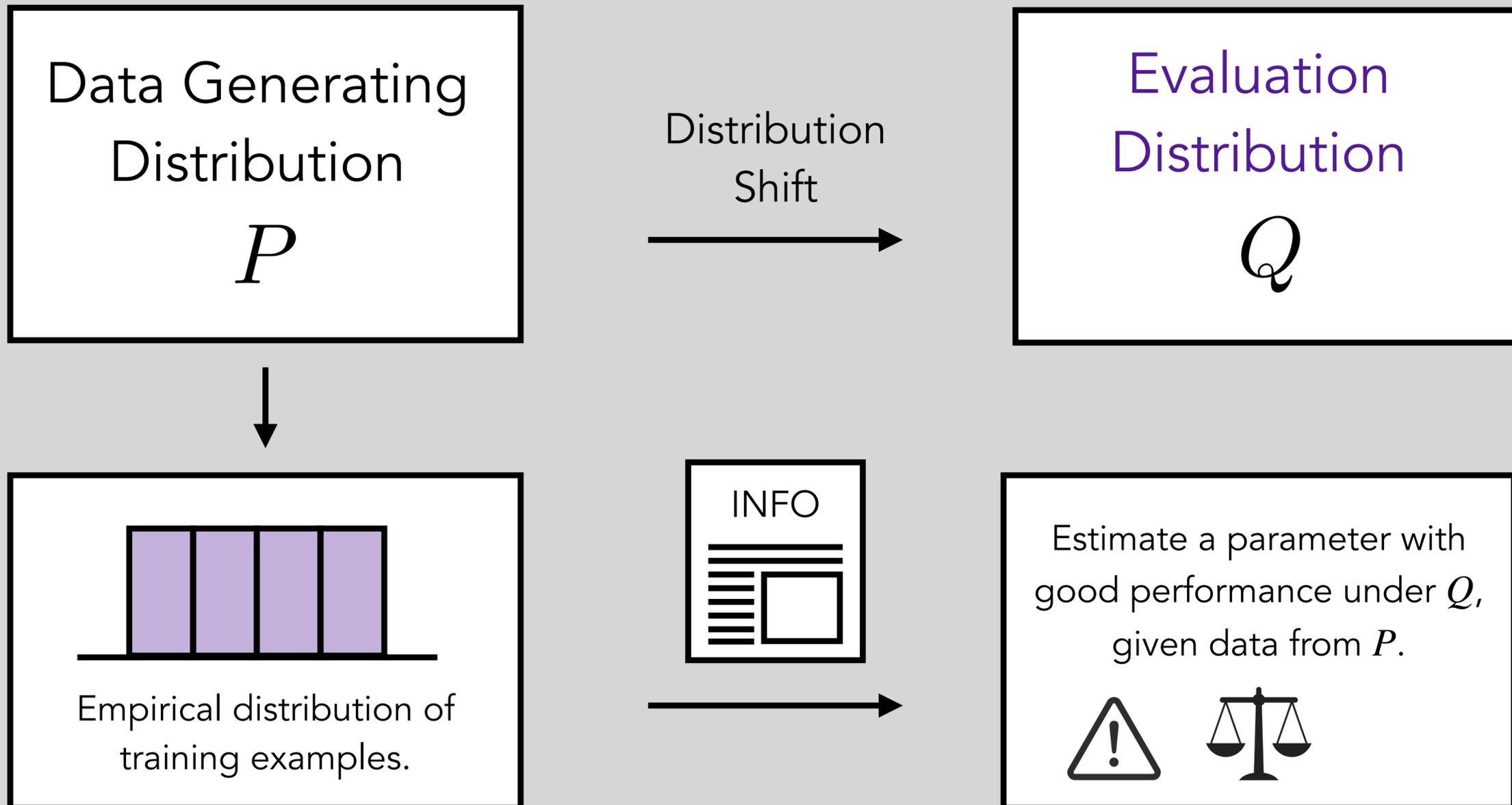


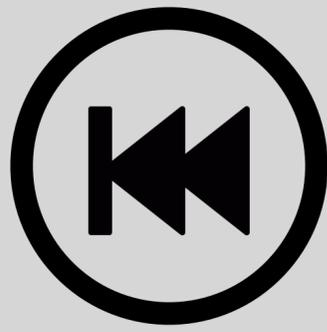

Empirical distribution of
training examples.

INFO




Estimate a parameter with
good performance under Q ,
given data from P .

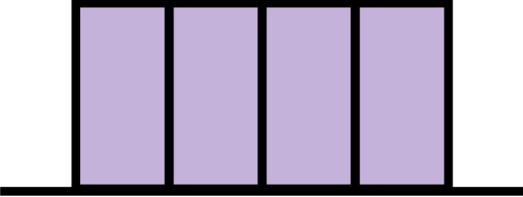


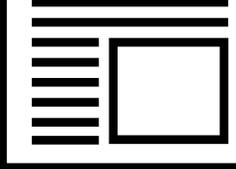
Data Generating
Distribution
 P

Distribution
Shift
→

Evaluation
Distribution
 Q




Empirical distribution of
training examples.

INFO




Estimate a parameter with
good performance under Q ,
given data from P .
  



$$\min_{\mathbf{x} \in \mathcal{X}} \max_{\mathbf{y} \in \mathcal{Y}} \langle \mathbf{y}, f(\mathbf{x}) \rangle - \psi(\mathbf{y}) + \phi(\mathbf{x})$$



$$\min_{\mathbf{x} \in \mathcal{X}} \max_{\mathbf{y} \in \mathcal{Y}} \underbrace{\langle \mathbf{y}, f(\mathbf{x}) \rangle}_{\text{Loss}} - \psi(\mathbf{y}) + \phi(\mathbf{x})$$

General Exam: Distributionally Robust Optimization

$$\mathcal{X} = \mathbb{R}^d$$

Model parameters

$$\mathcal{Y} = \Delta^{n-1}$$

Weights on training examples

$$f(\mathbf{x}) = (f_1(\mathbf{x}), \dots, f_n(\mathbf{x}))$$

Losses on training examples



$$\min_{\mathbf{x} \in \mathcal{X}} \max_{\mathbf{y} \in \mathcal{Y}} \underbrace{\langle \mathbf{y}, f(\mathbf{x}) \rangle - \psi(\mathbf{y}) + \phi(\mathbf{x})}$$

General Exam: Distributionally Robust Optimization

$$\mathcal{X} = \mathbb{R}^d$$

Model parameters

$$\mathcal{Y} = \Delta^{n-1}$$

Weights on training examples

$$f(\mathbf{x}) = (f_1(\mathbf{x}), \dots, f_n(\mathbf{x}))$$

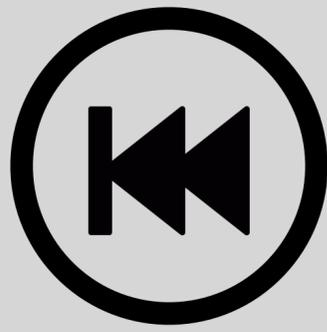
Losses on training examples

$$\phi(\mathbf{x}) = \frac{\mu}{2} \|\mathbf{x}\|_2^2$$

Regularizer

$$\psi(\mathbf{y}) = D_{\text{KL}}(\mathbf{y} \parallel \mathbf{1}/n)$$

Shift penalty



$$\min_{\mathbf{x} \in \mathcal{X}} \max_{\mathbf{y} \in \mathcal{Y}} \langle \mathbf{y}, f(\mathbf{x}) \rangle - \psi(\mathbf{y}) + \phi(\mathbf{x})$$

Example: Minimization with Functional Constraints



$$\min_{\mathbf{x} \in \mathcal{X}} \max_{\mathbf{y} \in \mathcal{Y}} \langle \mathbf{y}, f(\mathbf{x}) \rangle - \psi(\mathbf{y}) + \phi(\mathbf{x})$$

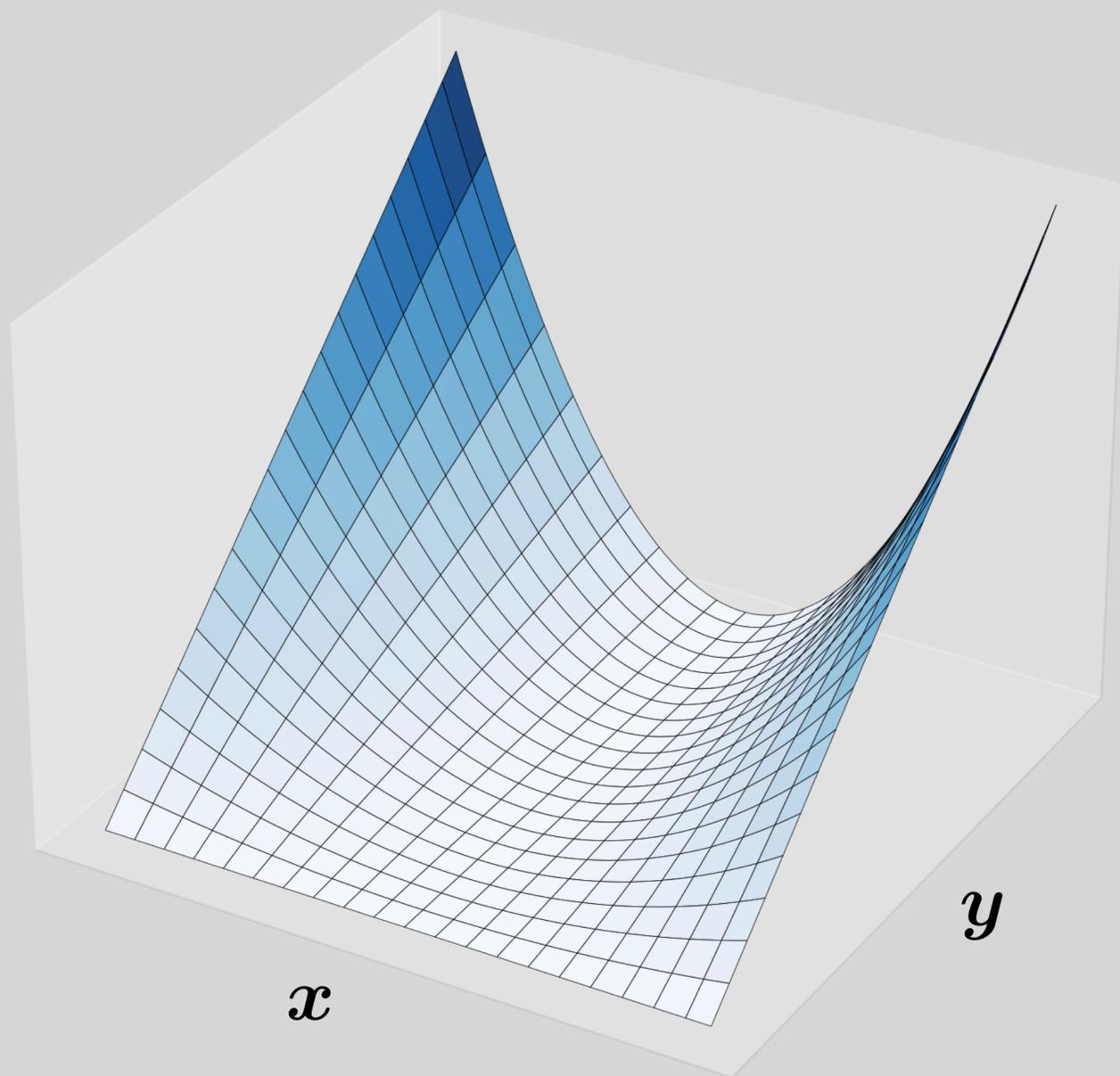
Example: Minimization with Functional Constraints

$$\min_{\mathbf{x} \in \mathcal{X}} \phi(\mathbf{x})$$

$$\text{s. t. } f_i(\mathbf{x}) \leq 0, \quad i = 1, \dots, n$$



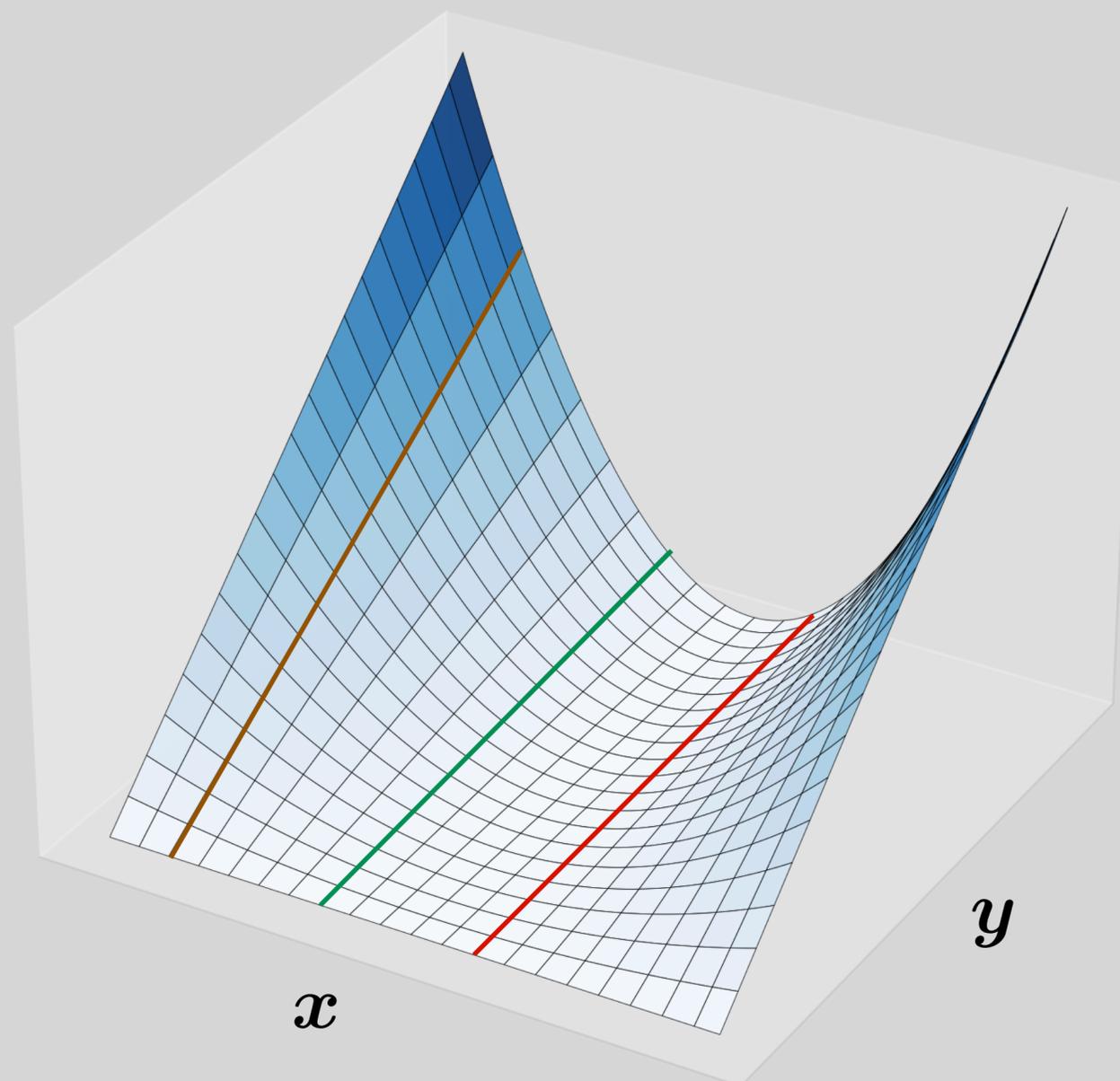
$$\min_{\mathbf{x} \in \mathcal{X}} \max_{\mathbf{y} \in \mathcal{Y}} \underbrace{\langle \mathbf{y}, f(\mathbf{x}) \rangle}_{\text{}} - \psi(\mathbf{y}) + \phi(\mathbf{x})$$





Semilinear Min-Max Problem

$$\min_{\mathbf{x} \in \mathcal{X}} \max_{\mathbf{y} \in \mathcal{Y}} \underbrace{\langle \mathbf{y}, f(\mathbf{x}) \rangle}_{\text{inner product}} - \psi(\mathbf{y}) + \phi(\mathbf{x})$$





Semilinear Min-Max Problem

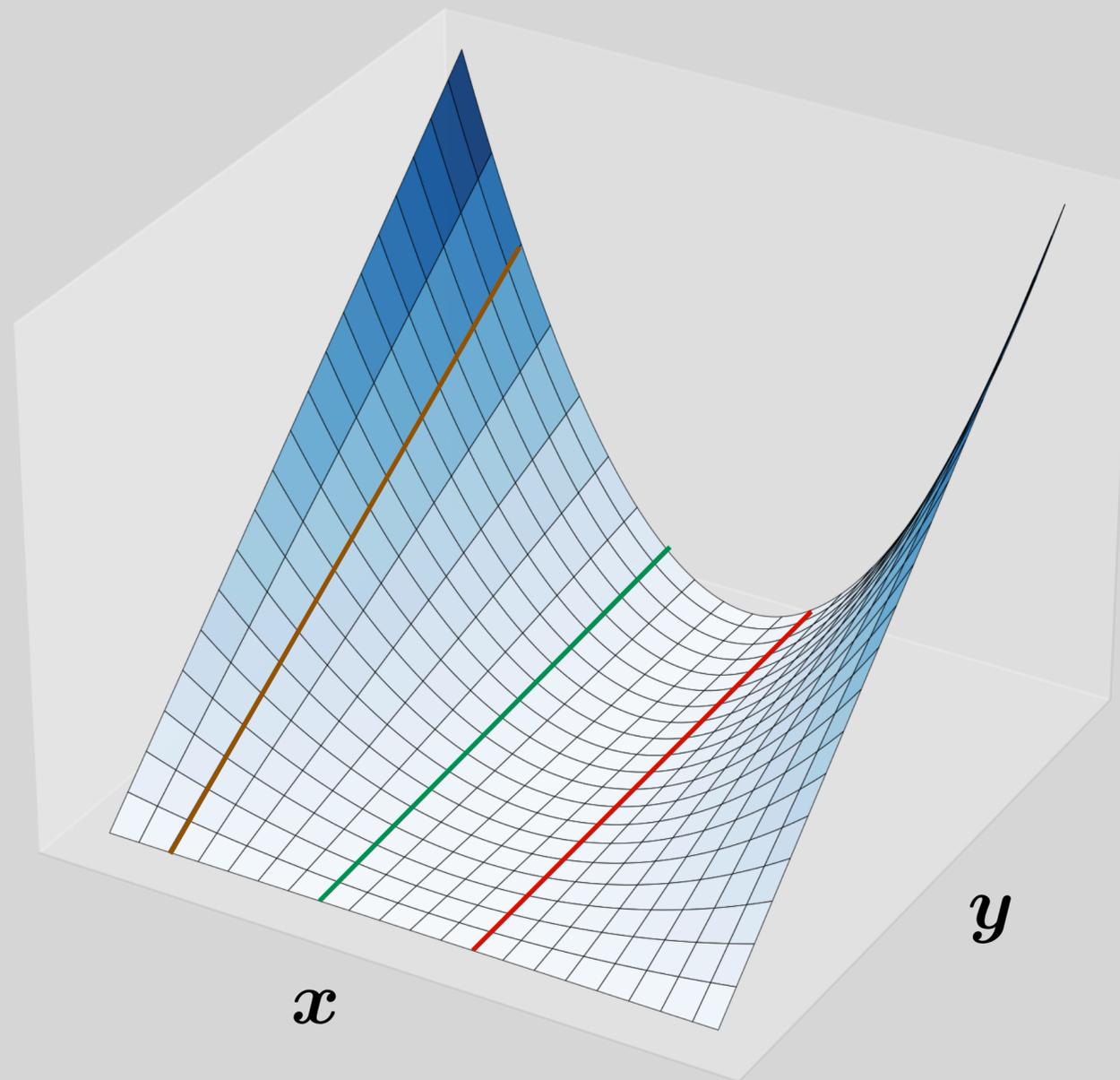
$$\langle \mathbf{y}, \mathbf{Ax} \rangle \subseteq$$

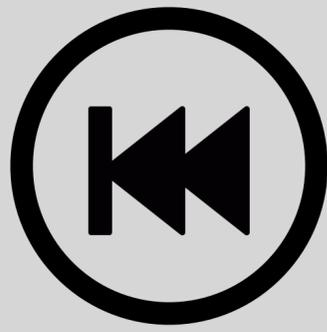
Bilinear

$$\min_{\mathbf{x} \in \mathcal{X}} \max_{\mathbf{y} \in \mathcal{Y}} \underbrace{\langle \mathbf{y}, f(\mathbf{x}) \rangle}_{\subseteq g(\mathbf{x}, \mathbf{y})} - \psi(\mathbf{y}) + \phi(\mathbf{x})$$

$$\subseteq g(\mathbf{x}, \mathbf{y})$$

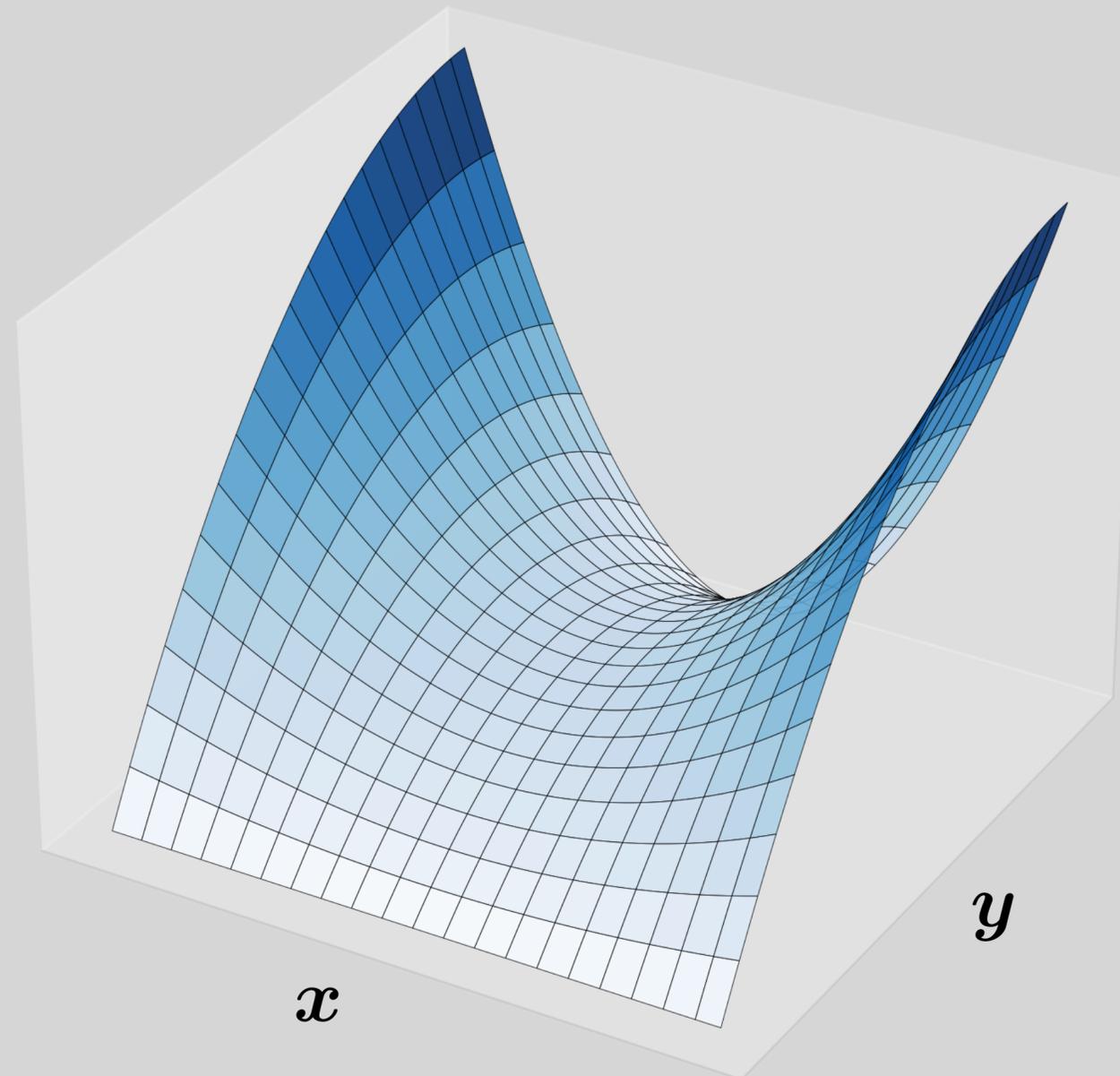
Nonbilinear

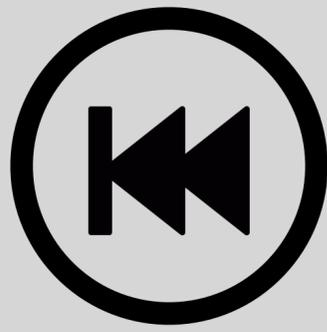




Semilinear Min-Max Problem

$$\min_{\mathbf{x} \in \mathcal{X}} \max_{\mathbf{y} \in \mathcal{Y}} \underbrace{\langle \mathbf{y}, f(\mathbf{x}) \rangle - \psi(\mathbf{y}) + \phi(\mathbf{x})}$$





Semilinear Min-Max Problem

$$\min_{\mathbf{x} \in \mathcal{X}} \max_{\mathbf{y} \in \mathcal{Y}} \langle \mathbf{y}, f(\mathbf{x}) \rangle - \psi(\mathbf{y}) + \phi(\mathbf{x})$$

What do we hope to achieve?



Semilinear Min-Max Problem

$$\min_{\mathbf{x} \in \mathcal{X}} \max_{\mathbf{y} \in \mathcal{Y}} \langle \mathbf{y}, f(\mathbf{x}) \rangle - \psi(\mathbf{y}) + \phi(\mathbf{x})$$

What do we hope to achieve?

Functions	Lipschitz	Smooth
f_1		
\vdots		
f_n		



Semilinear Min-Max Problem

$$\min_{\mathbf{x} \in \mathcal{X}} \max_{\mathbf{y} \in \mathcal{Y}} \langle \mathbf{y}, f(\mathbf{x}) \rangle - \psi(\mathbf{y}) + \phi(\mathbf{x})$$

What do we hope to achieve?

Functions	Lipschitz	Smooth
f_1	G_1	
\vdots	\vdots	
f_n	G_n	

$$|f_i(\mathbf{x}) - f_i(\mathbf{x}')| \leq G_i \|\mathbf{x} - \mathbf{x}'\|_{\mathcal{X}}$$



Semilinear Min-Max Problem

$$\min_{\mathbf{x} \in \mathcal{X}} \max_{\mathbf{y} \in \mathcal{Y}} \langle \mathbf{y}, f(\mathbf{x}) \rangle - \psi(\mathbf{y}) + \phi(\mathbf{x})$$

What do we hope to achieve?

Functions	Lipschitz	Smooth
f_1	G_1	L_1
\vdots	\vdots	\vdots
f_n	G_n	L_n

$$|f_i(\mathbf{x}) - f_i(\mathbf{x}')| \leq G_i \|\mathbf{x} - \mathbf{x}'\|_x$$

$$\|\nabla f_i(\mathbf{x}) - \nabla f_i(\mathbf{x}')\|_{x^*} \leq L_i \|\mathbf{x} - \mathbf{x}'\|_x$$



Semilinear Min-Max Problem

$$\min_{\mathbf{x} \in \mathcal{X}} \max_{\mathbf{y} \in \mathcal{Y}} \langle \mathbf{y}, f(\mathbf{x}) \rangle - \psi(\mathbf{y}) + \phi(\mathbf{x})$$

What do we hope to achieve?

Functions	Lipschitz	Smooth
f_1	G_1	L_1
\vdots	\vdots	\vdots
f_n	G_n	L_n

$$|f_i(\mathbf{x}) - f_i(\mathbf{x}')| \leq G_i \|\mathbf{x} - \mathbf{x}'\|_x$$

$$\|\nabla f_i(\mathbf{x}) - \nabla f_i(\mathbf{x}')\|_{x^*} \leq L_i \|\mathbf{x} - \mathbf{x}'\|_x$$

1. Can we use **arbitrary norms** (not only the Euclidean) to define (G_i, L_i) ?
2. Can we achieve the best **aggregation** of constants in the final complexity?

$$nG_{\max} \geq \sqrt{n \sum_{i=1}^n G_i^2} \geq \sum_{i=1}^n G_i$$



Semilinear Min-Max Problem

$$\min_{\mathbf{x} \in \mathcal{X}} \max_{\mathbf{y} \in \mathcal{Y}} \langle \mathbf{y}, f(\mathbf{x}) \rangle - \psi(\mathbf{y}) + \phi(\mathbf{x})$$

**Global Complexity Guarantee
(Convex/Concave):**

lower is better
↓

$$O\left(\frac{d\sqrt{(\sum_{i=1}^n \lambda_i) \cdot (\sum_{i=1}^n G_i^2 + \sum_{i=1}^n L_i^2)}}{\epsilon}\right)$$

$$\lambda_i = \sqrt{G_i^2 + L_i^2}$$

**Extragradient (Korpelevich '76)/
Mirror Prox (Nemirovski, '04) Baseline:**

$$O\left(\frac{nd\sqrt{G^2 + L^2}}{\epsilon}\right)$$

(Lipschitz/Smoothness constants of
the vector-valued function f).

Jelena Diakonikolas
University of Wisconsin-Madison



Zaid Harchaoui
University of Washington



Semilinear Min-Max Problem

$$\min_{\mathbf{x} \in \mathcal{X}} \max_{\mathbf{y} \in \mathcal{Y}} \langle \mathbf{y}, f(\mathbf{x}) \rangle - \psi(\mathbf{y}) + \phi(\mathbf{x})$$

Dual Linear Min-Max Optimization

Ronak Mehta¹ Jelena Diakonikolas² Zaid Harchaoui¹

¹University of Washington, Seattle ²University of Wisconsin, Madison

June 2, 2025

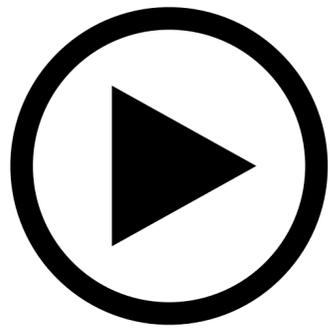
Abstract

We study a class of convex-concave min-max problems in which the coupled component of the objective is linear in at least one of the two decision vectors. This problem structure interpolates between the well-studied bilinearly coupled and relatively less-studied nonbilinearly coupled min-max problems. Relevant applications include distributionally robust optimization, fully composite optimization, and convex minimization with functional constraints. We prove the first complexity upper bounds for this class using an algorithm that combines randomized and cyclic coordinate-wise updates, which have previously been employed to achieve optimal complexities for bilinearly-coupled min-max problems. The analysis handles strongly convex and non-strongly convex components in a unified manner. We also provide lower bounds for particular parameter regimes of practical interest.

Jelena Diakonikolas
University of Wisconsin-Madison



Zaid Harchaoui
University of Washington



Data Generating
Distribution

P

Evaluation
Distribution

Q

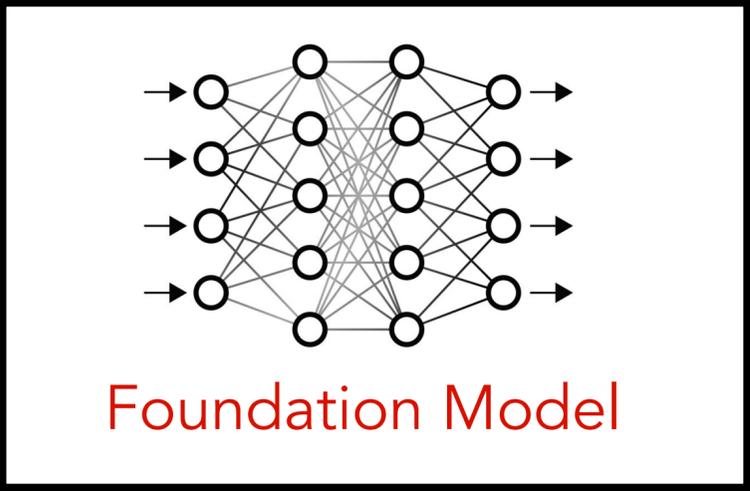
Data Generating
Distribution
 P

Setting, data,
and task all
change!



Evaluation
Distribution
 Q

Data Generating
Distribution
 P

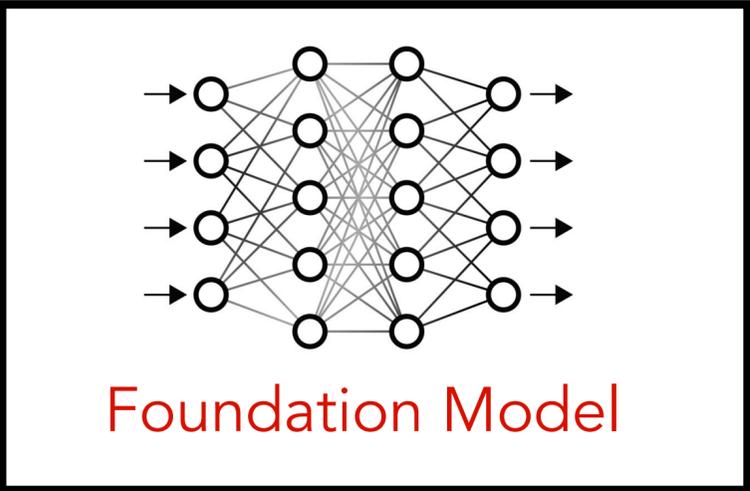


Setting, data,
and task all
change!



Evaluation
Distribution
 Q

Data Generating
Distribution
 P



Setting, data,
and task all
change!



Evaluation
Distribution
 Q



Generalize to unseen
tasks, *without* additional
training data.

The Mystery of Foundation Models

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 Ilya Sutskever¹

Self-Supervised Learning from Images with a Joint-Embedding Predictive Architecture

Mahmoud Assran^{1,2,3*} Quentin Duval¹ Ishan Misra¹ Piotr Bojanowski¹
Pascal Vincent¹ Michael Rabbat^{1,3} Yann LeCun^{1,4} Nicolas Ballas¹

¹Meta AI (FAIR) ²McGill University ³Mila, Quebec AI Institute ⁴New York University

GPT-4 Technical Report

DINOv2: Learning Robust Visual Features without Supervision

Maxime Oquab**, Timothée Darcet**, Théo Moutakanni**,
Marc Szafraniec*, Vasil Khalidov*, Pierre Fernandez, Daniel Haziza,
Alaaeldin El-Nouby, Mahmoud Assran, Nicolas Ballas, Wojciech Galuba,
Arman, Po-Yao Huang, Shang-Wen Li, Ishan Misra, Michael Rabbat,
Arman, Gabriel Synnaeve, Hu Xu, Hervé Jegou, Julien Mairal¹,
Patrick Labatut*, Armand Joulin*, Piotr Bojanowski*

Meta AI Research ¹Inria
*core team **equal contribution



DeepSeek-R1: Incentivizing Reasoning Capabilities with Reinforcement Learning

DeepSeek-AI

research@deepseek.com

GPT-4.



The Llama 3 Herd of Models

Llama Team, AI @ Meta¹

¹A detailed contributor list can be found in the appendix of this paper.

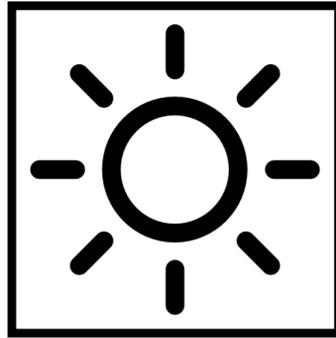
Modern artificial intelligence (AI) systems are powered by foundation models. This paper presents a new set of foundation models, called Llama 3. It is a herd of language models that natively support multilinguality, coding, reasoning, and tool usage. Our largest model is a dense Transformer with 405B parameters and a context window of up to 128K tokens. This paper presents an extensive empirical evaluation of Llama 3. We find that Llama 3 delivers comparable quality to leading language models such as GPT-4 on a plethora of tasks. We publicly release Llama 3, including pre-trained and post-trained versions of the 405B parameter language model and our Llama Guard 3 model for input and output safety. The paper also presents the results of experiments in which we integrate image, video, and speech capabilities into Llama 3 via a compositional approach. We observe this approach performs competitively with the state-of-the-art on image, video, and speech recognition tasks. The resulting models are not yet being broadly released as they are still under development.

Date: July 23, 2024

Website: <https://llama.meta.com/>

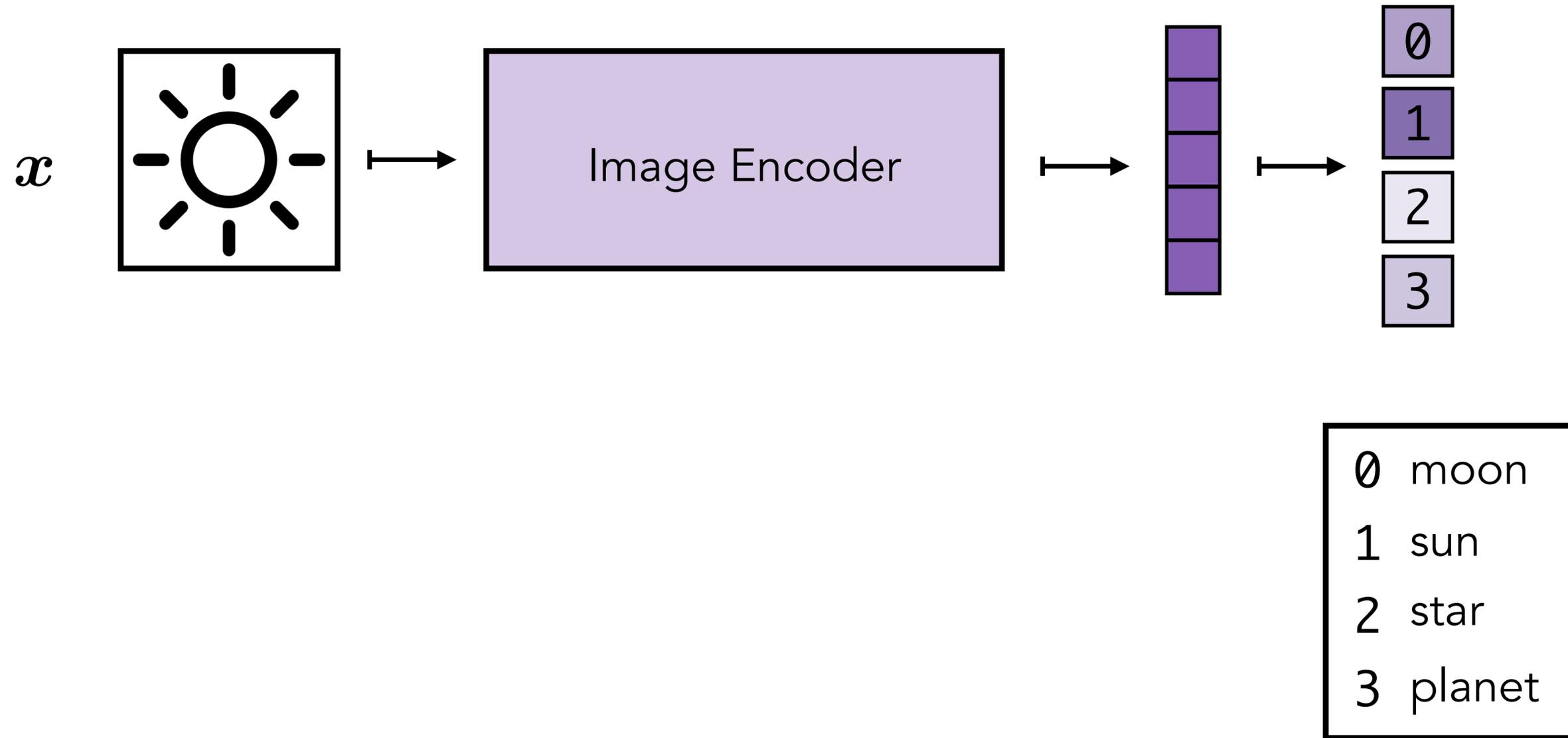
Supervised Learning

x



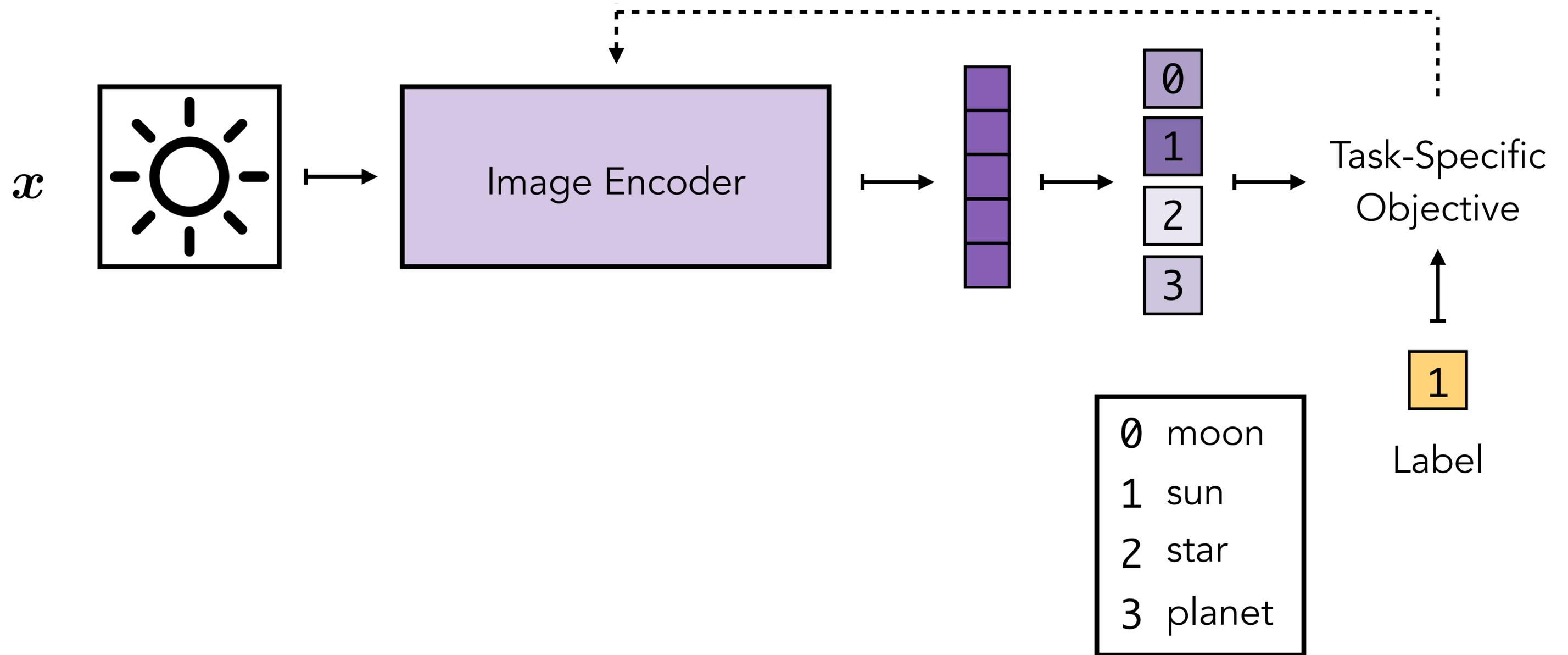
- 0 moon
- 1 sun
- 2 star
- 3 planet

Supervised Learning



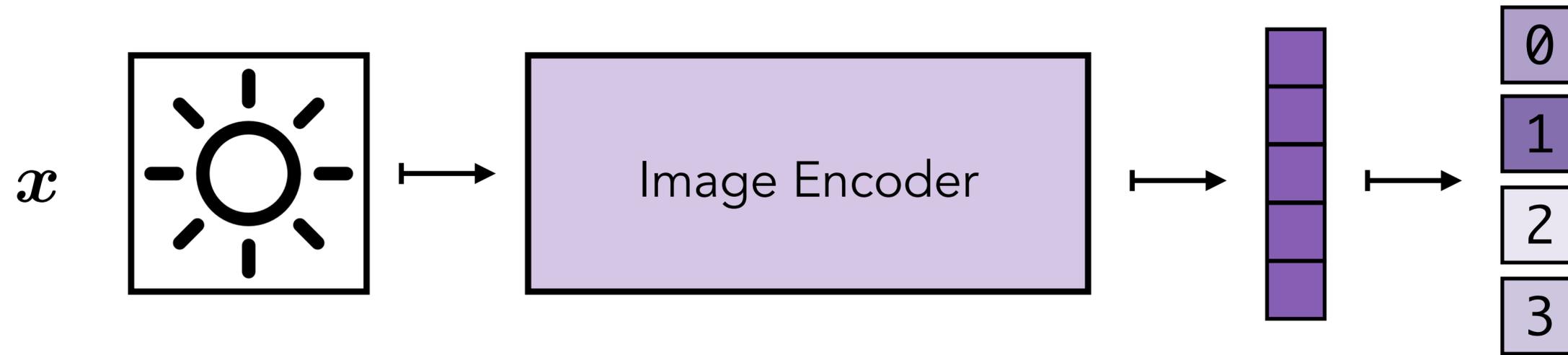
Training

Supervised Learning



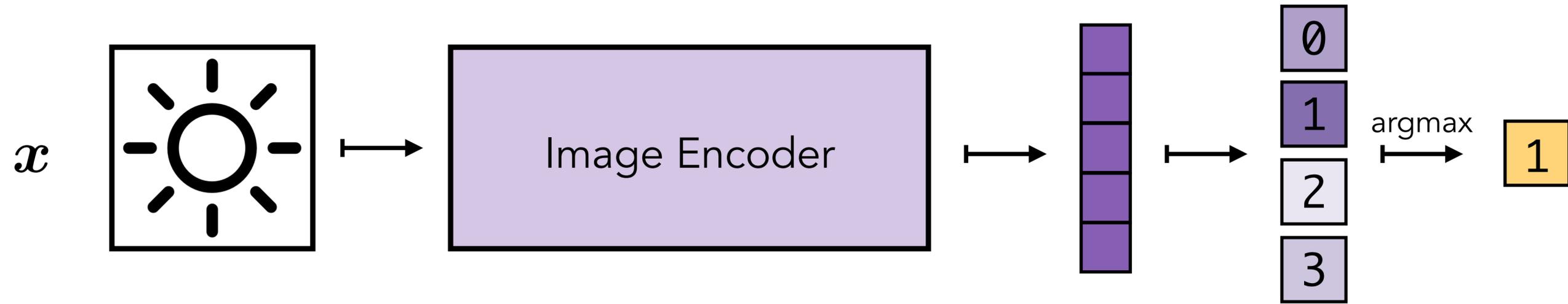
Training

Supervised Learning

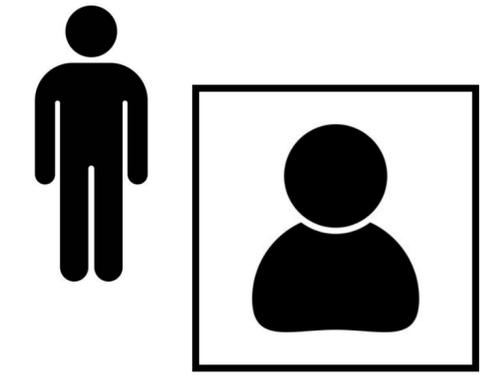
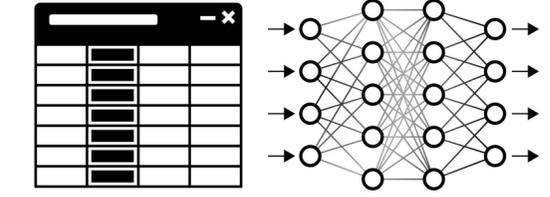
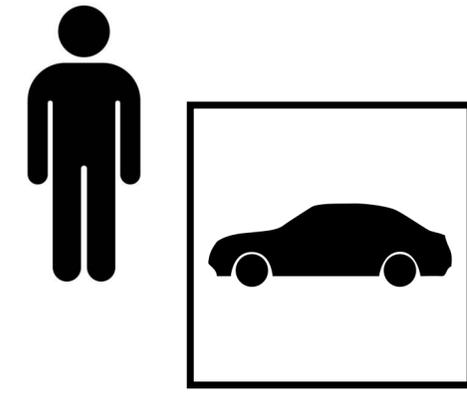
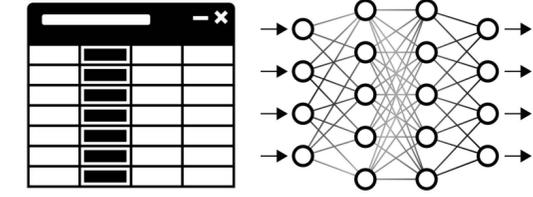
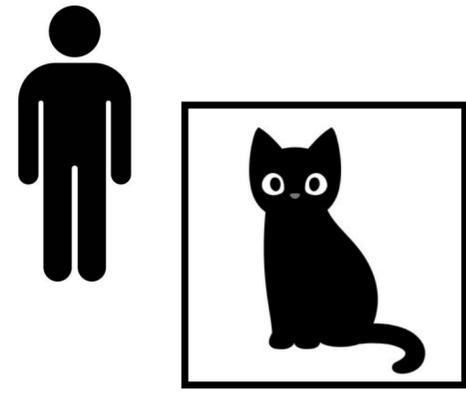
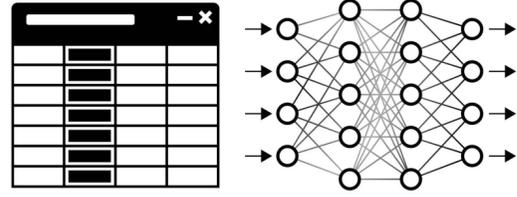


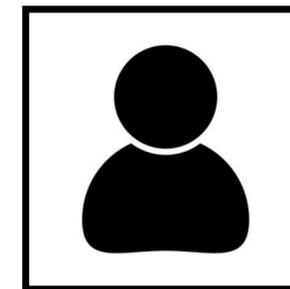
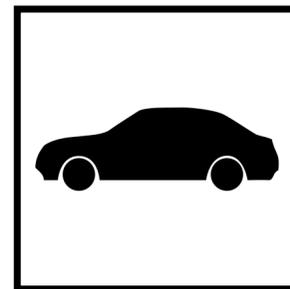
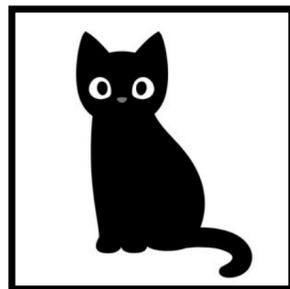
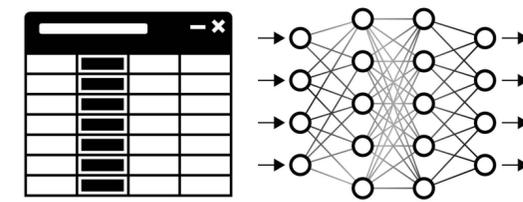
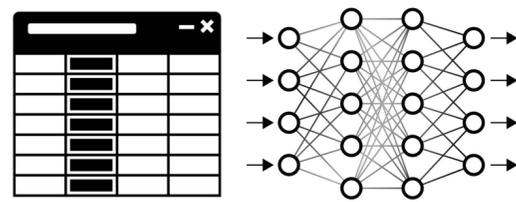
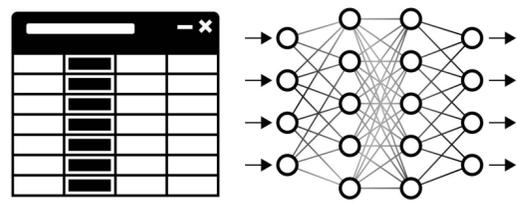
Evaluation

Supervised Learning

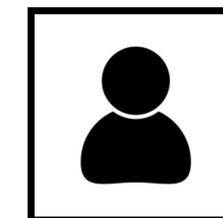
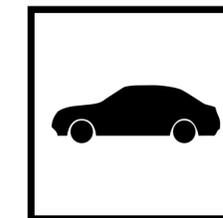
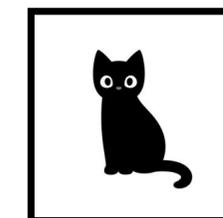
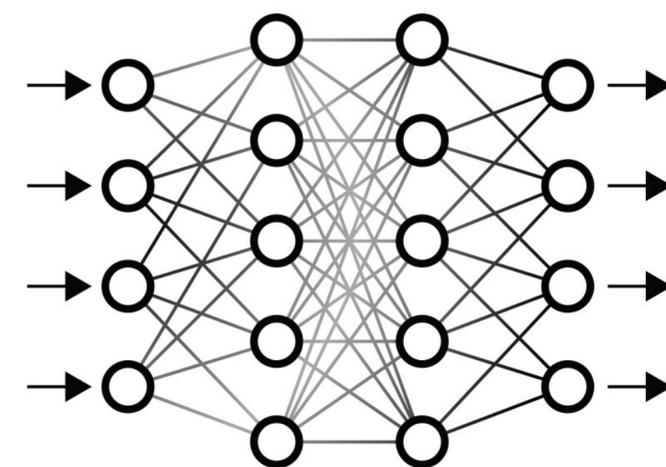
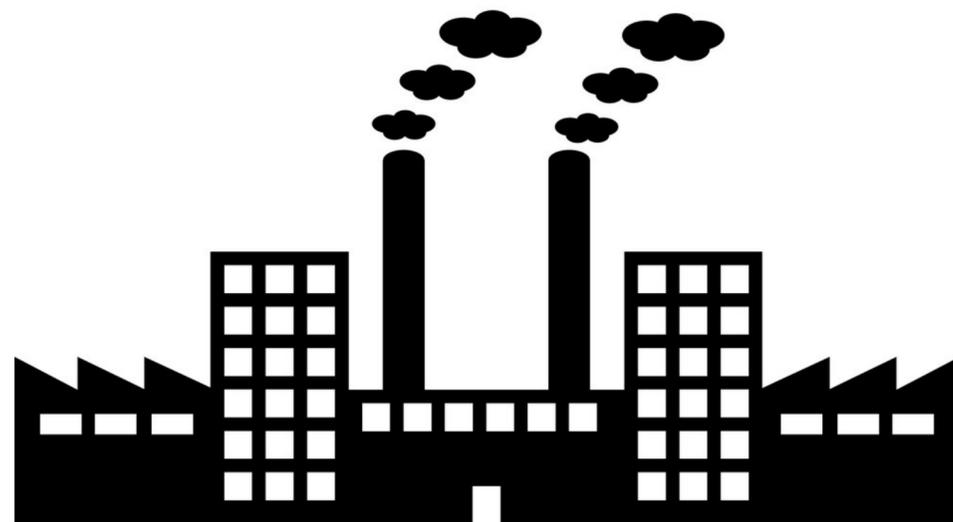
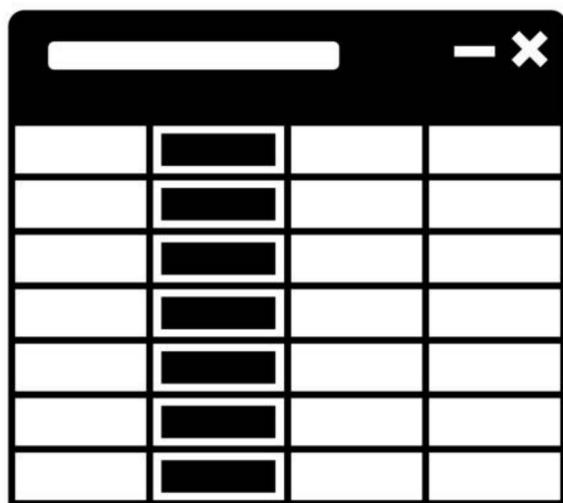


Evaluation

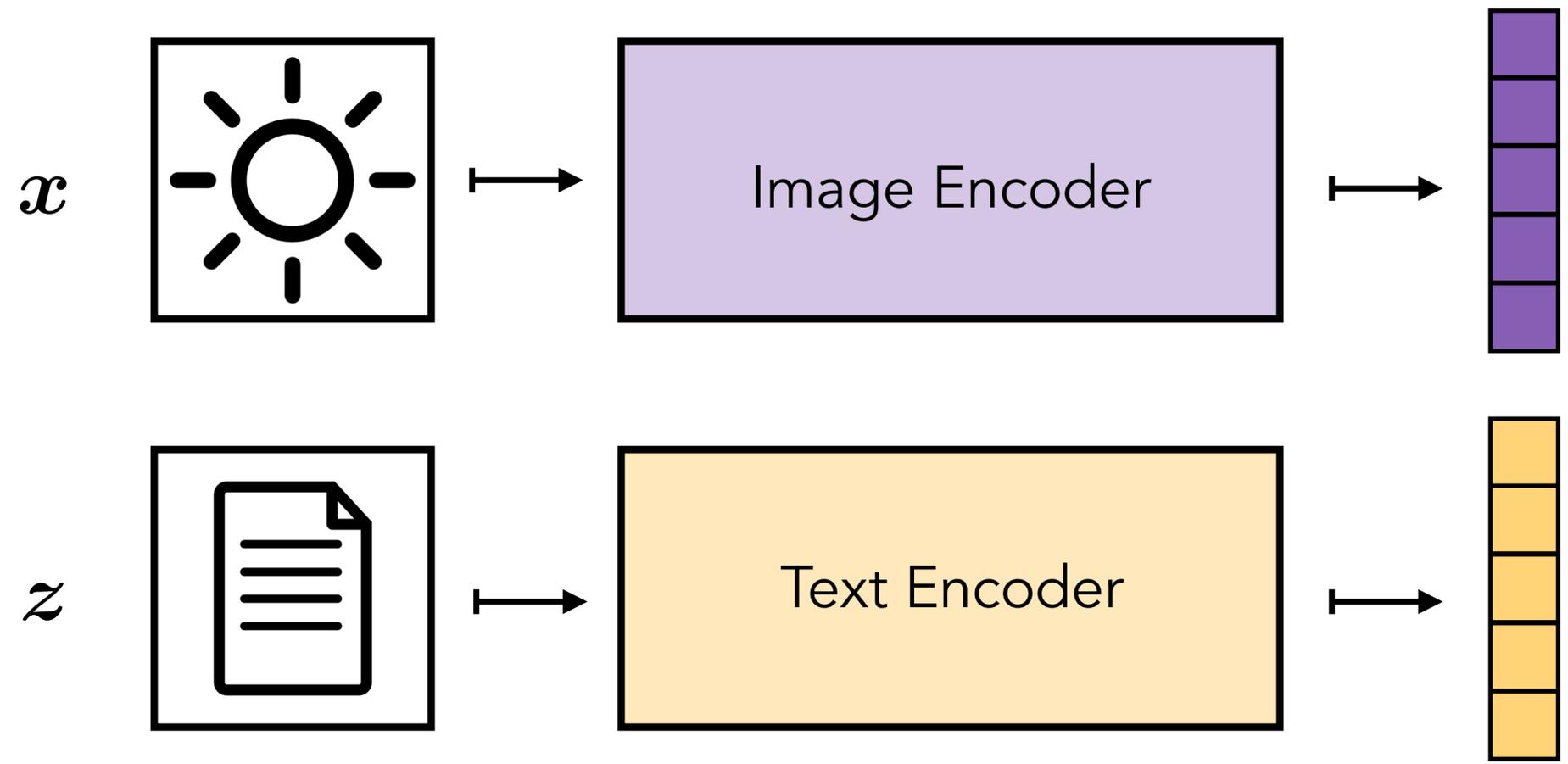




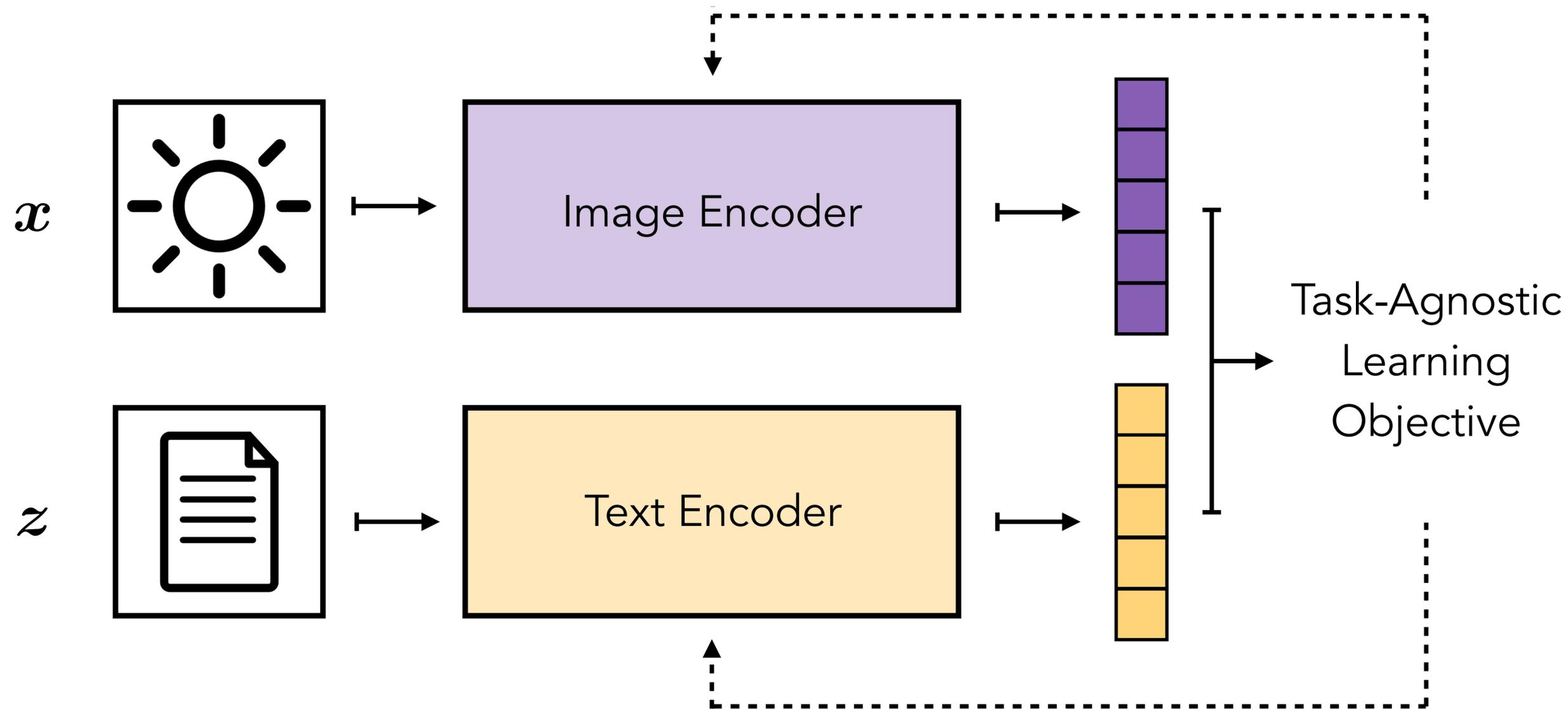
Foundation Modeling Pipeline



Foundation Modeling

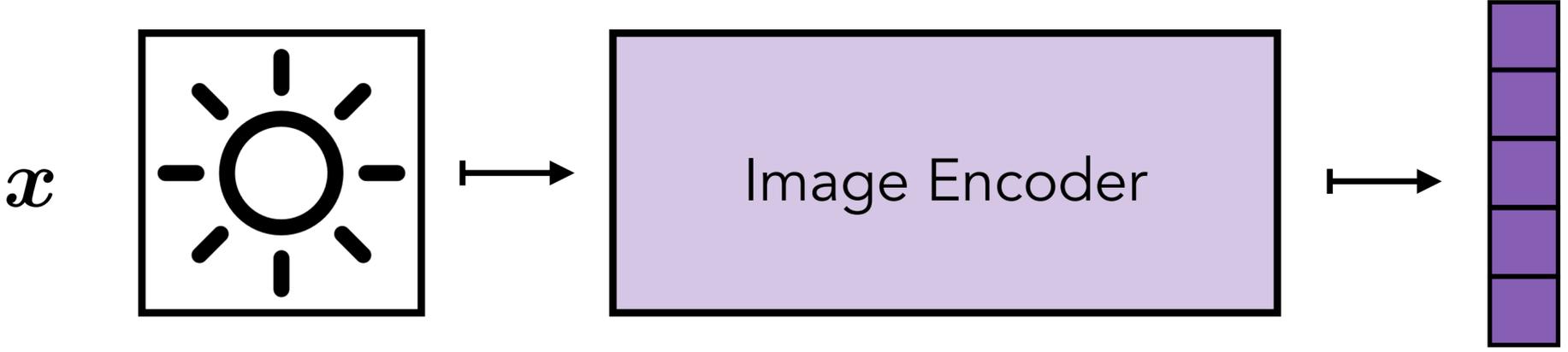


Foundation Modeling



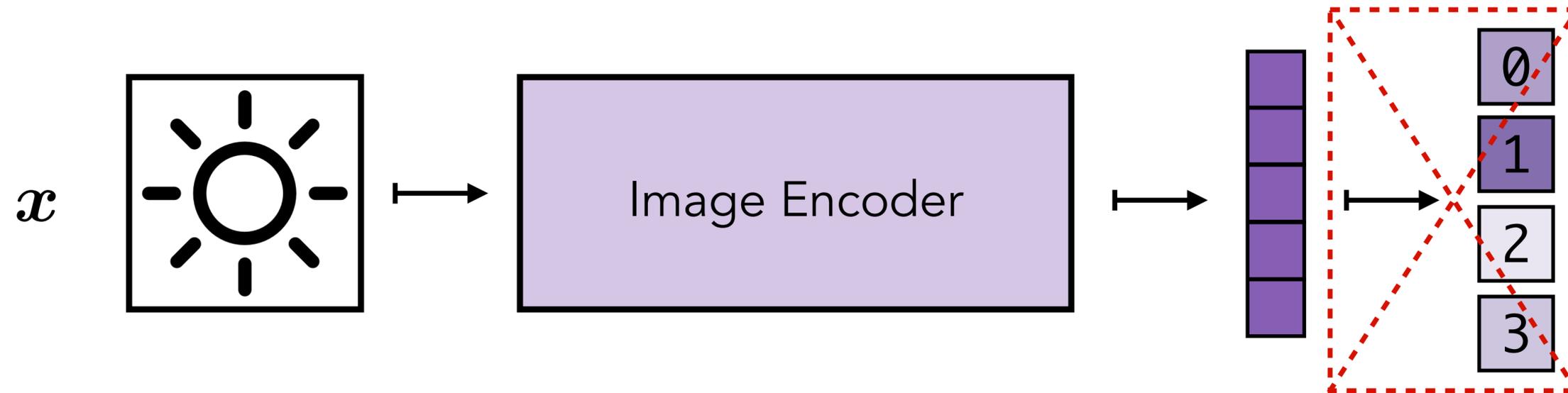
Pre-Training

Foundation Modeling



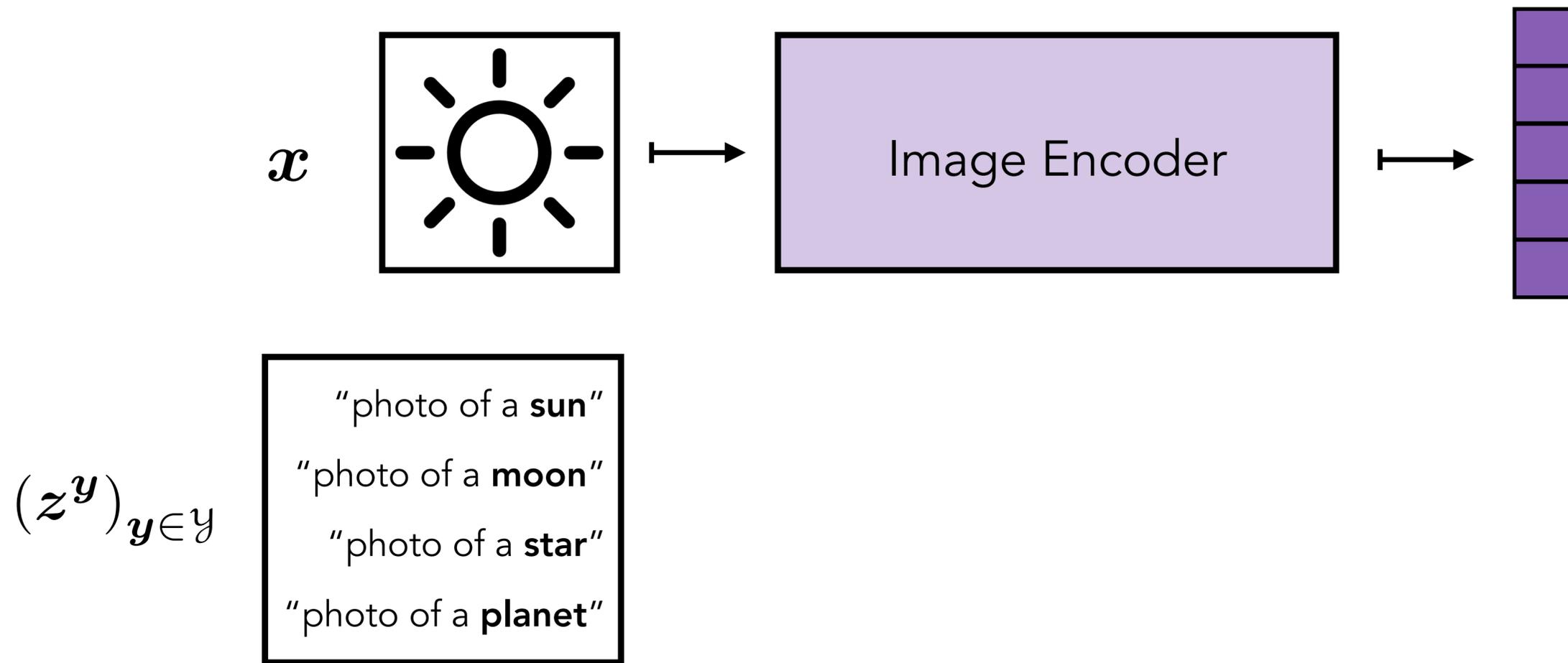
Evaluation

Foundation Modeling



Evaluation

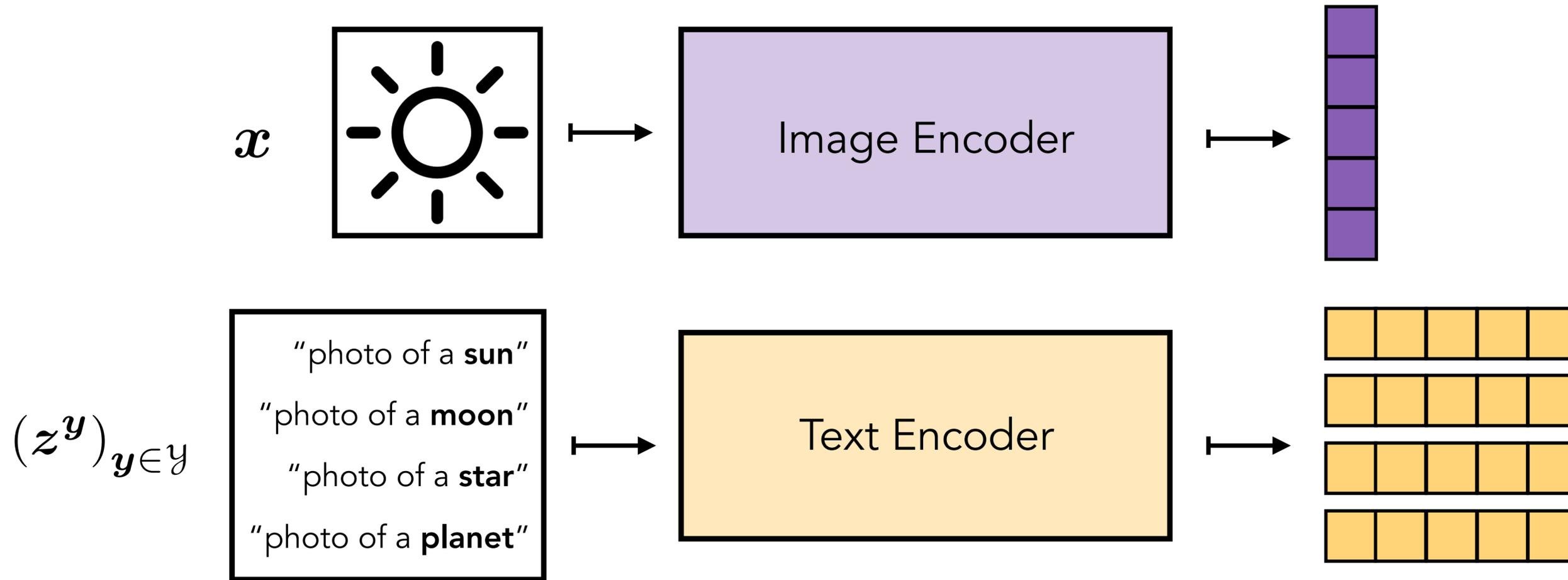
Foundation Modeling



Idea: Convert labels into pseudo-captions (prompts)

Evaluation

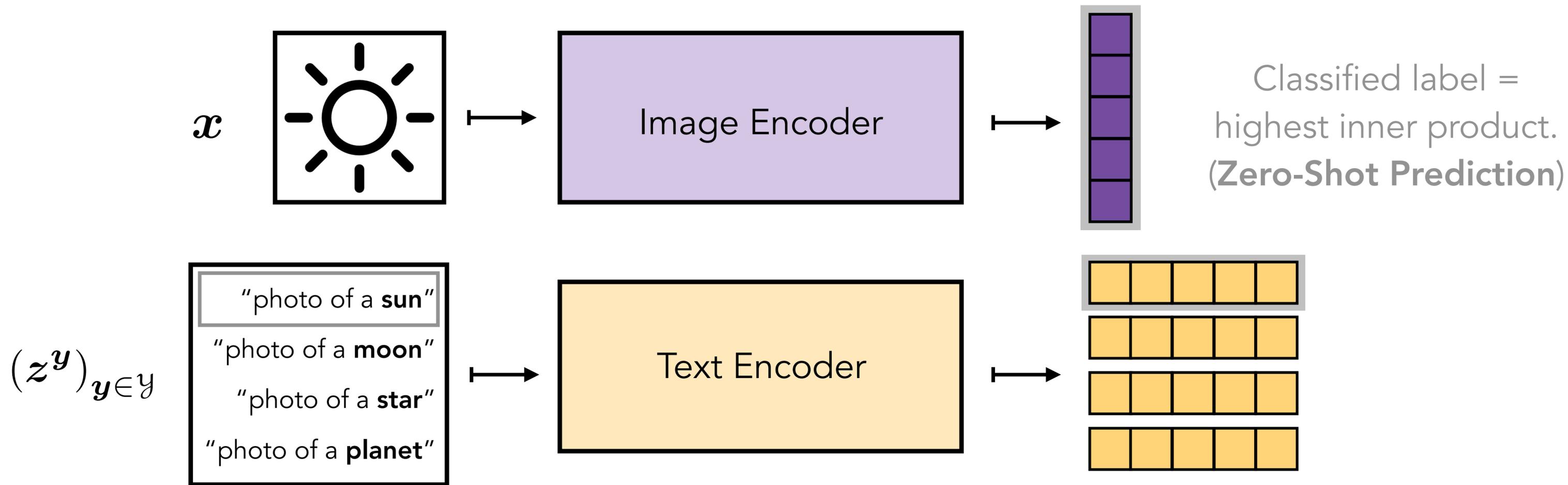
Foundation Modeling



Idea: Convert labels into pseudo-captions (prompts)

Evaluation

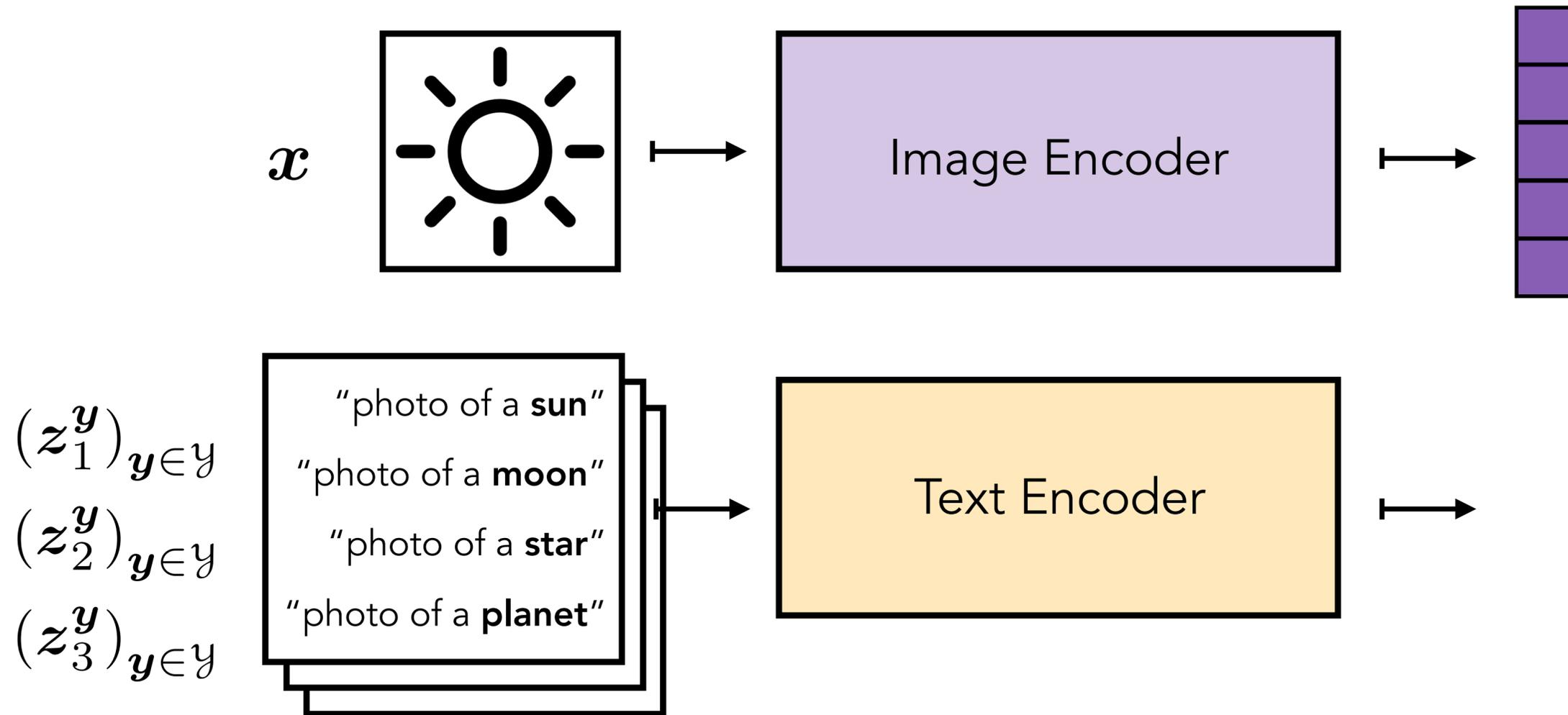
Foundation Modeling



Idea: Convert labels into pseudo-captions (prompts)

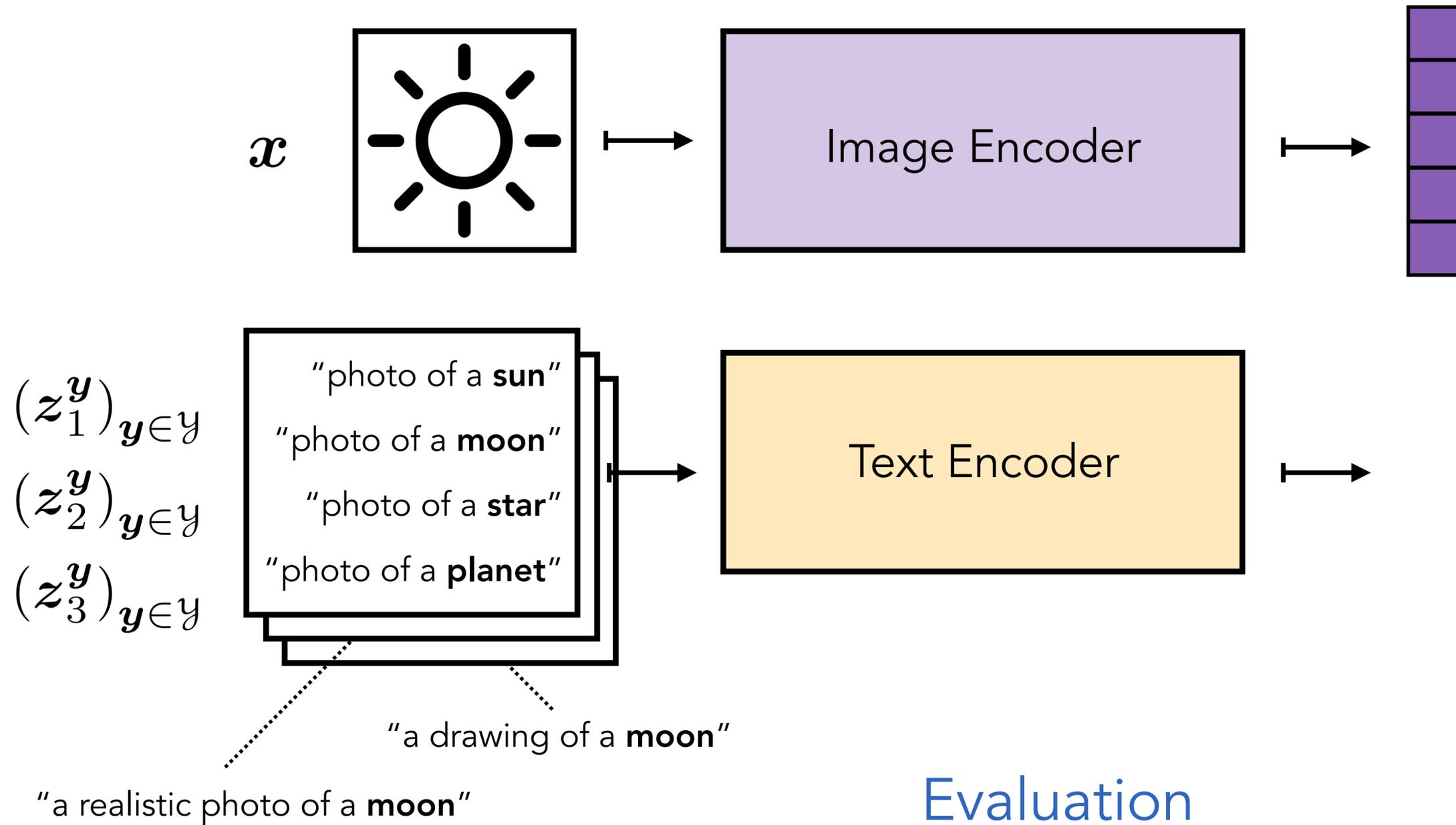
Evaluation

Foundation Modeling

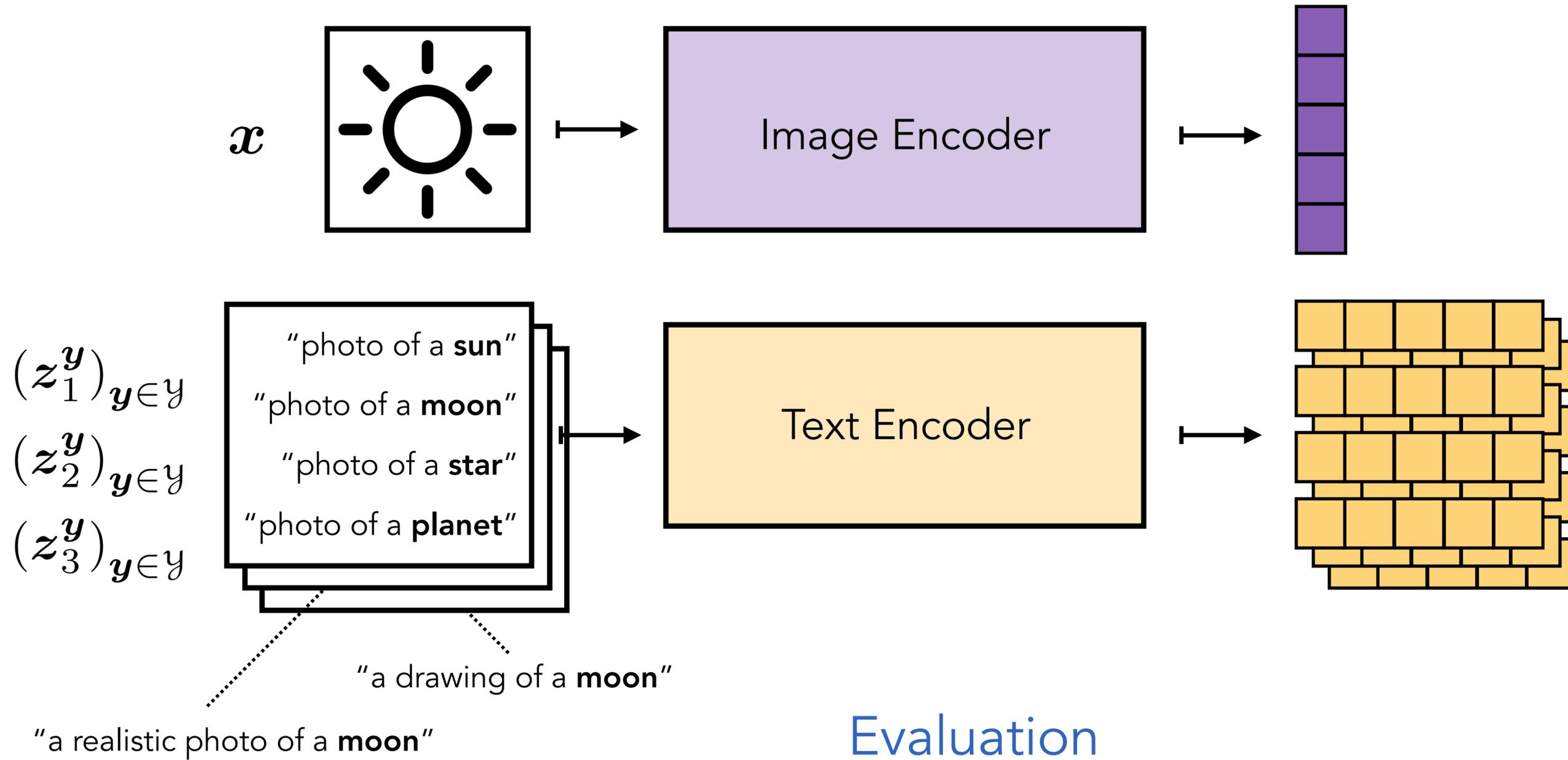


Evaluation

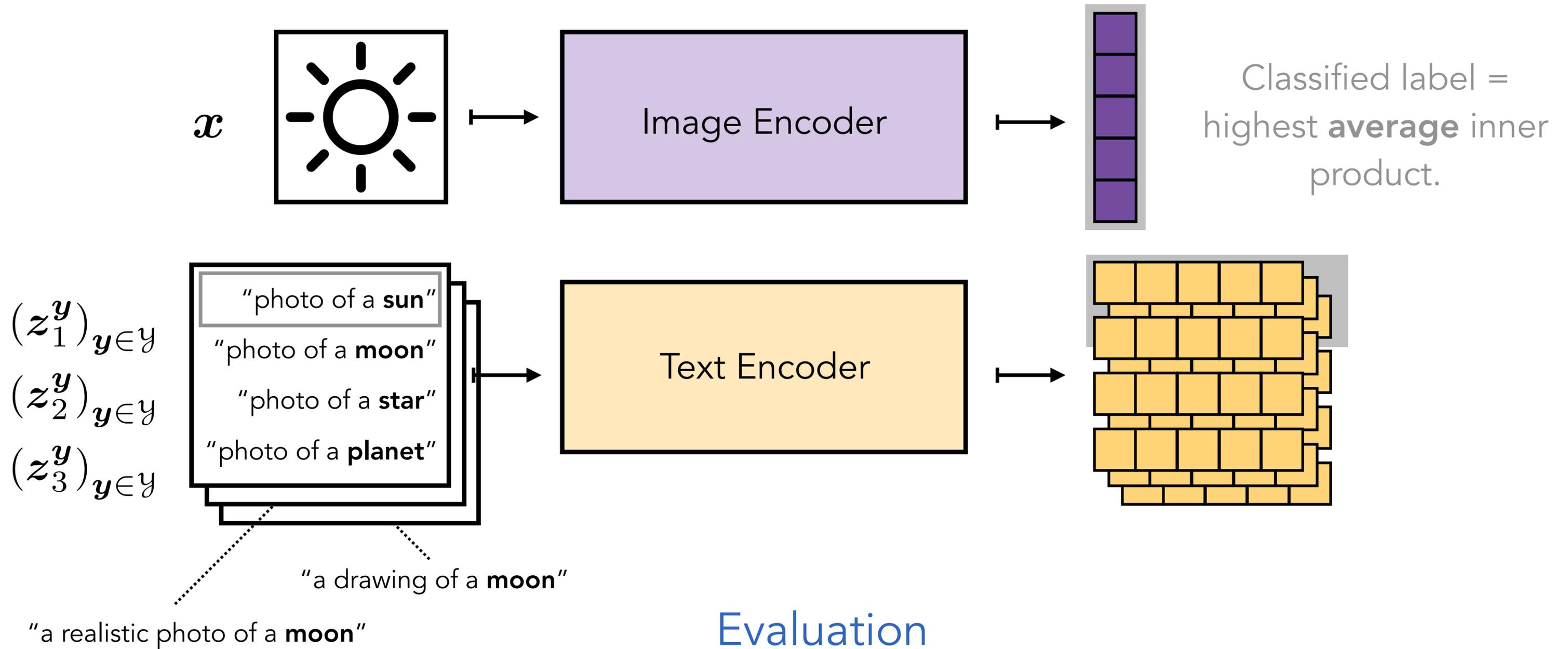
Foundation Modeling



Foundation Modeling



Foundation Modeling



Many empirical studies have confirmed that selecting pre-training data is arguably the most important part of foundation modeling. Specifically, **balancing** paired data to achieve certain marginal constraints is an effective recipe.

DATAComp:
In search of the next generation of multimodal datasets

Samir Yitzhak Gadre^{*2}, Gabriel Ilharco^{*1}, Alex Fang^{*1}, Jonathan Hayase¹, Georgios Smyrnis⁵, Thao Nguyen¹, Ryan Marten^{7,9}, Mitchell Wortsman¹, Dhruva Ghosh¹, Jieyu Zhang¹, Eyal Orgad³, Rahim Entezari¹⁰, Giannis Daras⁵, Sarah Pratt¹, Vivek Ramanujan¹, Yonatan Bitton¹¹, Kalyani Marathe¹, Stephen Mussmann¹, Richard Vencu⁶, Mehdi Cherti^{6,8}, Raniyu Krishna¹, Song²,
1, 3,

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

**DoReMi: Optimizing Data Mixtures
Speeds Up Language Model Pretraining**

Sang Michael Xie^{1,2}, Hieu Pham¹, Xuanyi Dong¹, Nan Du¹, Hanxiao Liu¹, Yifeng Lu¹, Percy Liang², Quoc V. Le¹, Tengyu Ma², and Adams Wei Yu¹

¹Google DeepMind
²Stanford University

Techniques for pre-training, prompting, and overall model reuse have advanced significantly in applications. Virtually no **generalization guarantees** (from the perspective of statistical learning theory) exist for this pipeline.

UNDERSTANDING TRANSFERABLE REPRESENTATION
LEARNING AND ZERO-SHOT TRANSFER IN CLIP

Zixiang Chen^{†*}, Yihe Deng^{†*}, Yuanzhi Li[°], Quanquan Gu[†]
[†]Department of Computer Science, University of California, Los Angeles

Language in a Bottle: Language Model Guided Concept Bottlenecks
for Interpretable Image Classification

Yue Yang, Artemis Panagopoulou, Shenghao Zhou, Daniel Jin,

Enhancing CLIP with GPT-4: Harnessing Visual Descriptions as Prompts

Mayug Maniparambil, Chris Vorster, Derek Molloy, Noel Murphy, Kevin McGuinness, Noel E. O'Connor
ML Labs, Dublin City University,
Dublin, Ireland

What does a platypus look like?
Generating customized prompts for zero-shot image classification

Sarah Pratt^{1*} Ian Covert¹ Rosanne Liu^{2,3} Ali Farhadi¹
¹University of Washington ²Google DeepMind ³ML Collective

Many empirical studies have confirmed that selecting pre-training data is arguably the most important part of foundation modeling. Specifically, **balancing** paired data to achieve certain marginal constraints is an effective recipe.

DATAComp:
In search of the next generation of multimodal datasets

Samir Yitzhak Gadre^{*2}, Gabriel Ilharco^{*1}, Alex Fang^{*1}, Jonathan Hayase¹, Georgios Smyrnis⁵, Thao Nguyen¹, Ryan Marten^{7,9}, Mitchell Wortsman¹, Dhruva Ghosh¹, Jieyu Zhang¹, Eyal Orgad³, Rahim Entezari¹⁰, Giannis Daras⁵, Sarah Pratt¹, Vivek Ramanujan¹, Yonatan Bitton¹¹, Kalyani Marathe¹, Stephen Muecmann¹, Richard Vencu⁶, Mehdi Cherti^{6,8}, Ranjay Krishna¹, Song², Song², Song²

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

**DoReMi: Optimizing Data Mixtures
Speeds Up Language Model Pretraining**

Sang Michael Xie^{1,2}, Hieu Pham¹, Xuanyi Dong¹, Nan Du¹, Hanxiao Liu¹, Yifeng Lu¹, Percy Liang², Quoc V. Le¹, Tengyu Ma², and Adams Wei Yu¹

¹Google DeepMind
²Stanford University

Techniques for pre-training, prompting, and overall model reuse have advanced significantly in applications. Virtually no **generalization guarantees** (from the perspective of statistical learning theory) exist for this pipeline.

UNDERSTANDING TRANSFERABLE REPRESENTATION LEARNING AND ZERO-SHOT TRANSFER IN CLIP

Zixiang Chen^{†*}, Yihe Deng^{†*}, Yuanzhi Li[°], Quanquan Gu[‡]
[‡]Department of Computer Science, University of California, Los Angeles

Language in a Bottle: Language Model Guided Concept Bottlenecks for Interpretable Image Classification

Yue Yang, Artemis Panagopoulou, Shenghao Zhou, Daniel Jin,

Enhancing CLIP with GPT-4: Harnessing Visual Descriptions as Prompts

Mayug Maniparambil, Chris Vorster, Derek Molloy, Noel Murphy, Kevin McGuinness, Noel E. O'Connor
ML Labs, Dublin City University, Dublin, Ireland

**What does a platypus look like?
Generating customized prompts for zero-shot image classification**

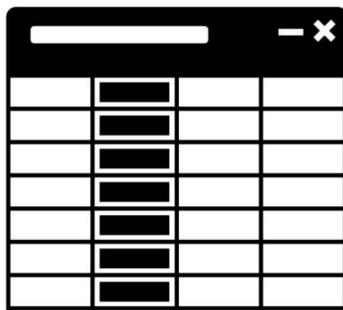
Sarah Pratt^{1*} Ian Covert¹ Rosanne Liu^{2,3} Ali Farhadi¹

¹University of Washington ²Google DeepMind ³ML Collective

Motivating Questions

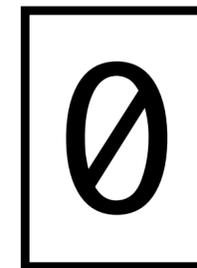
Foundation Modeling

What is the statistical effect of **balancing** methods on the pre-training and the resulting foundation model?



Zero-Shot Prediction

How do we analyze prompt-based **zero-shot prediction** as a statistical estimator, such as its comparison to direct supervision?

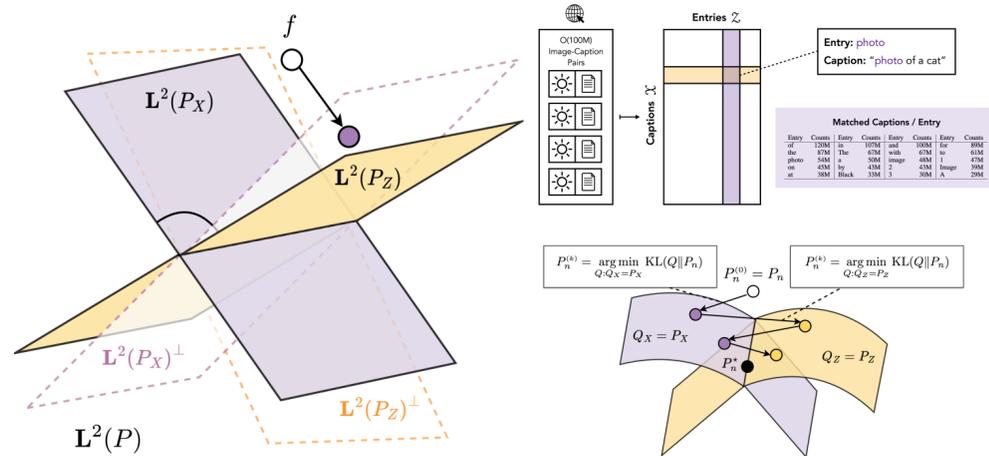


Research Contributions

Exact variance reduction analysis of balancing-based mean estimation!

Theorem (Liu, M., Pal, Harchaoui)

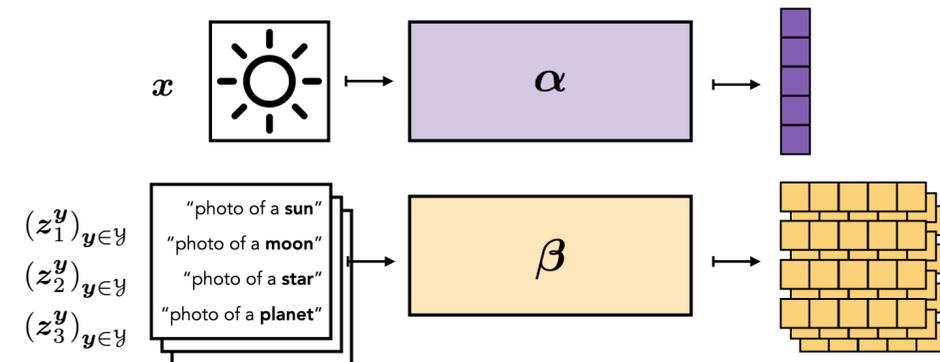
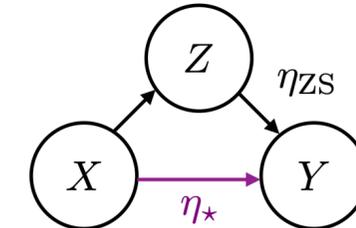
$$\mathbb{E}_P [(P_n^{(k)}(h) - P(h))^2] = \frac{\text{Var}(\dots \overset{k \text{ times}}{\mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp} h)}{n} + \tilde{O}\left(\frac{k^6}{n^{3/2}}\right)$$



Theoretical framework for obtaining generalization guarantees for two broad classes of zero-shot prediction models!

Thm. 1 (M., Harchaoui)

$$\|\eta_\star - \eta_{ZS}\|_{L^2(Q_X)}^2 \lesssim \mathbb{E}_{Q_Z} [I(X, Y|Z)] + \|g_Q - g_R\|_{L^2(Q_Z)}^2$$



Preliminaries

Balanced Pre-Training

Zero-Shot Prediction

Conclusion

\mathcal{X}

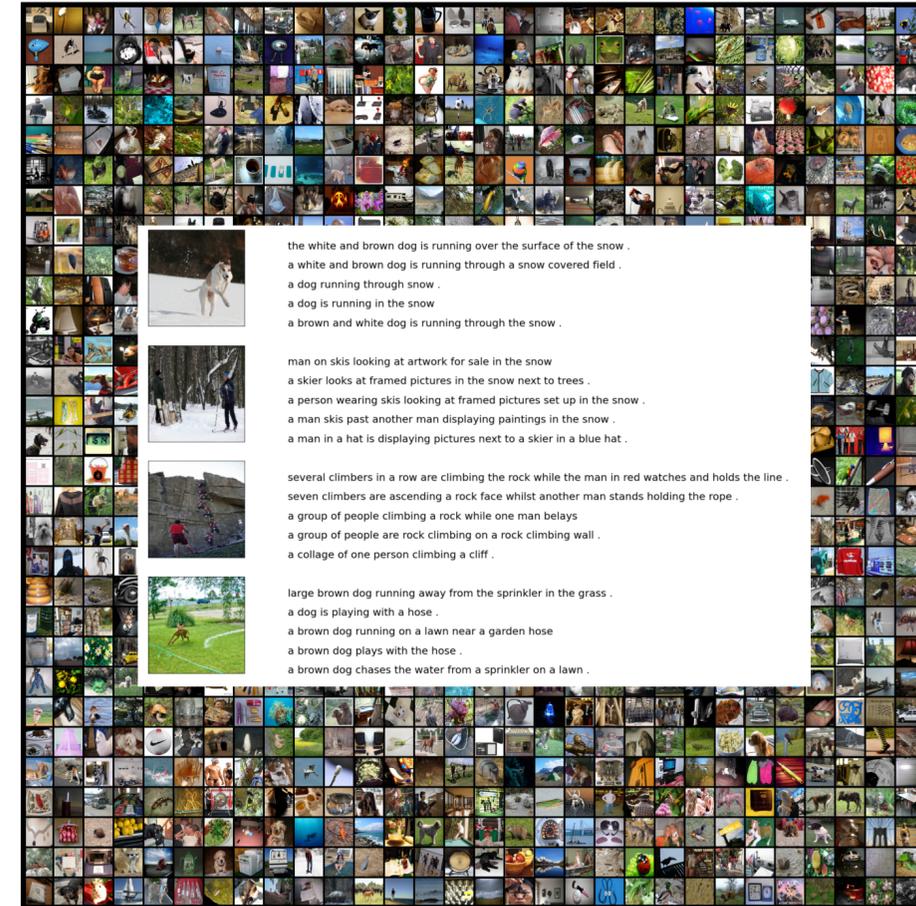
Sample Space (Images)

 \mathcal{Z}

Sample Space (Text)

$$P \equiv P_{\mathcal{X}, \mathcal{Z}}$$

Pre-Training Distribution



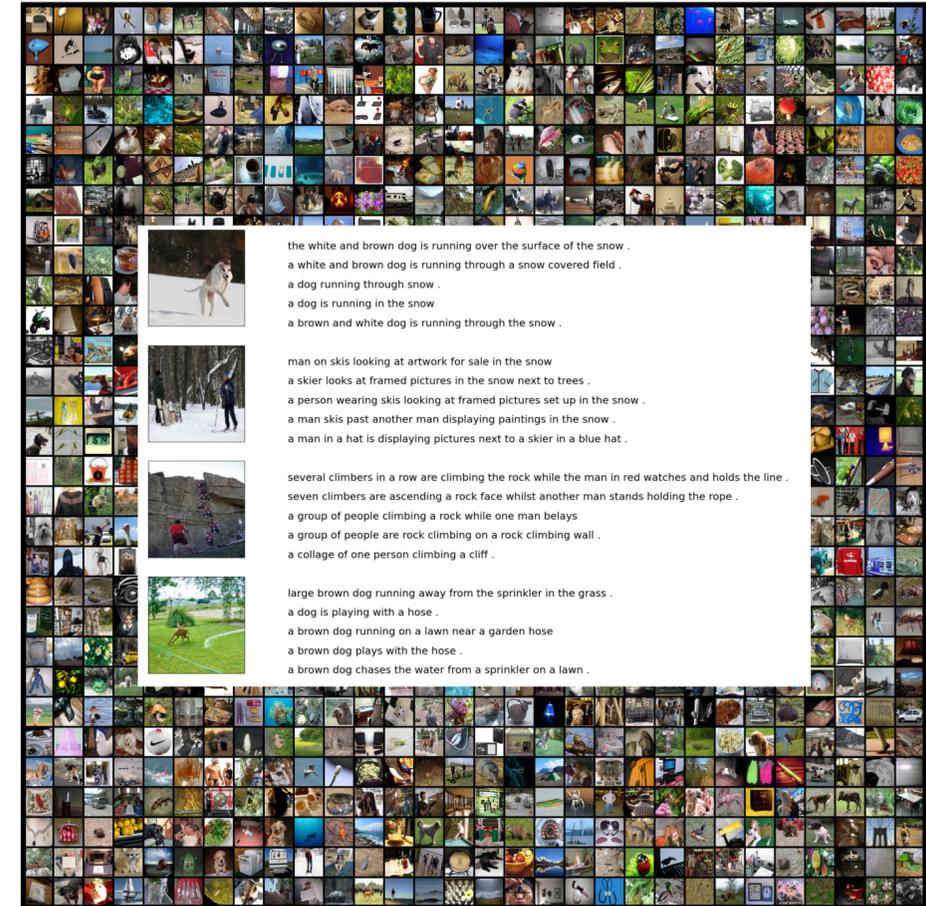
\mathcal{X} Sample Space (Images)

\mathcal{Z} Sample Space (Text)

$P \equiv P_{X,Z}$ Pre-Training Distribution

$$\mathbf{L}^2(P) = \left\{ f : \|f\|_{\mathbf{L}^2(P)}^2 := \int_{\mathcal{X} \times \mathcal{Z}} f^2(\mathbf{x}, \mathbf{z}) dP(\mathbf{x}, \mathbf{z}) < \infty \right\}$$

$\mathbf{L}^2(P_X), \mathbf{L}^2(P_Z)$ Defined Similarly



\mathcal{X} Sample Space (Images)

\mathcal{Z} Sample Space (Text)

$P \equiv P_{X,Z}$ Pre-Training Distribution

$$\mathbf{L}^2(P) = \left\{ f : \|f\|_{\mathbf{L}^2(P)}^2 := \int_{\mathcal{X} \times \mathcal{Z}} f^2(\mathbf{x}, \mathbf{z}) dP(\mathbf{x}, \mathbf{z}) < \infty \right\}$$

$\mathbf{L}^2(P_X), \mathbf{L}^2(P_Z)$ Defined Similarly

Conditional Mean

Projecting onto $\mathbf{L}^2(P_X), \mathbf{L}^2(P_Z)$

$$f \in \mathbf{L}^2(P) \mapsto \begin{cases} \mathbb{E}_P [f(X, Z) | Z] (\cdot) \in \mathbf{L}^2(P_Z) \\ \mathbb{E}_P [f(X, Z) | X] (\cdot) \in \mathbf{L}^2(P_X) \end{cases}$$

\mathcal{X} Sample Space (Images)

\mathcal{Z} Sample Space (Text)

$P \equiv P_{X,Z}$ Pre-Training Distribution

$$\mathbf{L}^2(P) = \left\{ f : \|f\|_{\mathbf{L}^2(P)}^2 := \int_{\mathcal{X} \times \mathcal{Z}} f^2(\mathbf{x}, \mathbf{z}) dP(\mathbf{x}, \mathbf{z}) < \infty \right\}$$

$\mathbf{L}^2(P_X), \mathbf{L}^2(P_Z)$ Defined Similarly

Conditional Mean

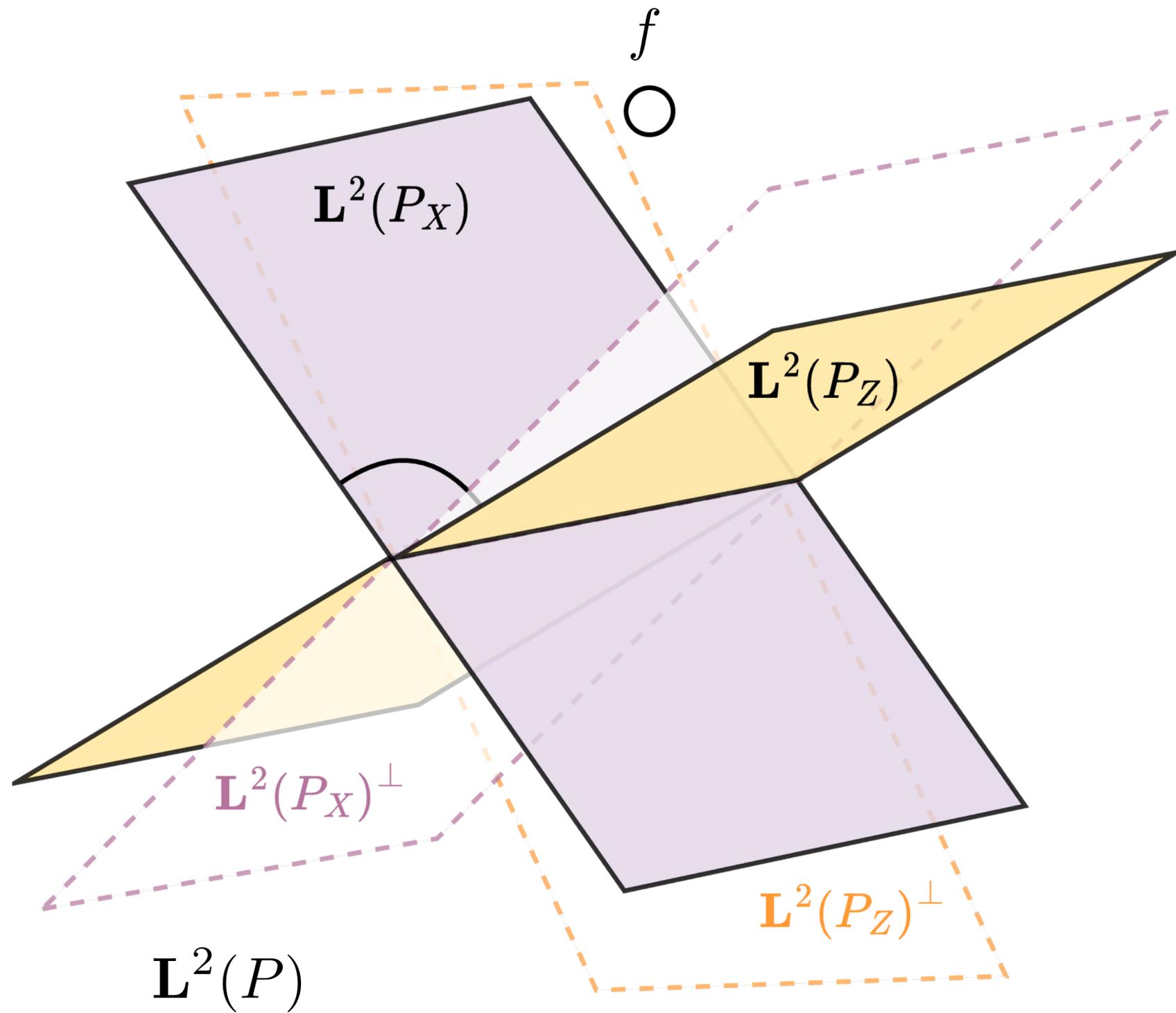
Projecting onto $\mathbf{L}^2(P_X), \mathbf{L}^2(P_Z)$

$$f \in \mathbf{L}^2(P) \mapsto \begin{cases} \mathbb{E}_P [f(X, Z)|Z] (\cdot) \in \mathbf{L}^2(P_Z) \\ \mathbb{E}_P [f(X, Z)|X] (\cdot) \in \mathbf{L}^2(P_X) \end{cases}$$

Conditional Centering

Projecting onto $\mathbf{L}^2(P_X)^\perp, \mathbf{L}^2(P_Z)^\perp$

$$f \in \mathbf{L}^2(P) \mapsto \begin{cases} f - \mathbb{E}_P [f(X, Z)|Z] (\cdot) \in \mathbf{L}^2(P) \\ f - \mathbb{E}_P [f(X, Z)|X] (\cdot) \in \mathbf{L}^2(P) \end{cases}$$



Conditional Mean

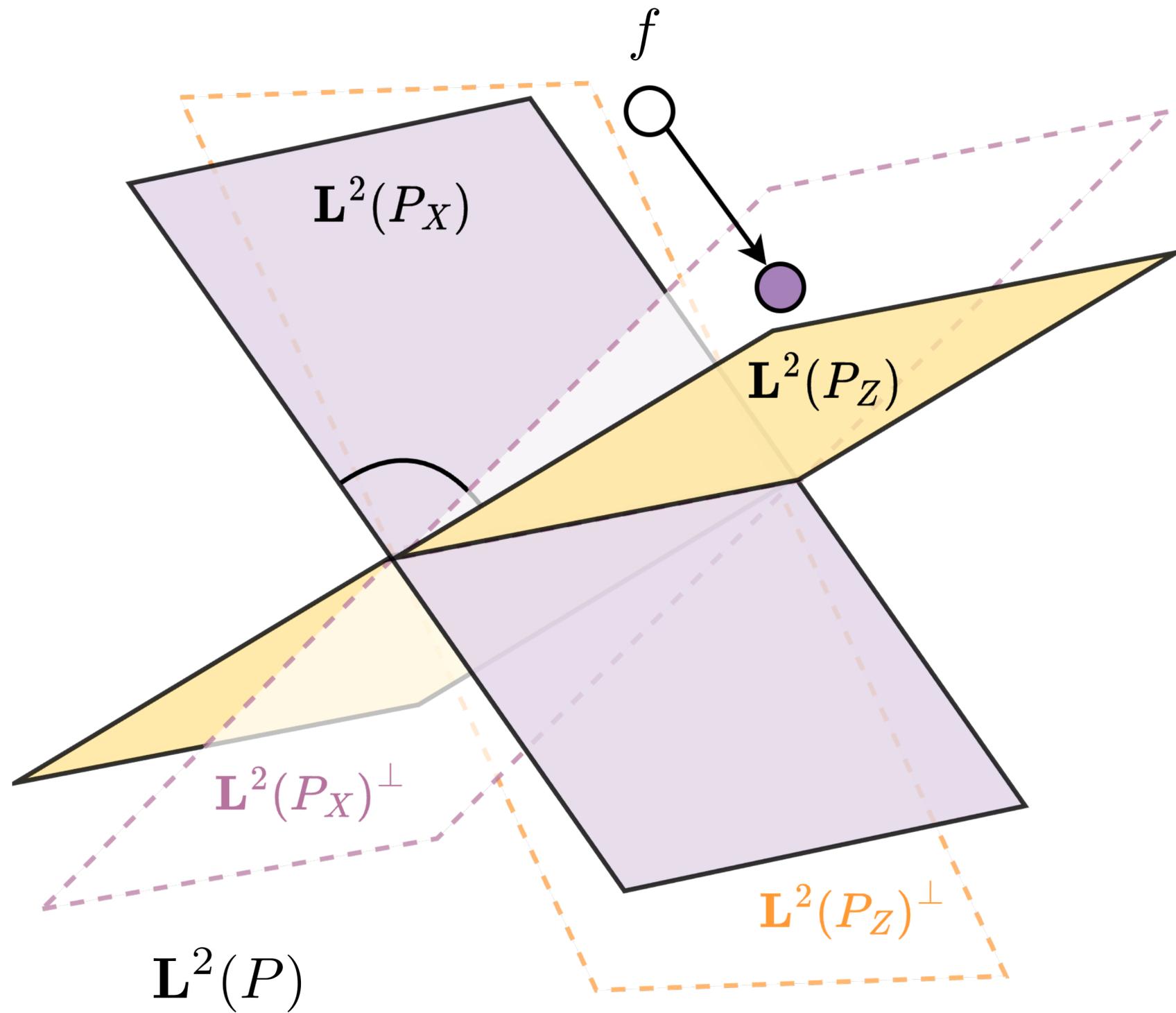
Projecting onto $\mathbf{L}^2(P_X), \mathbf{L}^2(P_Z)$

$$f \in \mathbf{L}^2(P) \mapsto \begin{cases} \mathbb{E}_P [f(X, Z)|Z] (\cdot) \in \mathbf{L}^2(P_Z) \\ \mathbb{E}_P [f(X, Z)|X] (\cdot) \in \mathbf{L}^2(P_X) \end{cases}$$

Conditional Centering

Projecting onto $\mathbf{L}^2(P_X)^\perp, \mathbf{L}^2(P_Z)^\perp$

$$f \in \mathbf{L}^2(P) \mapsto \begin{cases} f - \mathbb{E}_P [f(X, Z)|Z] (\cdot) \in \mathbf{L}^2(P) \\ f - \mathbb{E}_P [f(X, Z)|X] (\cdot) \in \mathbf{L}^2(P) \end{cases}$$



Conditional Mean

Projecting onto $\mathbf{L}^2(P_X), \mathbf{L}^2(P_Z)$

$$f \in \mathbf{L}^2(P) \mapsto \begin{cases} \mathbb{E}_P [f(X, Z)|Z] (\cdot) \in \mathbf{L}^2(P_Z) \\ \mathbb{E}_P [f(X, Z)|X] (\cdot) \in \mathbf{L}^2(P_X) \end{cases}$$

Conditional Centering

Projecting onto $\mathbf{L}^2(P_X)^\perp, \mathbf{L}^2(P_Z)^\perp$

$$f \in \mathbf{L}^2(P) \mapsto \begin{cases} f - \mathbb{E}_P [f(X, Z)|Z] (\cdot) \in \mathbf{L}^2(P) \\ f - \mathbb{E}_P [f(X, Z)|X] (\cdot) \in \mathbf{L}^2(P) \end{cases}$$

Conditional Mean Operator: Singular Value Decomposition

$$\mu_{X \leftarrow Z} : \mathbf{L}^2(P) \rightarrow \mathbf{L}^2(P_X)$$

$$[\mu_{X \leftarrow Z} f](\mathbf{x}) = \mathbb{E}_P [f(X, Z) | X](\mathbf{x})$$

Conditional Mean Operator: Singular Value Decomposition

$$\mu_{X \leftarrow Z} : \mathbf{L}^2(P_Z) \rightarrow \mathbf{L}^2(P_X)$$

$$[\mu_{X \leftarrow Z} g](\mathbf{x}) = \mathbb{E}_P [g(Z) | X](\mathbf{x})$$

Conditional Mean Operator: Singular Value Decomposition

$$\mu_{X \leftarrow Z} : \mathbf{L}^2(P_Z) \rightarrow \mathbf{L}^2(P_X)$$

$$[\mu_{X \leftarrow Z} g](\mathbf{x}) = \mathbb{E}_P [g(Z) | X](\mathbf{x})$$

$$\langle h, \mu_{X \leftarrow Z} g \rangle_{\mathbf{L}^2(P_X)} = \mathbb{E}_{P_X} [h(X) \mathbb{E}_P [g(Z) | X]] = \mathbb{E}_P [h(X) g(Z)]$$

Conditional Mean Operator: Singular Value Decomposition

$$\mu_{X \leftarrow Z} : \mathbf{L}^2(P_Z) \rightarrow \mathbf{L}^2(P_X)$$

$$[\mu_{X \leftarrow Z} g](\mathbf{x}) = \mathbb{E}_P [g(Z) | X](\mathbf{x})$$

$$\langle h, \mu_{X \leftarrow Z} g \rangle_{\mathbf{L}^2(P_X)} = \mathbb{E}_{P_X} [h(X) \mathbb{E}_P [g(Z) | X]] = \mathbb{E}_P [h(X) g(Z)] = \mathbb{E}_{P_Z} [\mathbb{E}_P [h(X) | Z] g(Z)] = \langle \mu_{Z \leftarrow X} h, g \rangle_{\mathbf{L}^2(P_Z)}$$

The conditional mean of Z given X and vice versa are adjoint!

Conditional Mean Operator: Singular Value Decomposition

$$\mu_{X \leftarrow Z} : \mathbf{L}^2(P_Z) \rightarrow \mathbf{L}^2(P_X)$$

$$[\mu_{X \leftarrow Z} g](\mathbf{x}) = \mathbb{E}_P [g(Z) | X](\mathbf{x})$$

$$\langle h, \mu_{X \leftarrow Z} g \rangle_{\mathbf{L}^2(P_X)} = \mathbb{E}_{P_X} [h(X) \mathbb{E}_P [g(Z) | X]] = \mathbb{E}_P [h(X) g(Z)] = \mathbb{E}_{P_Z} [\mathbb{E}_P [h(X) | Z] g(Z)] = \langle \mu_{Z \leftarrow X} h, g \rangle_{\mathbf{L}^2(P_Z)}$$

The conditional mean of Z given X and vice versa are adjoint!

Basis of $\mathbf{L}^2(P_X)$: $\alpha_1, \alpha_2, \dots$

$$\mu_{X \leftarrow Z} \beta_i = s_i \alpha_i$$

Basis of $\mathbf{L}^2(P_Z)$: β_1, β_2, \dots

$$\mu_{Z \leftarrow X} \alpha_i = s_i \beta_i$$

Conditional Mean Operator: Singular Value Decomposition

$$\mu_{X \leftarrow Z} = \begin{bmatrix} | & | & \cdots \\ \alpha_1 & \alpha_2 & \cdots \\ | & | & \cdots \end{bmatrix} \begin{bmatrix} s_1 & & \\ & s_2 & \\ & & \ddots \end{bmatrix} \begin{bmatrix} - & \beta_1 & - \\ - & \beta_2 & - \\ & \vdots & \end{bmatrix}$$

$$\mu_{Z \leftarrow X} = \begin{bmatrix} | & | & \cdots \\ \beta_1 & \beta_2 & \cdots \\ | & | & \cdots \end{bmatrix} \begin{bmatrix} s_1 & & \\ & s_2 & \\ & & \ddots \end{bmatrix} \begin{bmatrix} - & \alpha_1 & - \\ - & \alpha_2 & - \\ & \vdots & \end{bmatrix}$$

The conditional mean of Z given X and vice versa are adjoint!

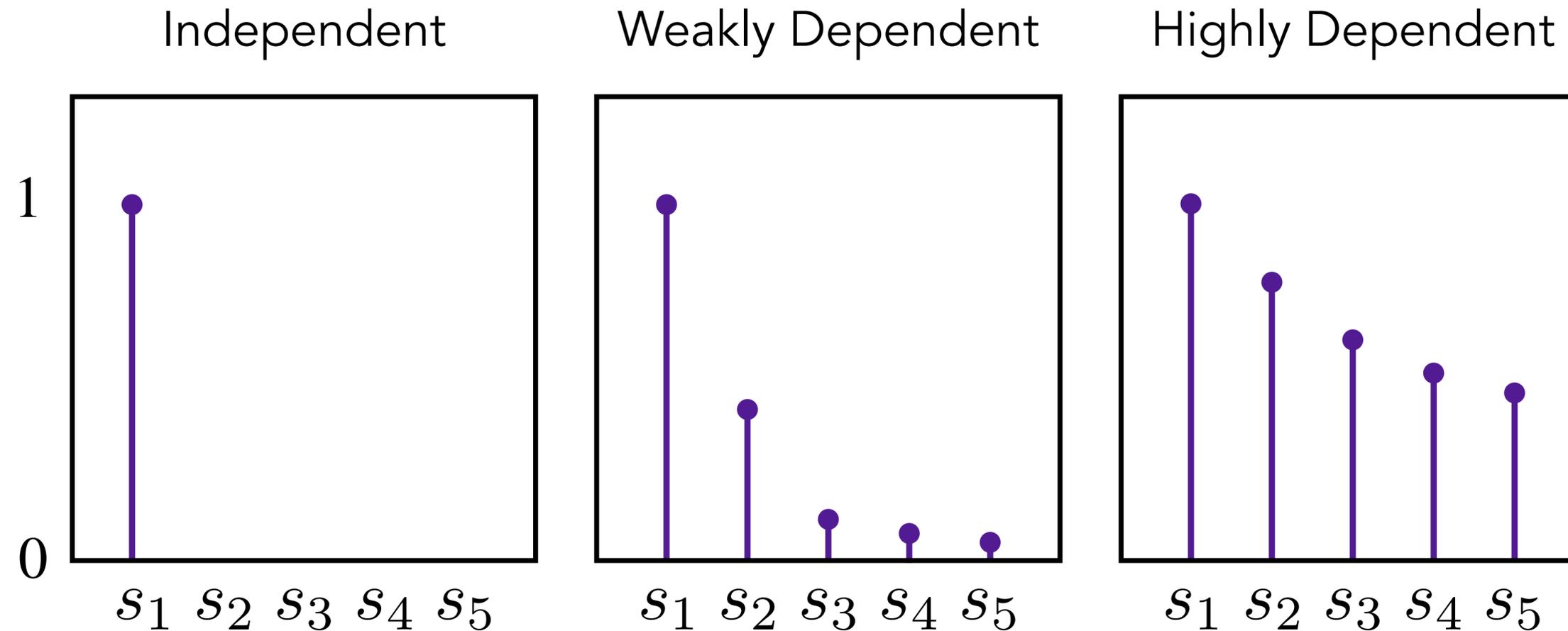
Basis of $\mathbf{L}^2(P_X)$: $\alpha_1, \alpha_2, \dots$

$$\mu_{X \leftarrow Z} \beta_i = s_i \alpha_i$$

Basis of $\mathbf{L}^2(P_Z)$: β_1, β_2, \dots

$$\mu_{Z \leftarrow X} \alpha_i = s_i \beta_i$$

Conditional Mean Operator: Singular Value Decomposition



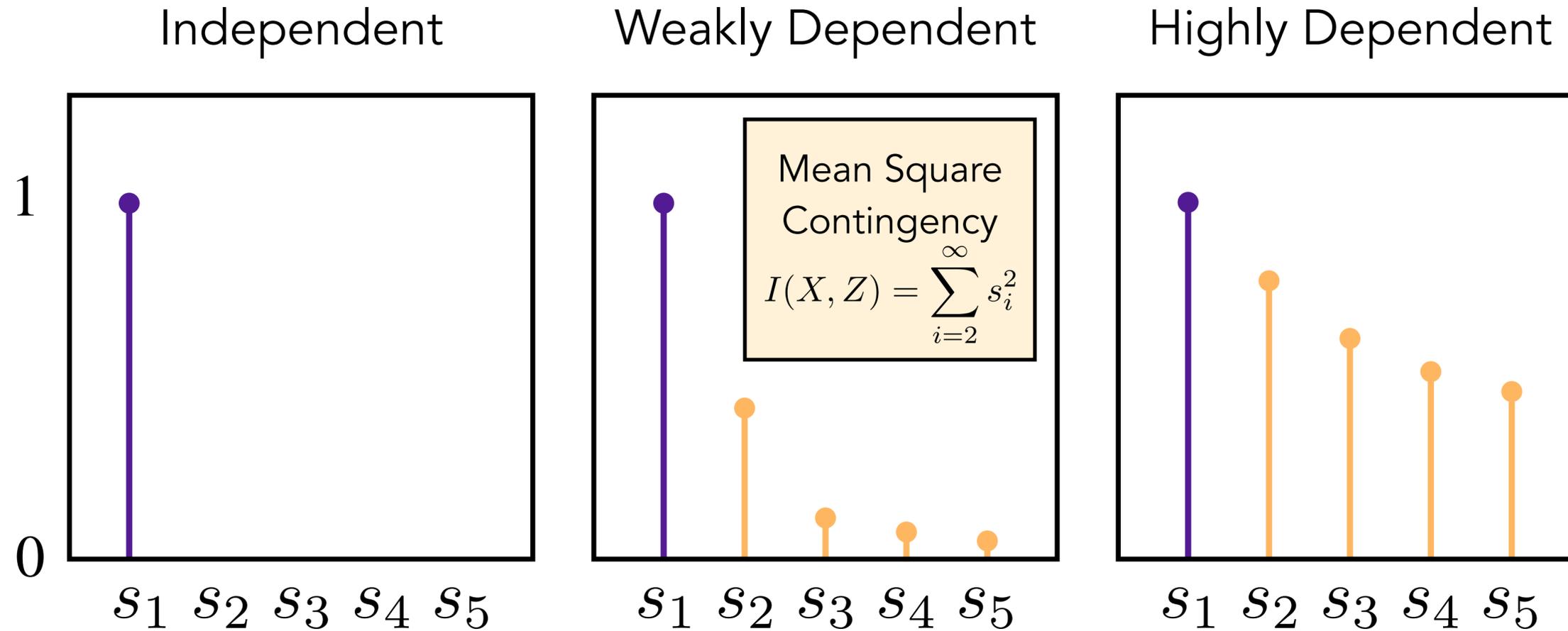
Basis of $\mathbf{L}^2(P_X)$: $\alpha_1, \alpha_2, \dots$

$$\mu_{X \leftarrow Z} \beta_i = s_i \alpha_i$$

Basis of $\mathbf{L}^2(P_Z)$: β_1, β_2, \dots

$$\mu_{Z \leftarrow X} \alpha_i = s_i \beta_i$$

Conditional Mean Operator: Singular Value Decomposition



Basis of $\mathbf{L}^2(P_X)$: $\alpha_1, \alpha_2, \dots$

$$\mu_{X \leftarrow Z} \beta_i = s_i \alpha_i$$

Basis of $\mathbf{L}^2(P_Z)$: β_1, β_2, \dots

$$\mu_{Z \leftarrow X} \alpha_i = s_i \beta_i$$

Conditional Mean \iff Information Density

$$R(\mathbf{x}, \mathbf{z}) := \frac{dP_{X,Z}}{dP_X P_Z}(\mathbf{x}, \mathbf{z}) = \frac{dP_{Z|X=\mathbf{x}}}{dP_Z}(\mathbf{z})$$

Conditional Mean \iff Information Density

$$R(\mathbf{x}, \mathbf{z}) := \frac{dP_{X,Z}}{dP_X P_Z}(\mathbf{x}, \mathbf{z}) = \frac{dP_{Z|X=\mathbf{x}}}{dP_Z}(\mathbf{z})$$

① $[\mu_{X \leftarrow Z} g](\mathbf{x}) = \mathbb{E}_{P_{X,Z}} [g(Z)|X](\mathbf{x}) = \mathbb{E}_{P_Z} [g(Z)R(\mathbf{x}, Z)]$

Conditional Mean \iff Information Density

$$R(\mathbf{x}, \mathbf{z}) := \frac{dP_{X,Z}}{dP_X P_Z}(\mathbf{x}, \mathbf{z}) = \frac{dP_{Z|X=\mathbf{x}}}{dP_Z}(\mathbf{z})$$

① $[\mu_{X \leftarrow Z} g](\mathbf{x}) = \mathbb{E}_{P_{X,Z}} [g(Z) | X](\mathbf{x}) = \mathbb{E}_{P_Z} [g(Z) R(\mathbf{x}, Z)]$

Conditional Mean \iff Information Density

$$R(\boldsymbol{x}, \boldsymbol{z}) := \frac{dP_{X,Z}}{dP_X P_Z}(\boldsymbol{x}, \boldsymbol{z}) = \frac{dP_{Z|X=\boldsymbol{x}}}{dP_Z}(\boldsymbol{z})$$

① $[\mu_{X \leftarrow Z} g](\boldsymbol{x}) = \mathbb{E}_{P_{X,Z}} [g(Z) | X](\boldsymbol{x}) = \mathbb{E}_{P_Z} [g(Z) R(\boldsymbol{x}, Z)]$

Conditional Mean \iff Information Density

$$R(\mathbf{x}, \mathbf{z}) := \frac{dP_{X,Z}}{dP_X P_Z}(\mathbf{x}, \mathbf{z}) = \frac{dP_{Z|X=\mathbf{x}}}{dP_Z}(\mathbf{z})$$

1 $[\mu_{X \leftarrow Z} g](\mathbf{x}) = \mathbb{E}_{P_{X,Z}} [g(Z)|X](\mathbf{x}) = \mathbb{E}_{P_Z} [g(Z)R(\mathbf{x}, Z)]$

2

Conditional Mean \iff Information Density

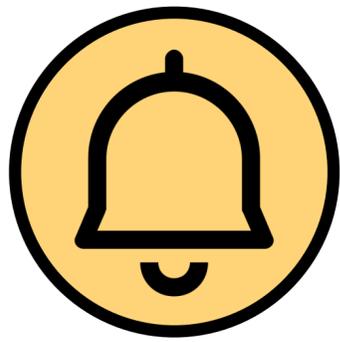
$$R(\mathbf{x}, \mathbf{z}) := \frac{dP_{X,Z}}{dP_X P_Z}(\mathbf{x}, \mathbf{z}) = \frac{dP_{Z|X=\mathbf{x}}}{dP_Z}(\mathbf{z})$$

- 1 $[\mu_{X \leftarrow Z} g](\mathbf{x}) = \mathbb{E}_{P_{X,Z}} [g(Z)|X](\mathbf{x}) = \mathbb{E}_{P_Z} [g(Z)R(\mathbf{x}, Z)]$
- 2 $I(X, Z) = \mathbb{E}_{P_{X,Z}} [(R(X, Z) - 1)^2]$

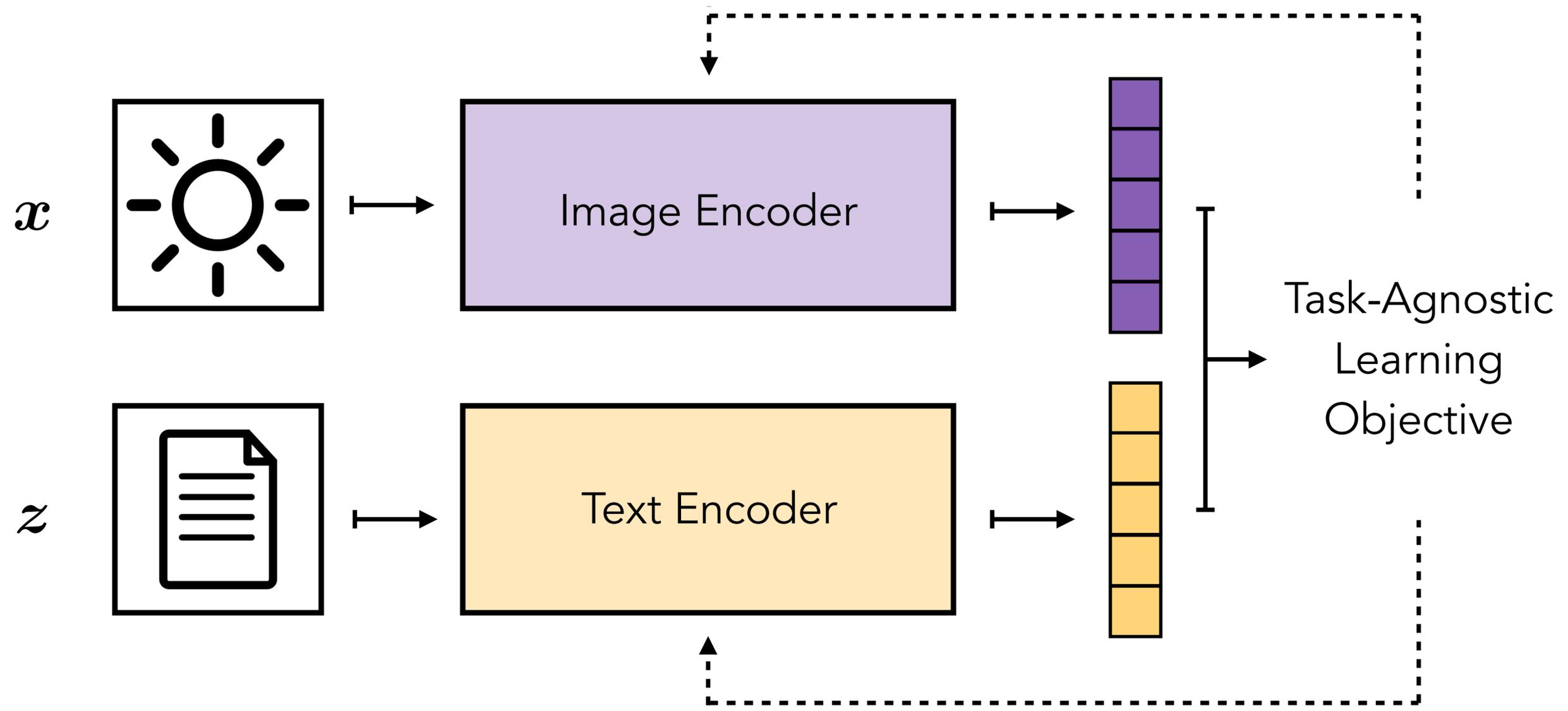
Conditional Mean \iff Information Density

$$R(\mathbf{x}, \mathbf{z}) := \frac{dP_{X,Z}}{dP_X P_Z}(\mathbf{x}, \mathbf{z}) = \frac{dP_{Z|X=\mathbf{x}}}{dP_Z}(\mathbf{z})$$

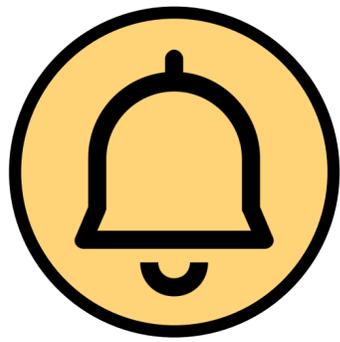
- ① $[\mu_{X \leftarrow Z} g](\mathbf{x}) = \mathbb{E}_{P_{X,Z}} [g(Z)|X](\mathbf{x}) = \mathbb{E}_{P_Z} [g(Z)R(\mathbf{x}, Z)]$
- ② $I(X, Z) = \mathbb{E}_{P_{X,Z}} [(R(X, Z) - 1)^2] = \sum_{i=2}^{\infty} s_i^2$



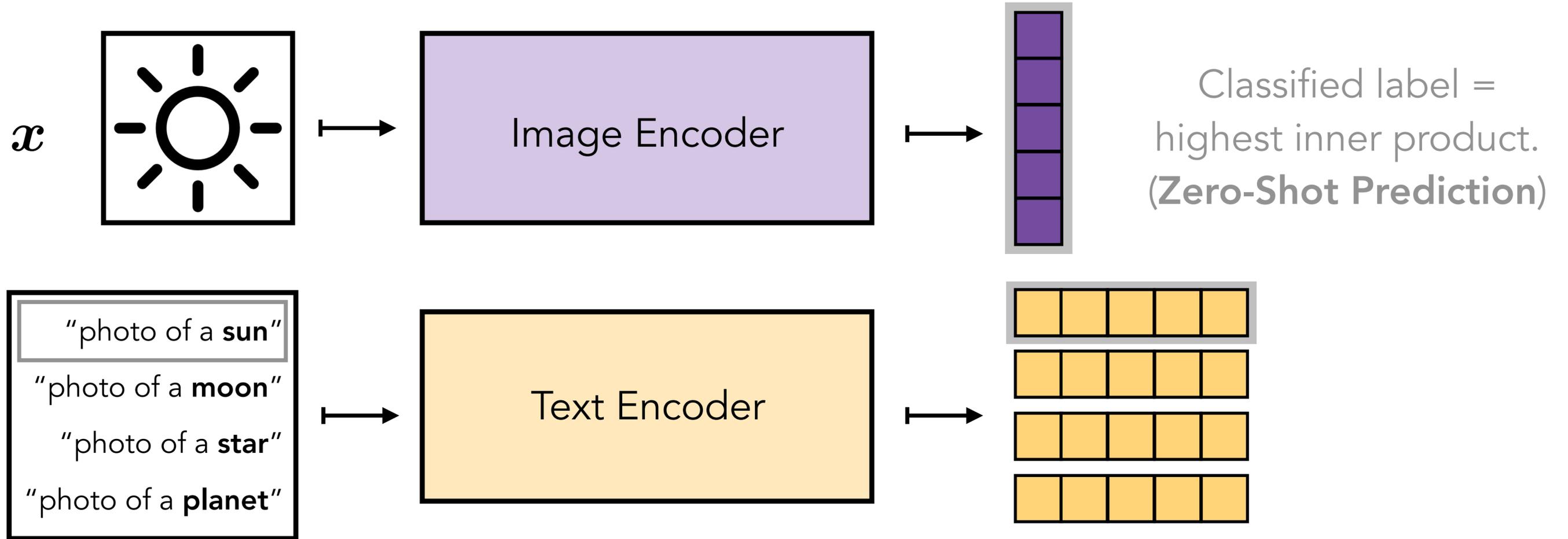
Foundation Modeling



Pre-Training



Foundation Modeling



Idea: Convert labels into pseudo-captions (prompts)

Evaluation

Preliminaries

Balanced Pre-Training

Zero-Shot Prediction

Conclusion



Lang Liu
University of
Washington



Soumik Pal
University of
Washington



Zaid Harchaoui
University of
Washington

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

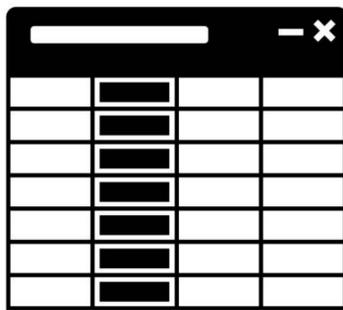
Lang Liu* **Ronak Mehta*** **Soumik Pal** **Zaid Harchaoui**
University of Washington

NeurIPS 2024

Motivating Questions

Foundation Modeling

What is the statistical effect of **balancing** methods on the pre-training and the resulting foundation model?



Zero-Shot Prediction

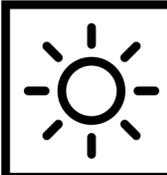
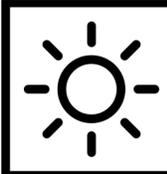
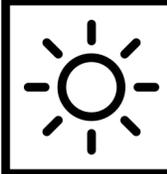
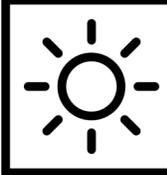
How do we analyze prompt-based **zero-shot prediction** as a statistical estimator, such as its comparison to direct supervision?



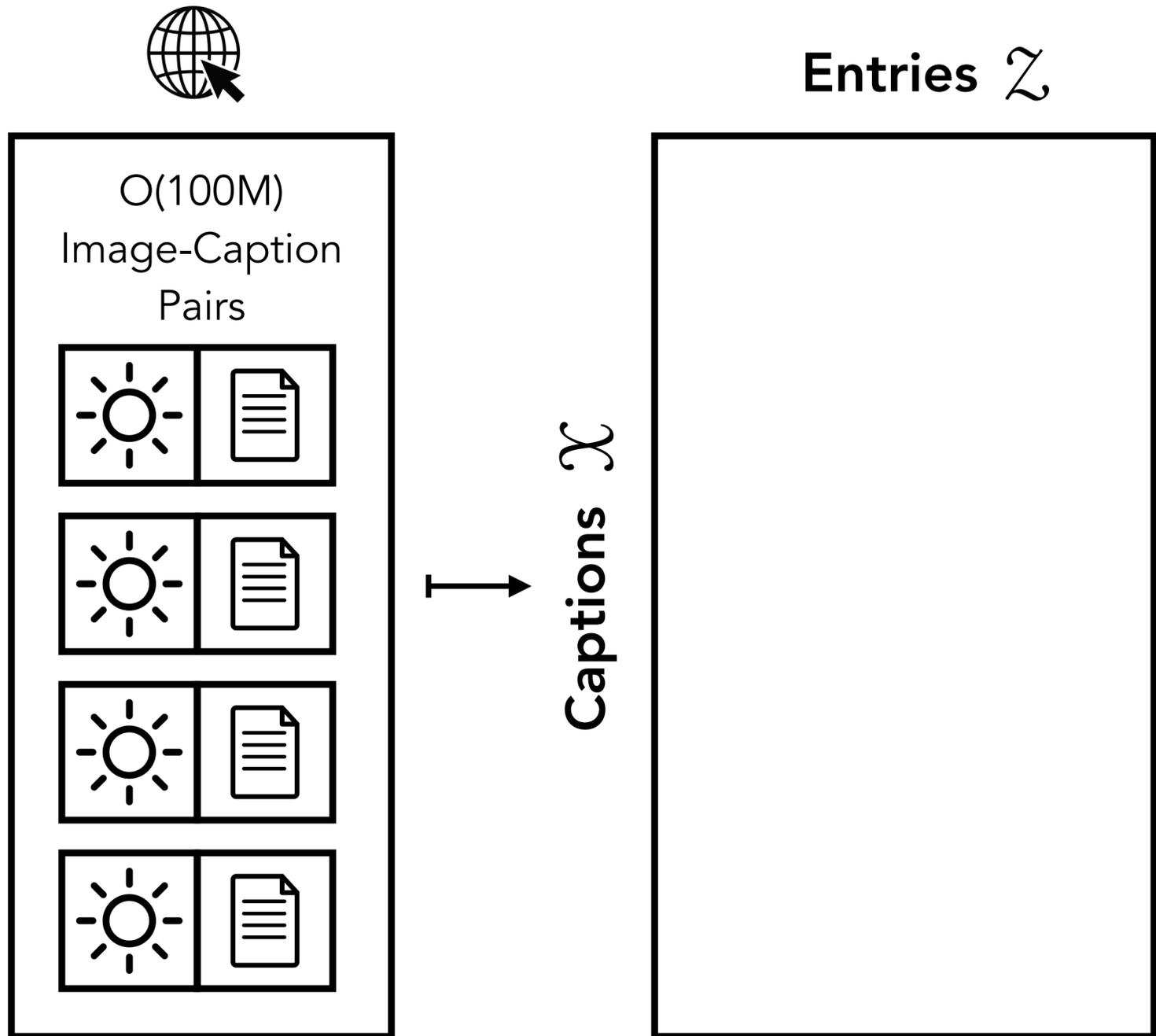
Pre-Training Data Curation: Balancing Keyword Distributions



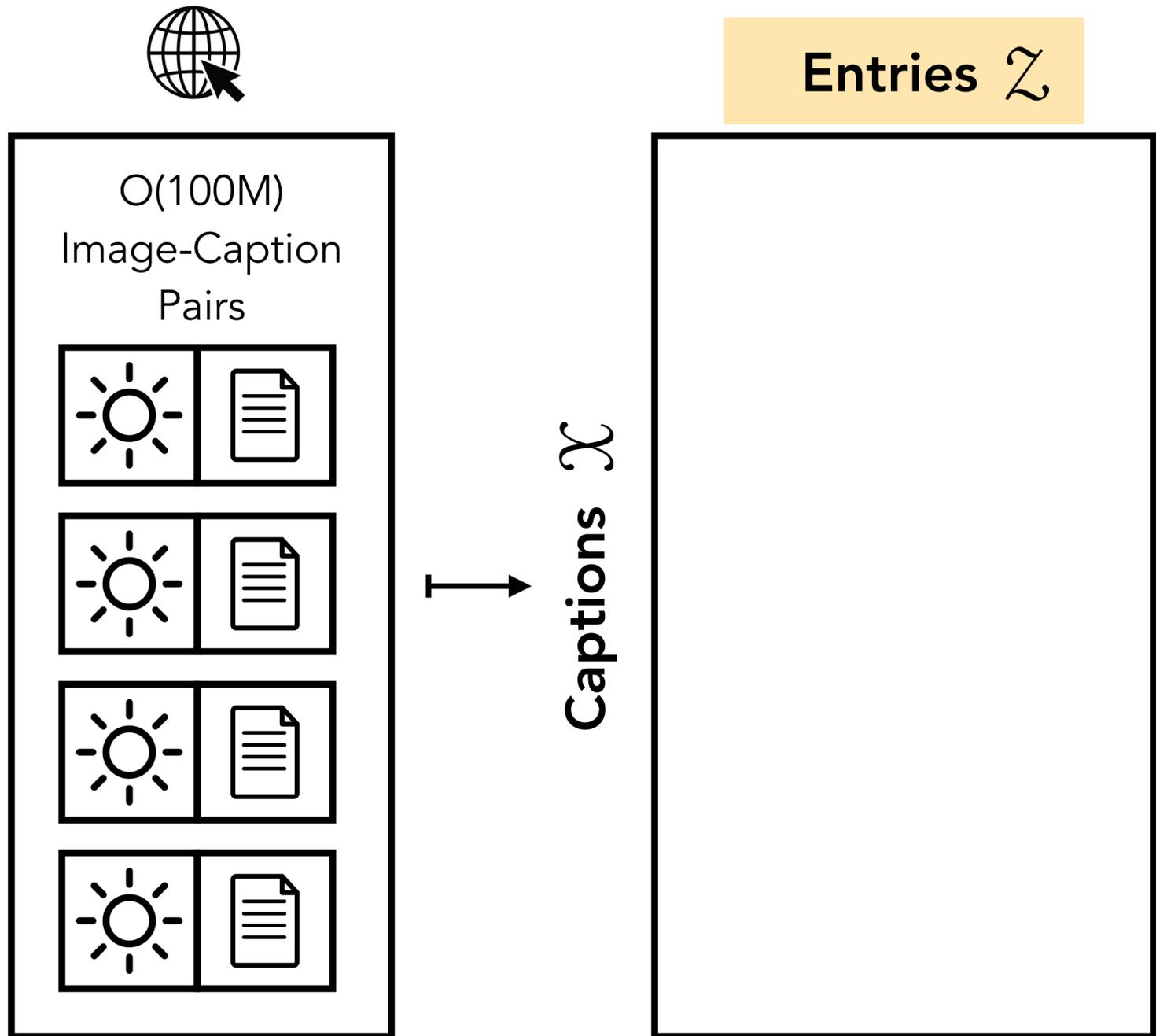
O(100M)
Image-Caption
Pairs

Pre-Training Data Curation: Balancing Keyword Distributions



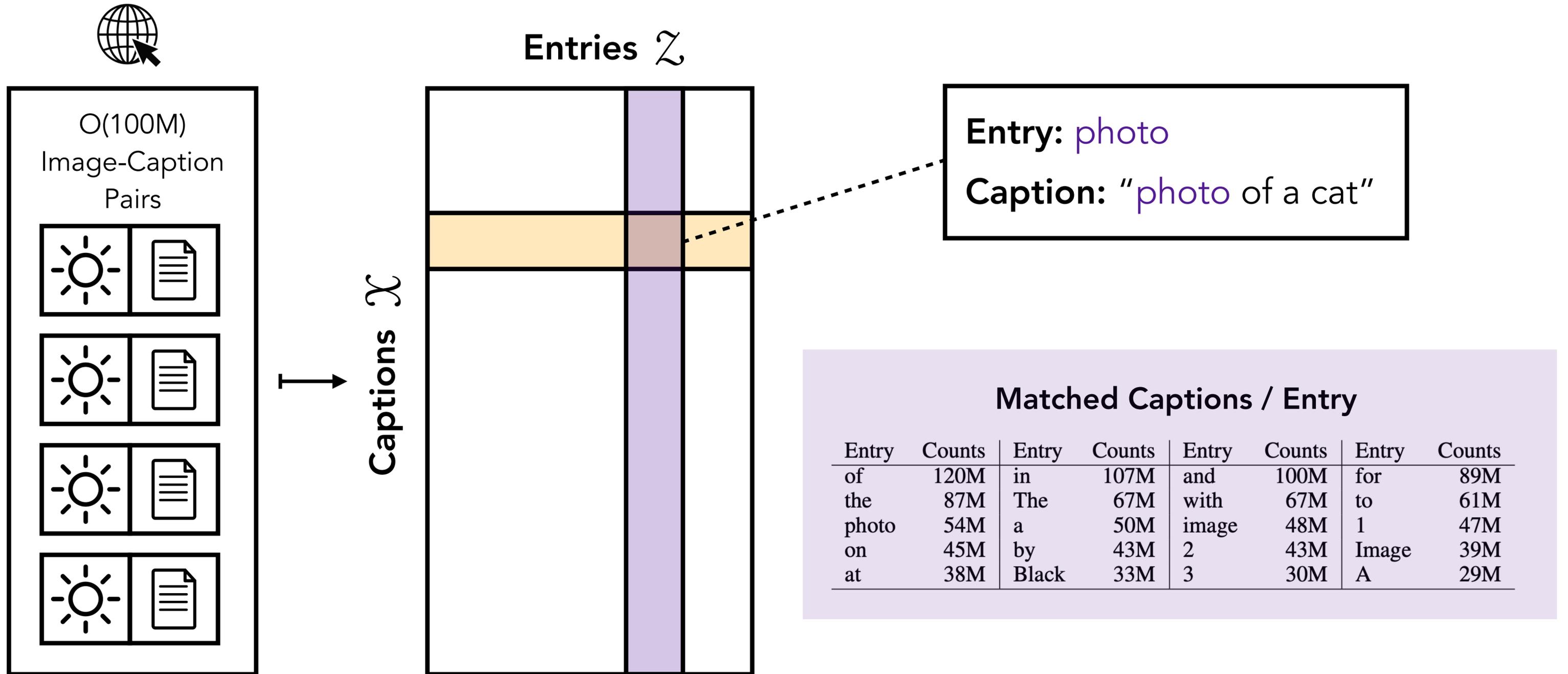
Pre-Training Data Curation: Balancing Keyword Distributions



Matched Captions / Entry

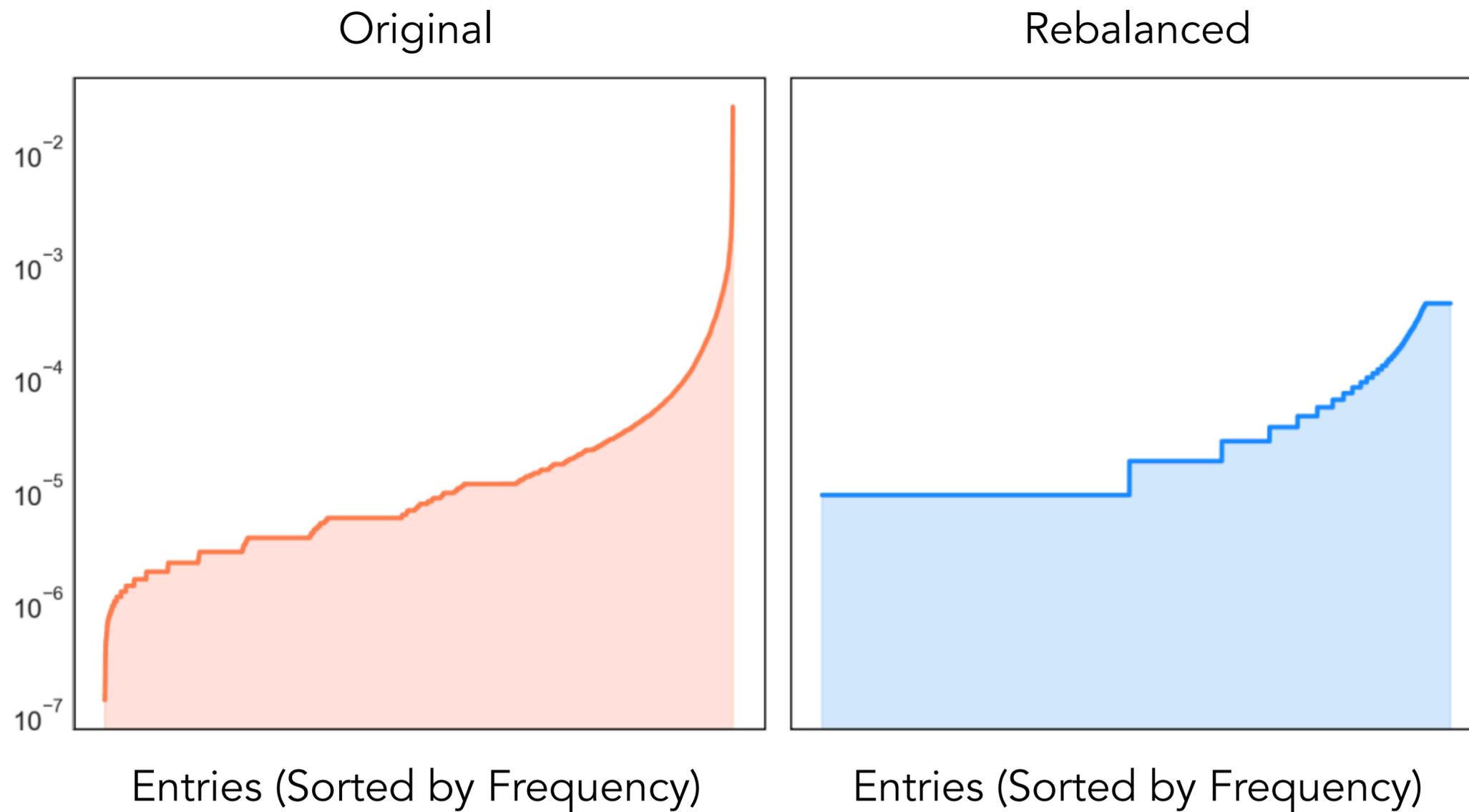
Entry	Counts	Entry	Counts	Entry	Counts	Entry	Counts
of	120M	in	107M	and	100M	for	89M
the	87M	The	67M	with	67M	to	61M
photo	54M	a	50M	image	48M	1	47M
on	45M	by	43M	2	43M	Image	39M
at	38M	Black	33M	3	30M	A	29M

Pre-Training Data Curation: Balancing Keyword Distributions



Pre-Training Data Curation: Balancing Keyword Distributions

Histogram of Entries in Pre-Training Set



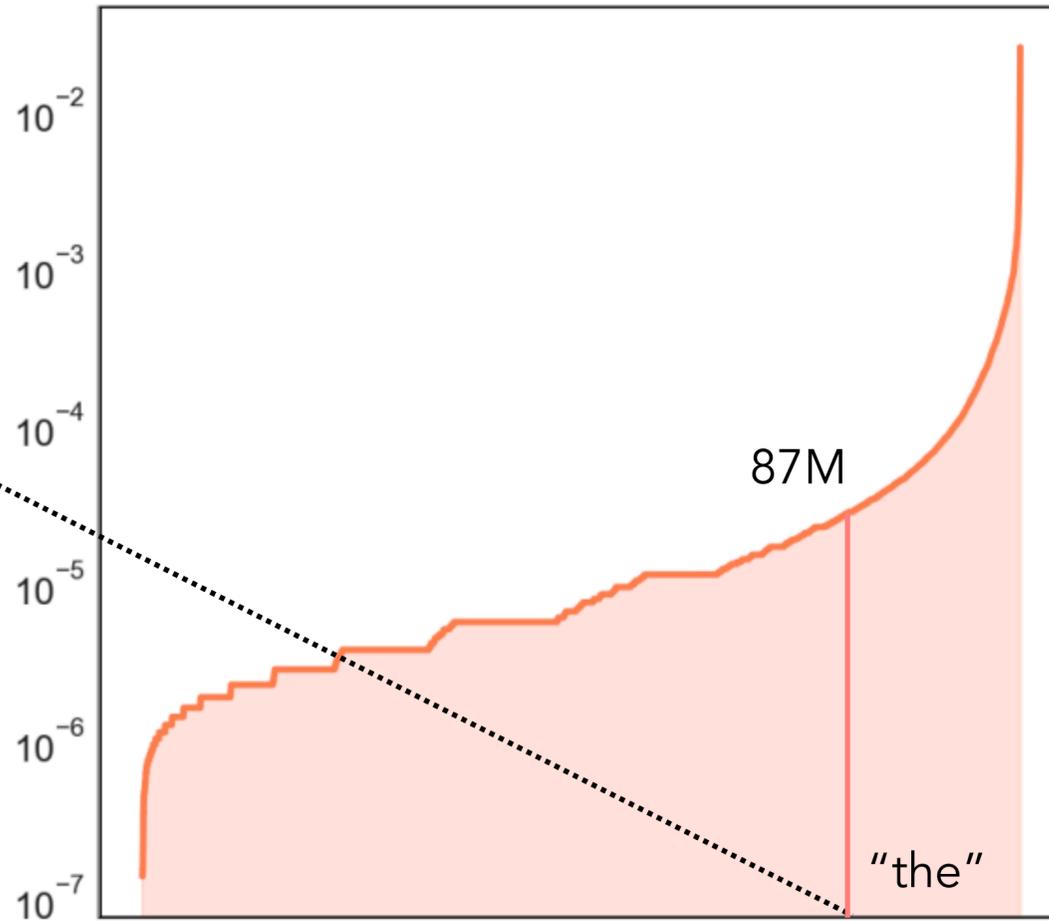
Pre-Training Data Curation: Balancing Keyword Distributions

Histogram of Entries in Pre-Training Set

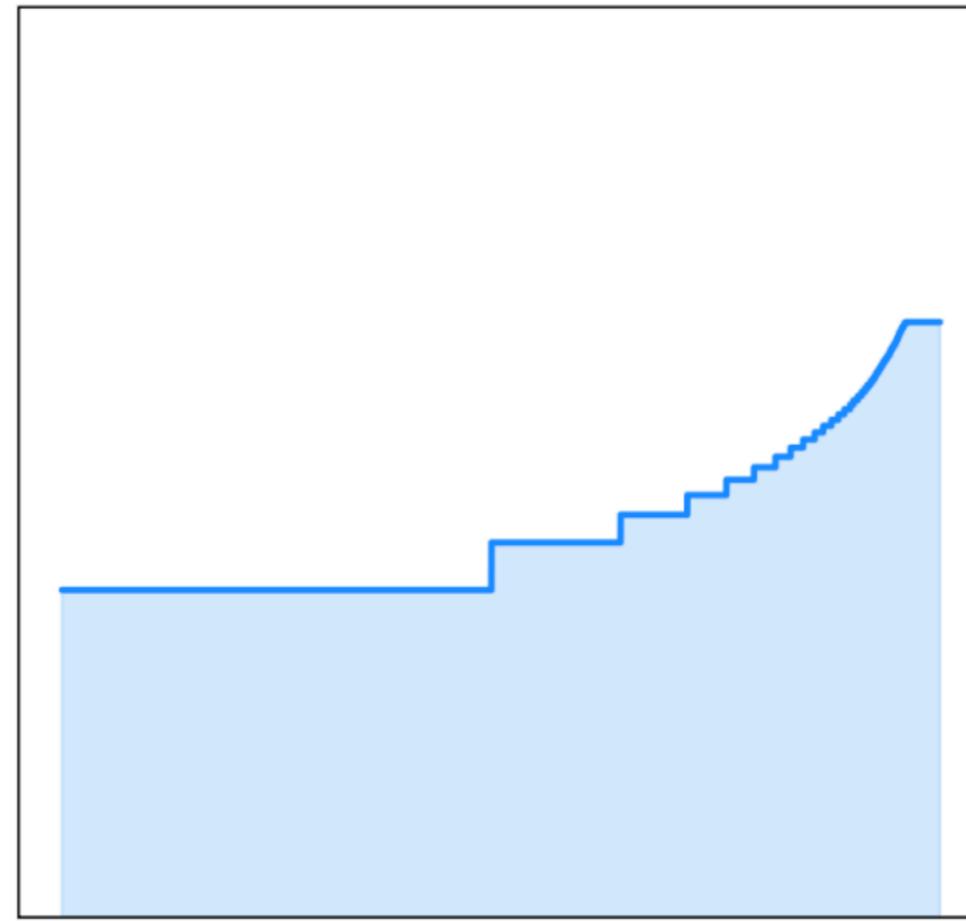
Original

Rebalanced

Entry	Counts
of	120M
the	87M
photo	54M
on	45M
at	38M



Entries (Sorted by Frequency)

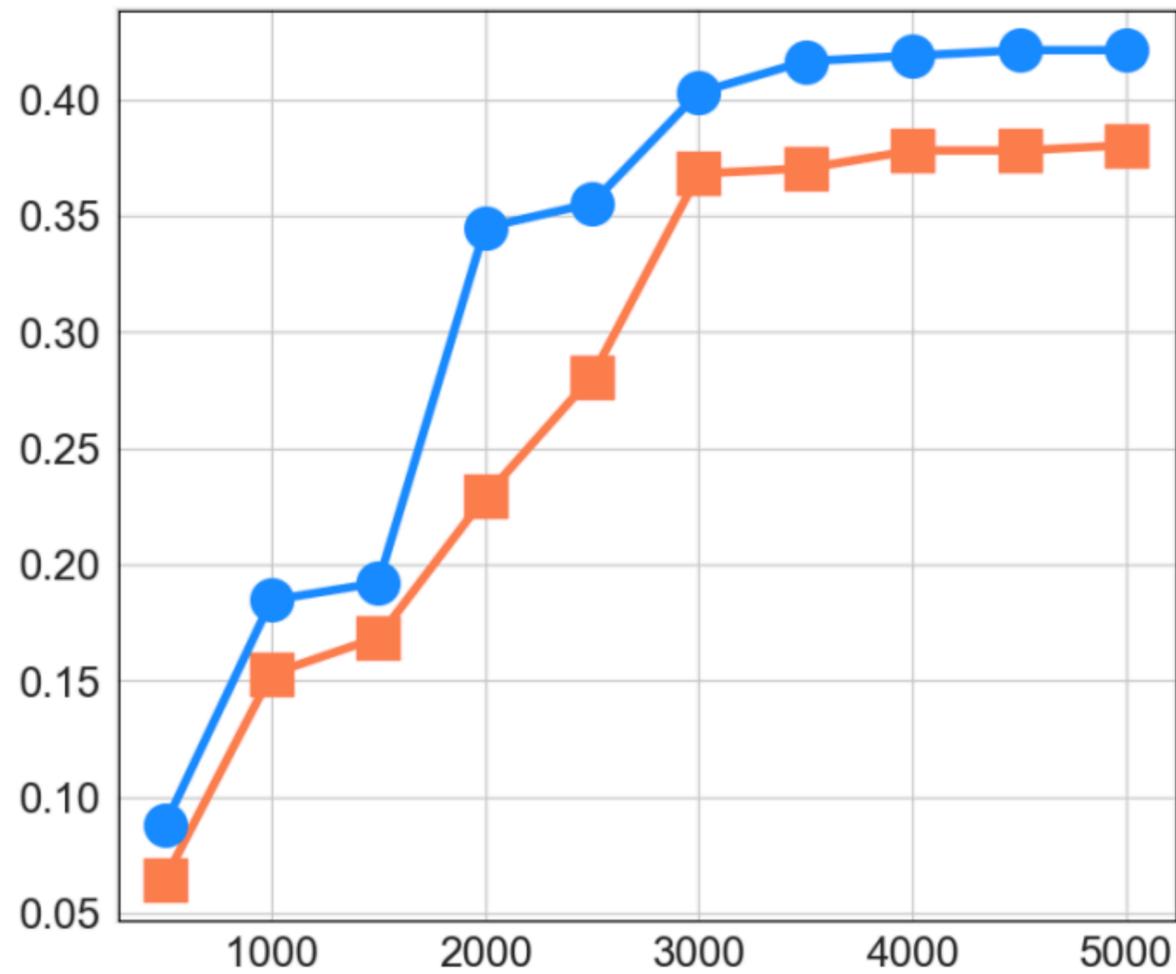


Entries (Sorted by Frequency)

Pre-Training Data Curation: Balancing Keyword Distributions

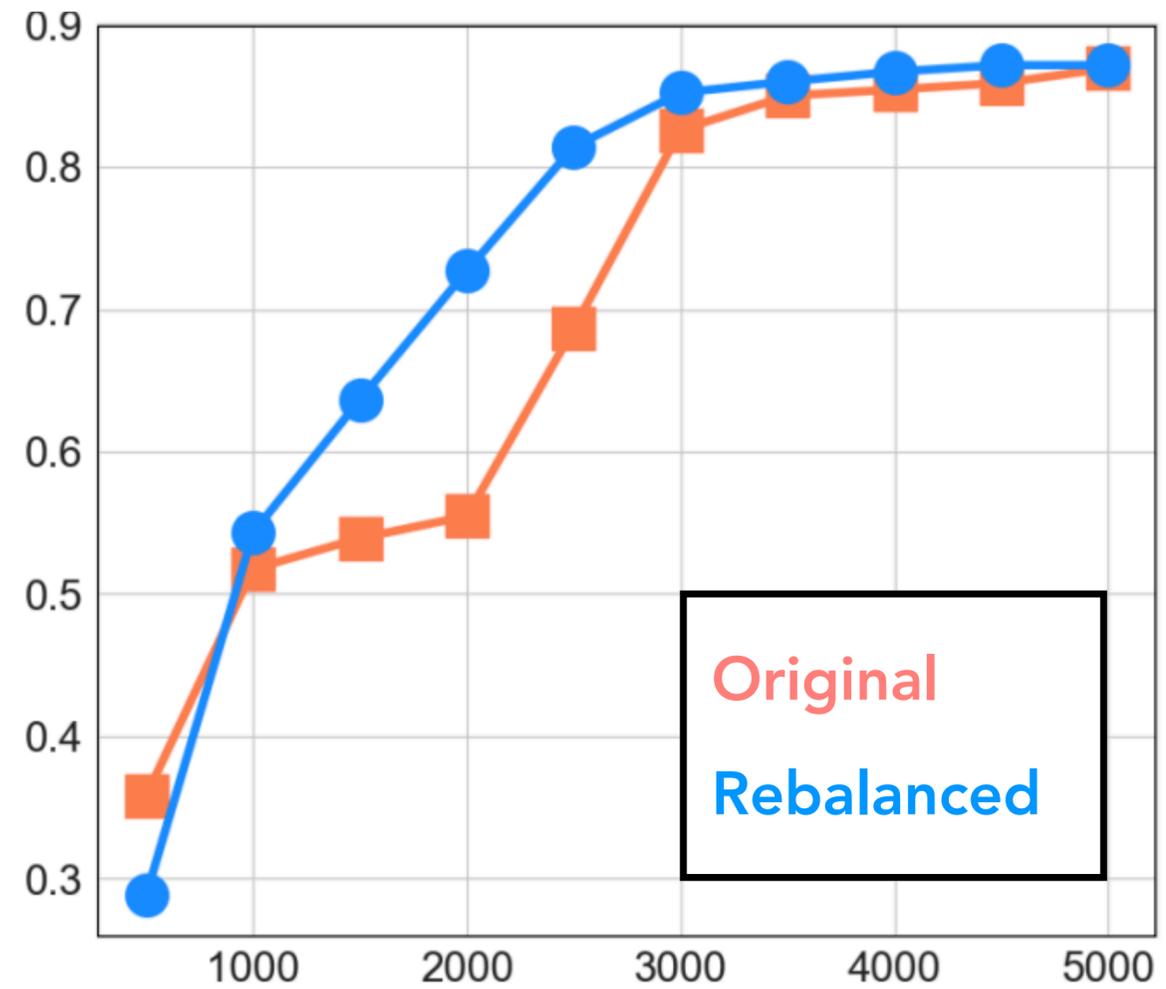
Zero-Shot Accuracy of Models with Different Pre-Training Data

Evaluated on CIFAR-100



Training Iterations

Evaluated on STL-10



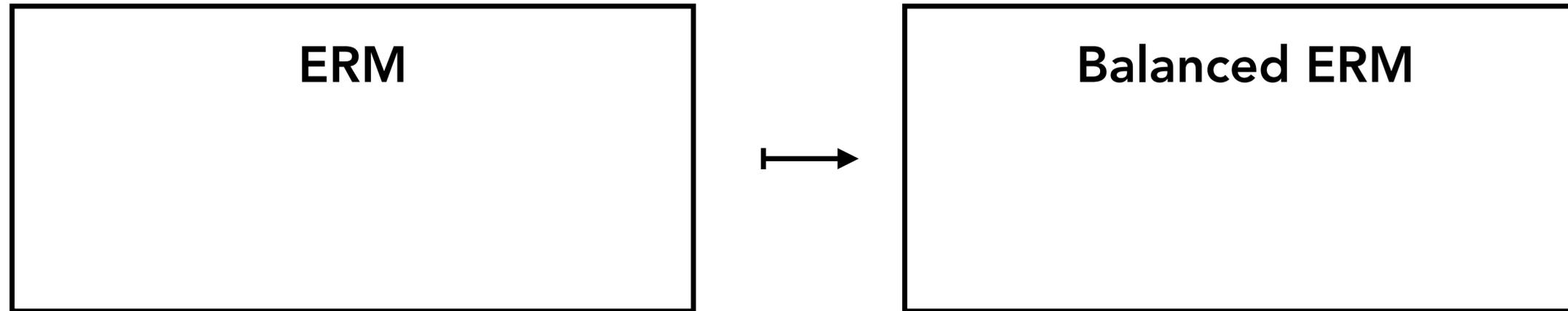
Training Iterations

Original
Rebalanced

Pre-Training Data Curation: Balancing Keyword Distributions

How should we interpret this empirically effective procedure statistically?

Empirical Risk Minimization with **Balancing**

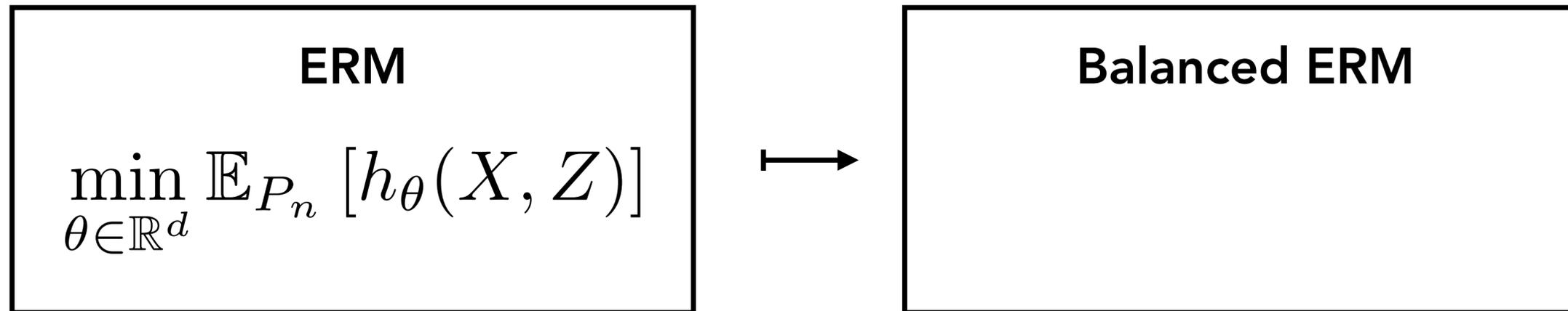


$$(X_1, Z_1), \dots, (X_n, Z_n) \stackrel{\text{i.i.d}}{\sim} P$$

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

$$h_\theta(\mathbf{x}, \mathbf{z}) = \text{loss on pair } (\mathbf{x}, \mathbf{z}) \text{ of model with parameter } \theta$$

Empirical Risk Minimization with **Balancing**

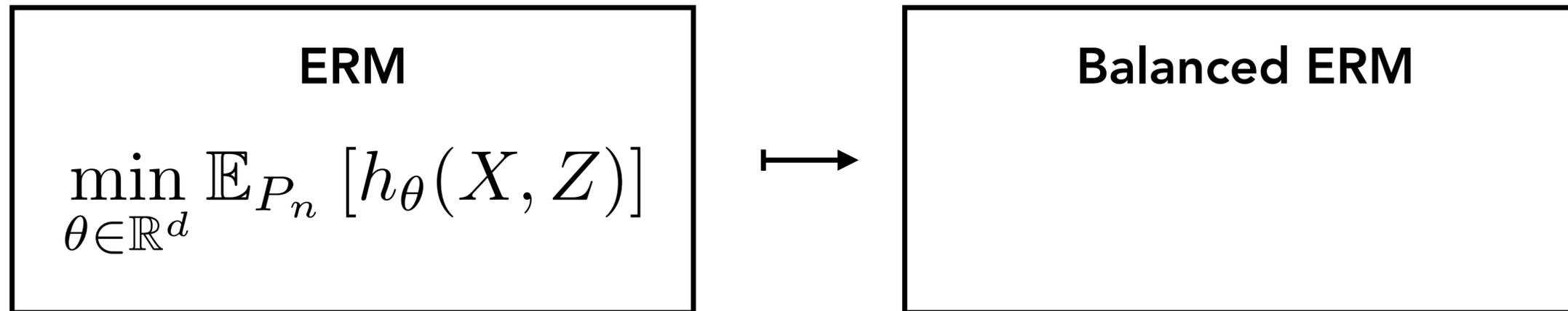


$$(X_1, Z_1), \dots, (X_n, Z_n) \stackrel{\text{i.i.d}}{\sim} P$$

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

$$h_{\theta}(\mathbf{x}, \mathbf{z}) = \text{loss on pair } (\mathbf{x}, \mathbf{z}) \text{ of model with parameter } \theta$$

Empirical Risk Minimization with **Balancing**



$$(X_1, Z_1), \dots, (X_n, Z_n) \stackrel{\text{i.i.d}}{\sim} P$$

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

$$h_{\theta}(\mathbf{x}, \mathbf{z}) = \text{loss on pair } (\mathbf{x}, \mathbf{z}) \text{ of model with parameter } \theta$$

known marginals (P_X, P_Z)

Empirical Risk Minimization with **Balancing**

ERM

$$\min_{\theta \in \mathbb{R}^d} \mathbb{E}_{P_n} [h_{\theta}(X, Z)]$$



Balanced ERM

$$\min_{\theta \in \mathbb{R}^d} \mathbb{E}_{P_n^*} [h_{\theta}(X, Z)]$$

$$(X_1, Z_1), \dots, (X_n, Z_n) \stackrel{\text{i.i.d.}}{\sim} P$$

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

$$h_{\theta}(\mathbf{x}, \mathbf{z}) = \text{loss on pair } (\mathbf{x}, \mathbf{z}) \text{ of model with parameter } \theta$$

known marginals (P_X, P_Z)

P_n^* = A distribution on $\mathcal{X} \times \mathcal{Z}$ that is close to P_n , but agrees with the known marginals.

Empirical Risk Minimization with **Balancing**

ERM

$$\min_{\theta \in \mathbb{R}^d} \mathbb{E}_{P_n} [h_{\theta}(X, Z)]$$



Balanced ERM

$$\min_{\theta \in \mathbb{R}^d} \mathbb{E}_{P_n^*} [h_{\theta}(X, Z)]$$

$$(X_1, Z_1), \dots, (X_n, Z_n) \stackrel{\text{i.i.d.}}{\sim} P$$

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

$$h_{\theta}(\mathbf{x}, \mathbf{z}) = \text{loss on pair } (\mathbf{x}, \mathbf{z}) \text{ of model with parameter } \theta$$

known marginals (P_X, P_Z)

$$P_n^* = \arg \min_{Q \in \text{Coup}(P_X, P_Z)} \text{KL}(Q \| P_n)$$

$$\text{Coup}(P_X, P_Z) = \{Q : Q_X = P_X, Q_Z = P_Z\}$$

How do we perform **balancing** in practice?

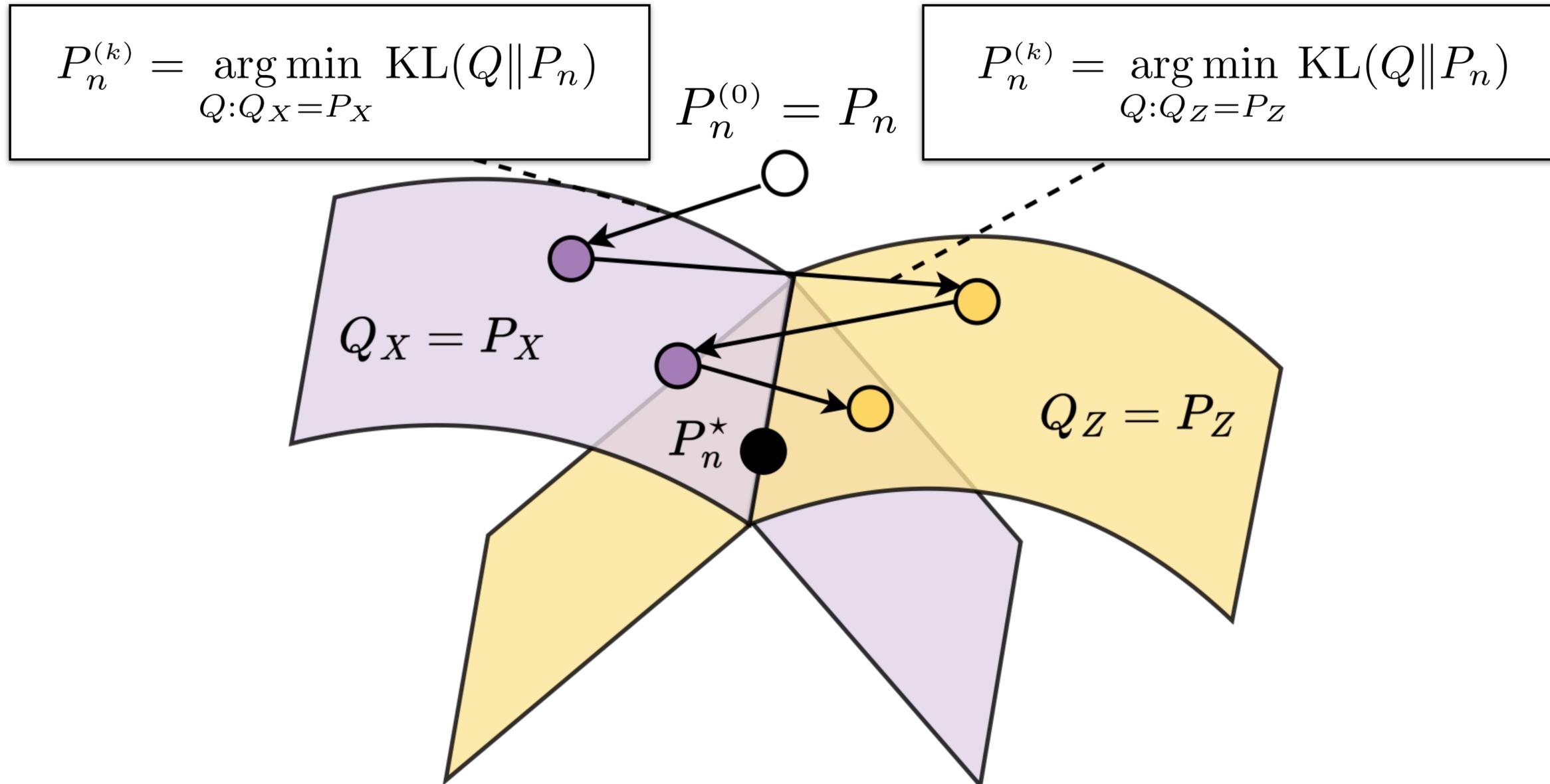
(Sinkhorn Iterations, Iterative Proportional Fitting, Raking Ratio Estimation)

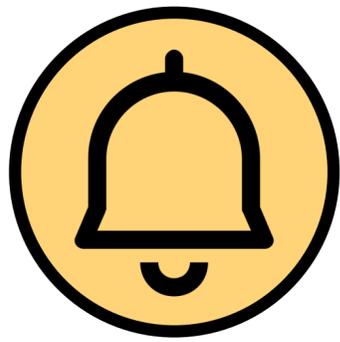
How do we perform **balancing** in practice?

(Sinkhorn Iterations, Iterative Proportional Fitting, Raking Ratio Estimation)

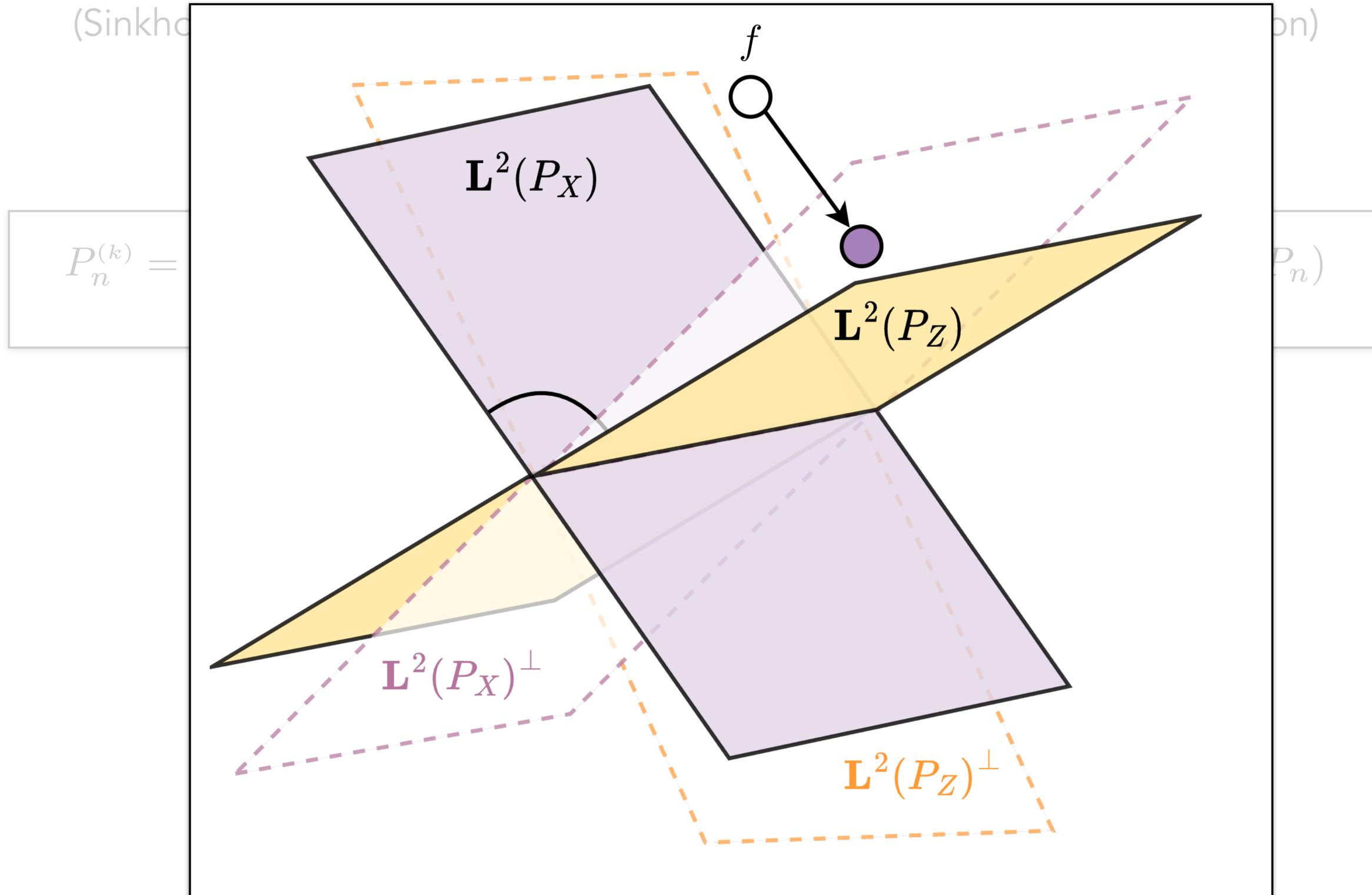
Odd Iterations

Even Iterations





How do we perform **balancing** in practice?



How do we perform **balancing** in practice?

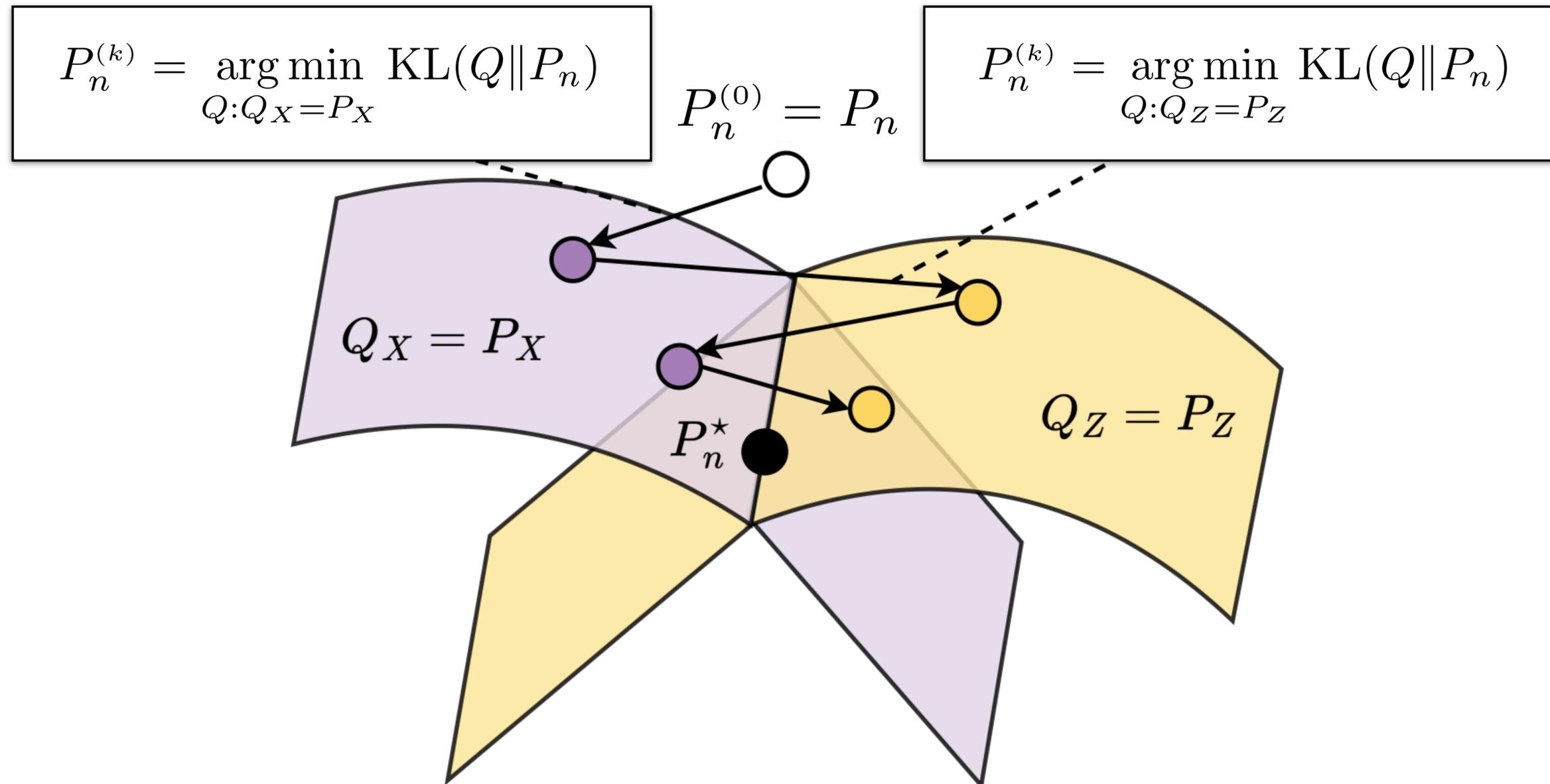
(Sinkhorn Iterations, Iterative Proportional Fitting, Raking Ratio Estimation)

Odd Iterations

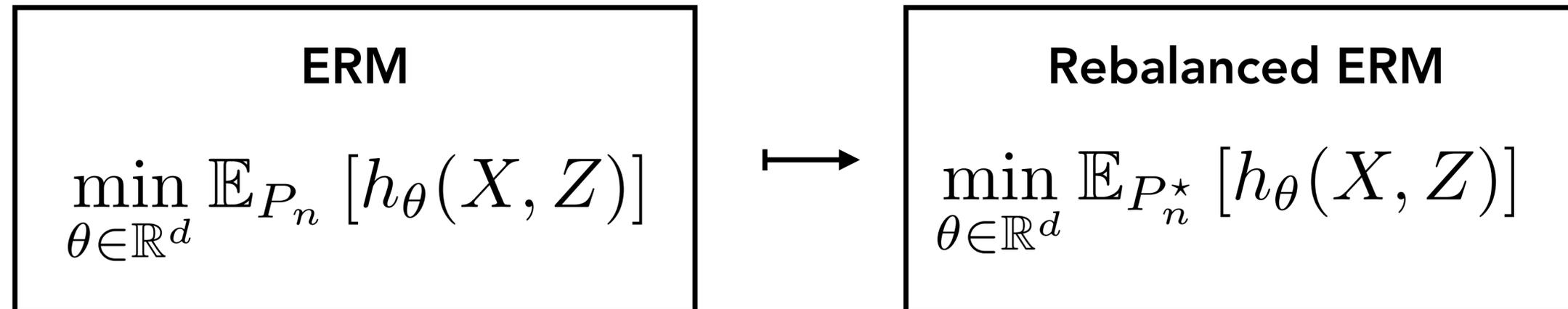
$$P_n^{(k)} = \arg \min_{Q: Q_X = P_X} \text{KL}(Q \| P_n)$$

Even Iterations

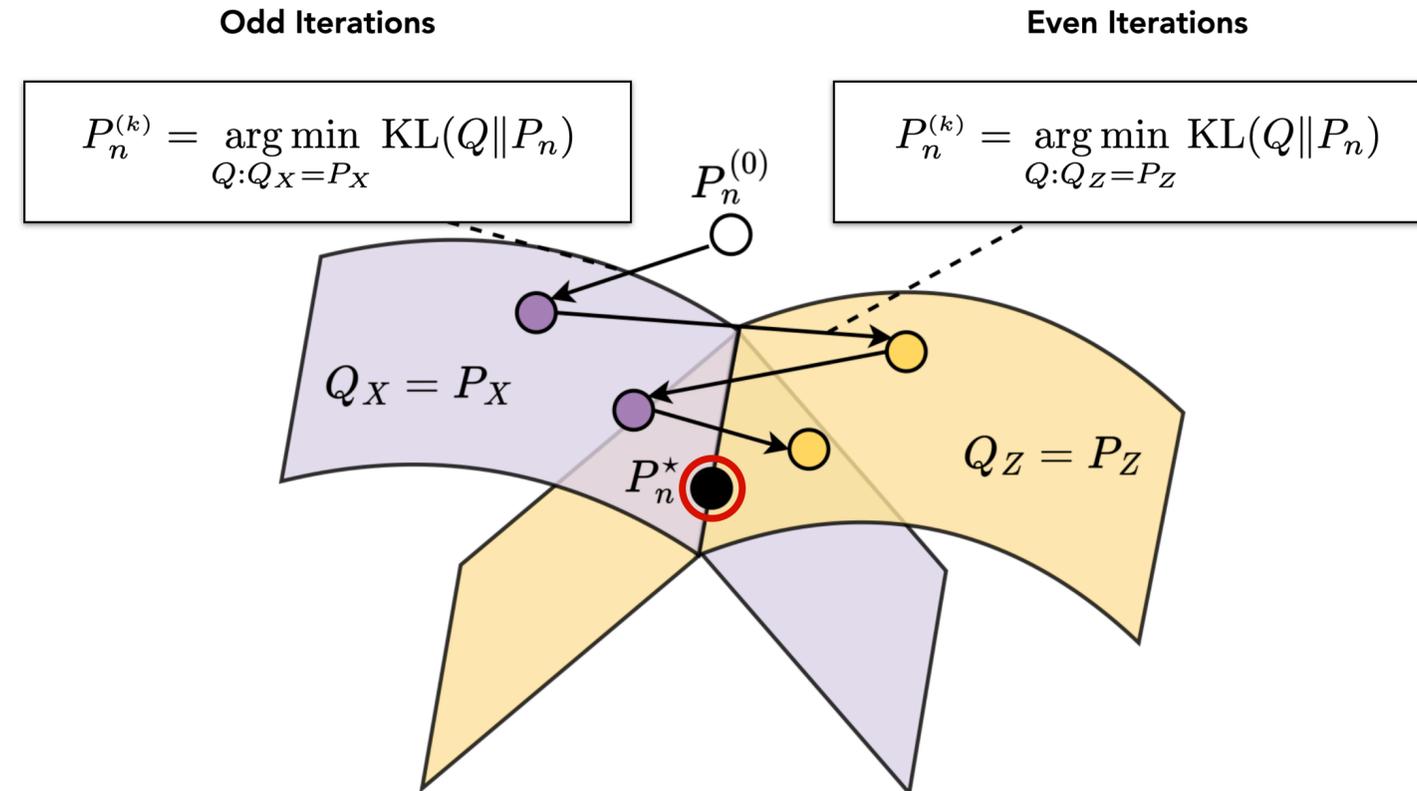
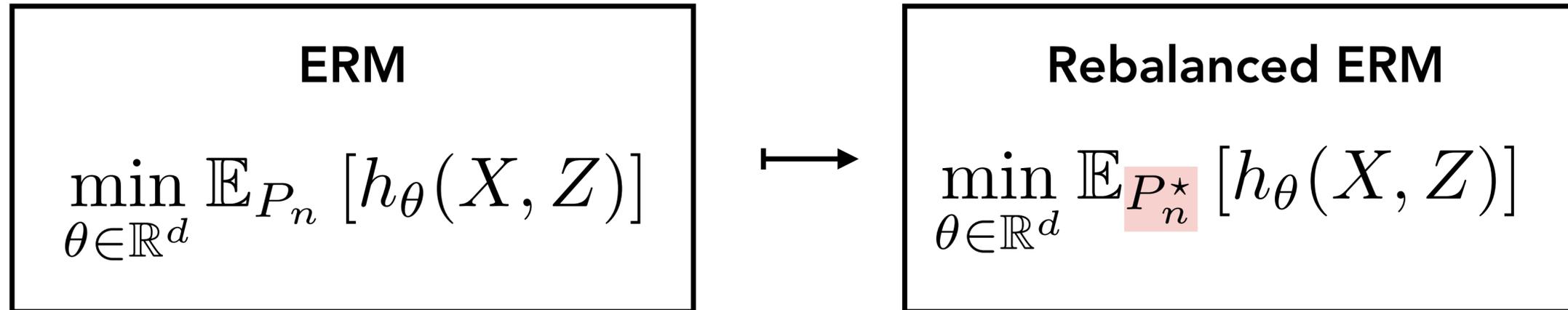
$$P_n^{(k)} = \arg \min_{Q: Q_Z = P_Z} \text{KL}(Q \| P_n)$$



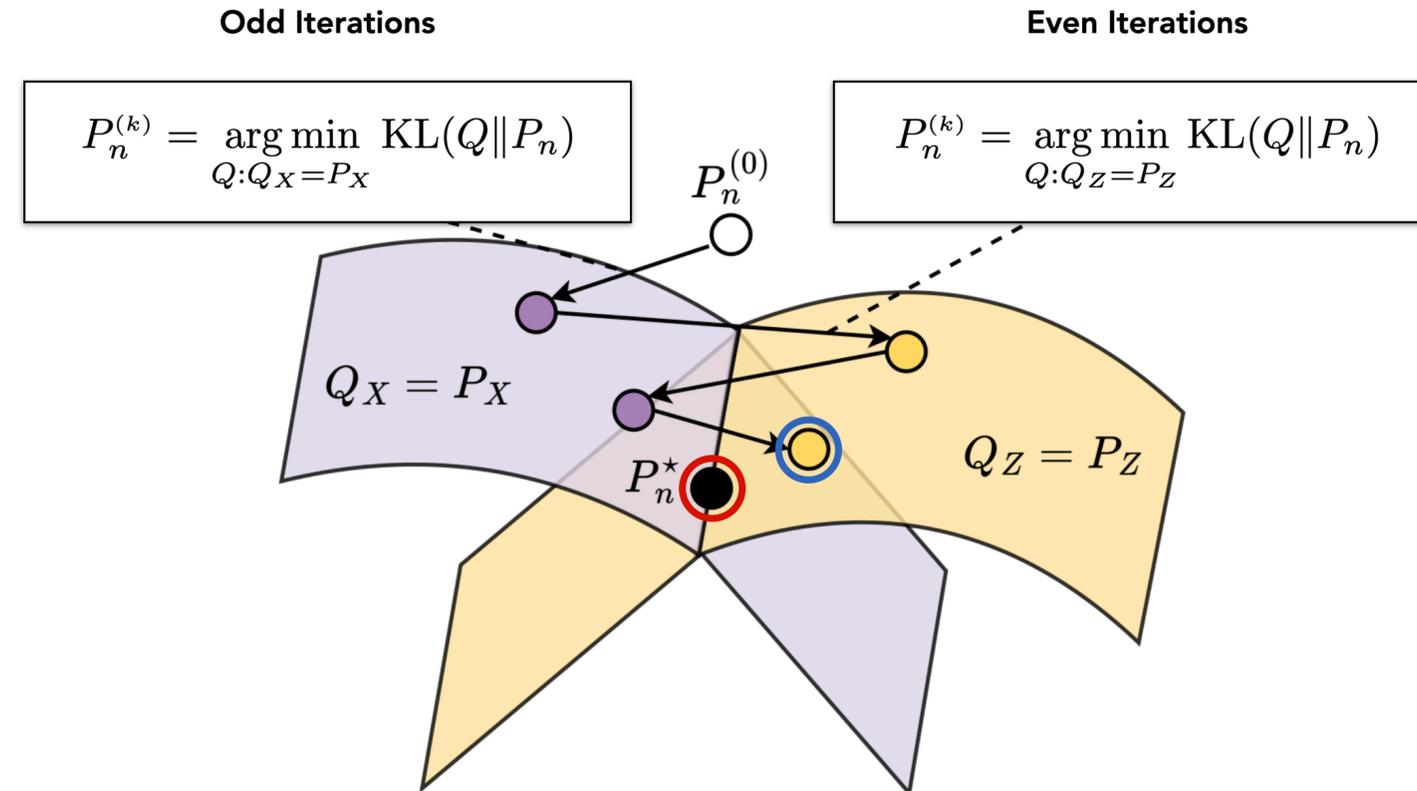
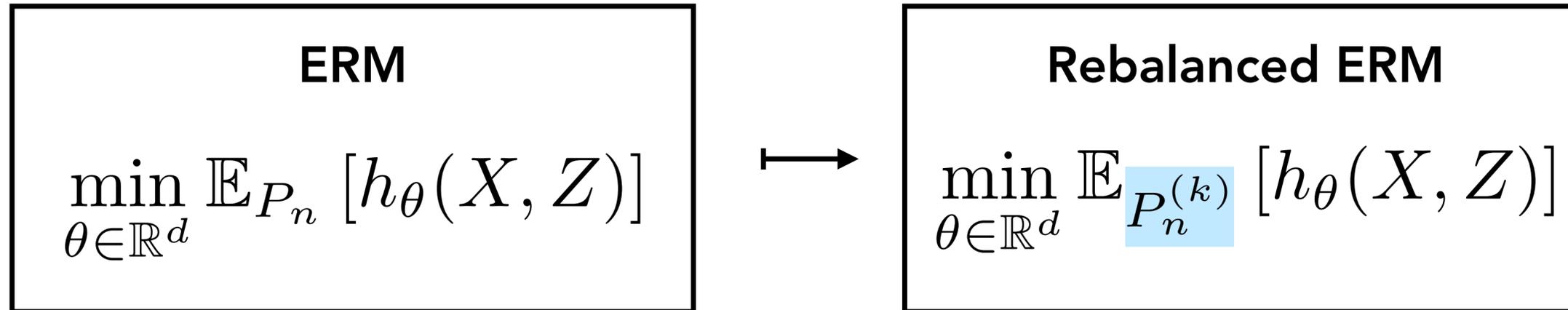
Empirical Risk Minimization with **Balancing**



Empirical Risk Minimization with **Balancing**



Empirical Risk Minimization with **Balancing**



Empirical Risk Minimization with **Balancing**

$$\text{ERM}$$
$$\min_{\theta \in \mathbb{R}^d} \mathbb{E}_{P_n} [h_{\theta}(X, Z)]$$



$$\text{Rebalanced ERM}$$
$$\min_{\theta \in \mathbb{R}^d} \mathbb{E}_{P_n^{(k)}} [h_{\theta}(X, Z)]$$

$$= P_n^{(k)}(h) \stackrel{?}{\approx} P(h)$$

We hide the dependence on θ
and consider point-wise
estimation for a fixed $h \equiv h_{\theta}$.

Empirical Risk Minimization with **Balancing**

ERM

$$\min_{\theta \in \mathbb{R}^d} \mathbb{E}_{P_n} [h_{\theta}(X, Z)]$$



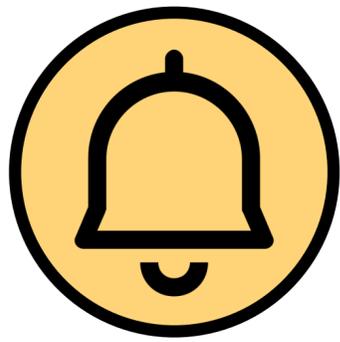
Rebalanced ERM

$$\min_{\theta \in \mathbb{R}^d} \mathbb{E}_{P_n^{(k)}} [h_{\theta}(X, Z)]$$

$$= P_n^{(k)}(h) \stackrel{?}{\approx} P(h)$$

We measure the benefit of balancing via **variance/MSE reduction** for estimating the expectation of a fixed test function.

$$\mathbb{E}_P \left[\left(P_n^{(k)}(h) - P(h) \right)^2 \right] \leq ? < \frac{\text{Var}(h)}{n}$$



Empirical Risk Minimization with **Balancing**

Conditional Mean

$$\mu_{X \leftarrow Z} : \mathbf{L}^2(P) \rightarrow \mathbf{L}^2(P_X)$$

$$[\mu_{X \leftarrow Z} h](\mathbf{x}) = \mathbb{E}_P [h(X, Z) | X](\mathbf{x})$$

Conditional Centering

$$\mu_{X \leftarrow Z}^\perp : \mathbf{L}^2(P) \rightarrow \mathbf{L}^2(P)$$

$$[\mu_{X \leftarrow Z}^\perp h](\mathbf{x}, \mathbf{z}) = h(\mathbf{x}, \mathbf{z}) - [\mu_{X \leftarrow Z} h](\mathbf{x})$$

$$= P_n^{(k)}(h) \approx P(h)$$

We measure the benefit of balancing via **variance/MSE reduction** for estimating the expectation of a fixed test function.

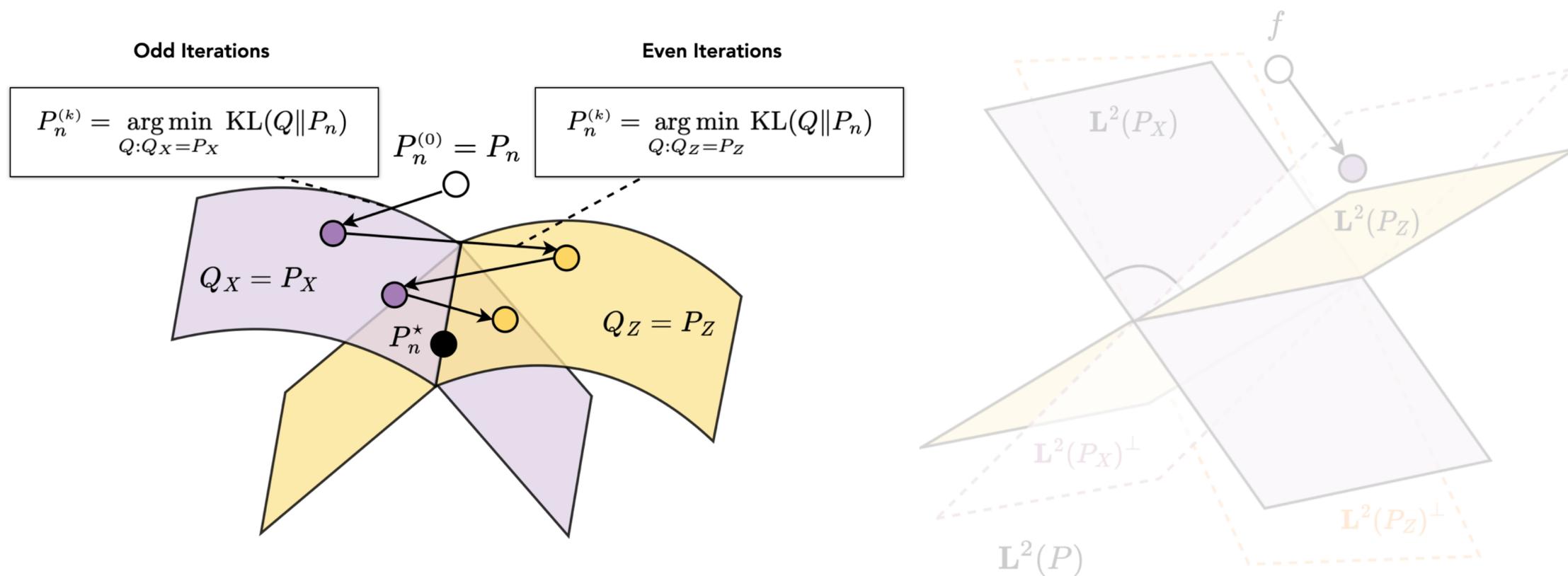
$$\mathbb{E}_P \left[(P_n^{(k)}(h) - P(h))^2 \right] \leq ? < \frac{\text{Var}(h)}{n}$$

Theorem (Liu, M., Pal, Harchaoui)

$$\mathbb{E}_P \left[\left(P_n^{(k)}(h) - P(h) \right)^2 \right] = \frac{\text{Var}(\overbrace{\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp}^{k \text{ times}} h)}{n} + \tilde{O} \left(\frac{k^6}{n^{3/2}} \right)$$

Theorem (Liu, M., Pal, Harchaoui)

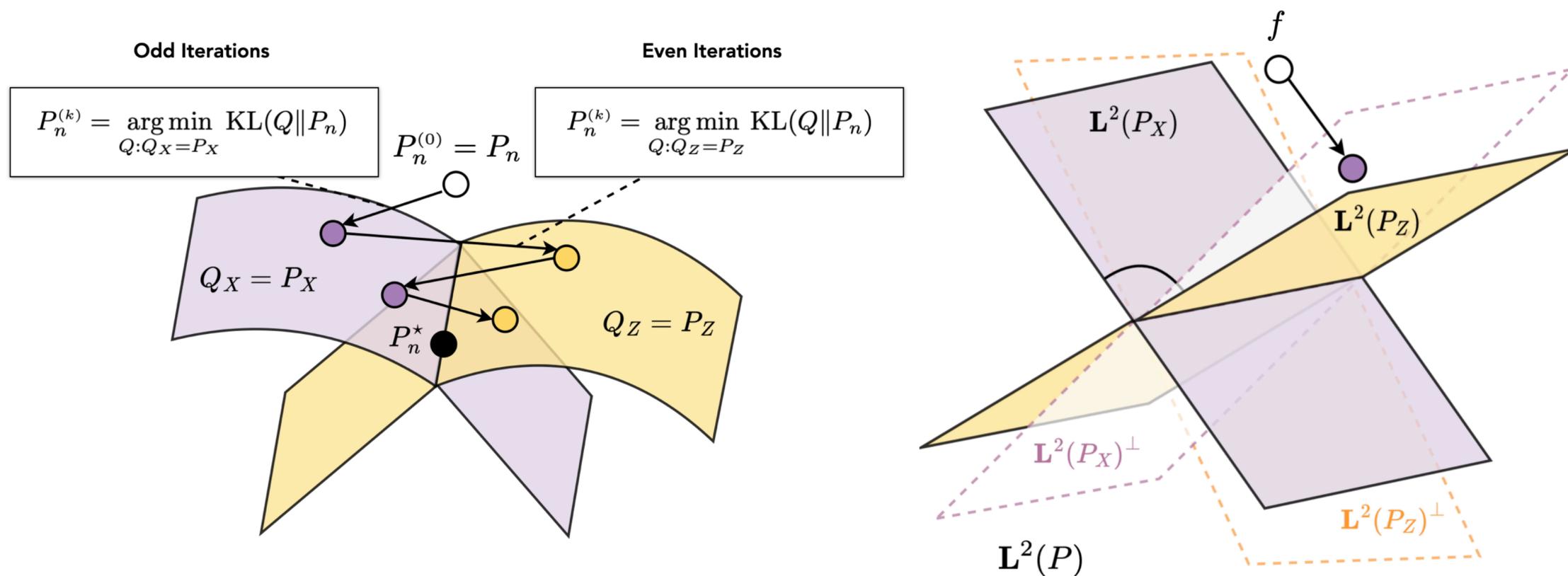
$$\mathbb{E}_P \left[\left(P_n^{(k)}(h) - P(h) \right)^2 \right] = \frac{\text{Var}(\overbrace{\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp}^{k \text{ times}} h)}{n} + \tilde{O} \left(\frac{k^6}{n^{3/2}} \right)$$



Nonlinear Projections \mapsto Linear Projections \mapsto Variance Reduction

Theorem (Liu, M., Pal, Harchaoui)

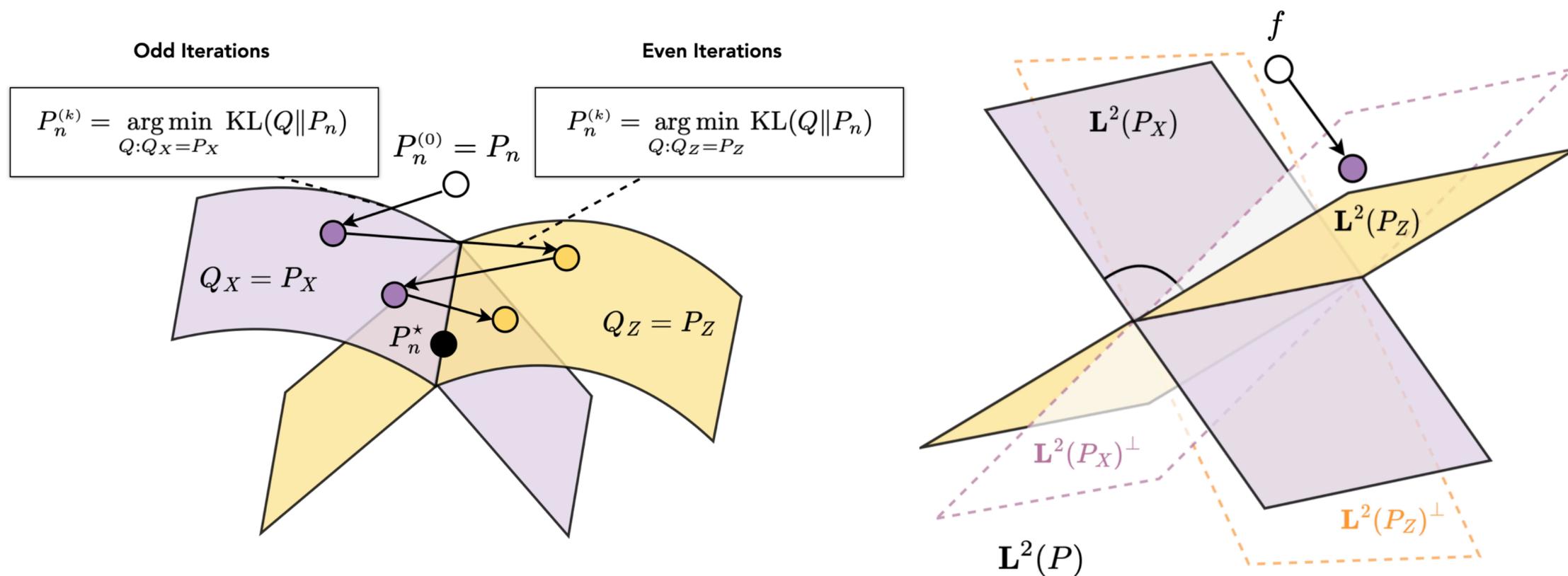
$$\mathbb{E}_P \left[\left(P_n^{(k)}(h) - P(h) \right)^2 \right] = \frac{\text{Var}(\overbrace{\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp}^{k \text{ times}} h)}{n} + \tilde{O} \left(\frac{k^6}{n^{3/2}} \right)$$



Nonlinear Projections \mapsto Linear Projections \mapsto Variance Reduction

Theorem (Liu, M., Pal, Harchaoui)

$$\mathbb{E}_P \left[\left(P_n^{(k)}(h) - P(h) \right)^2 \right] = \frac{\text{Var}(\overbrace{\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp}^{k \text{ times}} h)}{n} + \tilde{O} \left(\frac{k^6}{n^{3/2}} \right)$$



Nonlinear Projections \mapsto Linear Projections \mapsto Variance Reduction

Theorem (Liu, M., Pal, Harchaoui)

$$\mathbb{E}_P \left[\left(P_n^{(k)}(h) - P(h) \right)^2 \right] = \frac{\text{Var}(\overbrace{\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp}^{k \text{ times}} h)}{n} + \tilde{O} \left(\frac{k^6}{n^{3/2}} \right)$$

$$[P_n^{(k)} - P](h) = [P_n^{(k)} - P](\mu_{X \leftarrow Z}^\perp h) + [P_n^{(k)} - P](\mu_{Z \leftarrow X}^\perp h)$$

Theorem (Liu, M., Pal, Harchaoui)

$$\mathbb{E}_P \left[\left(P_n^{(k)}(h) - P(h) \right)^2 \right] = \frac{\text{Var}(\overbrace{\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp}^{k \text{ times}} h)}{n} + \tilde{O} \left(\frac{k^6}{n^{3/2}} \right)$$

$$[P_n^{(k)} - P](h) = [P_n^{(k)} - P](\mu_{X \leftarrow Z}^\perp h) + \overbrace{[P_n^{(k)} - P](\mu_{X \leftarrow Z} h)}{=0}$$

The X -marginal of both distributions match (assuming k odd)

Theorem (Liu, M., Pal, Harchaoui)

$$\mathbb{E}_P \left[\left(P_n^{(k)}(h) - P(h) \right)^2 \right] = \frac{\text{Var}(\overbrace{\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp}^{k \text{ times}} h)}{n} + \tilde{O} \left(\frac{k^6}{n^{3/2}} \right)$$

$$[P_n^{(k)} - P](h) = [P_n^{(k)} - P](\mu_{X \leftarrow Z}^\perp h) + \overbrace{[P_n^{(k)} - P](\mu_{X \leftarrow Z}^\perp h)}{=0}$$

Theorem (Liu, M., Pal, Harchaoui)

$$\mathbb{E}_P \left[(P_n^{(k)}(h) - P(h))^2 \right] = \frac{\text{Var}(\overbrace{\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp}^{k \text{ times}} h)}{n} + \tilde{O} \left(\frac{k^6}{n^{3/2}} \right)$$

$$\begin{aligned} [P_n^{(k)} - P](h) &= [P_n^{(k)} - P](\mu_{X \leftarrow Z}^\perp h) + \overbrace{[P_n^{(k)} - P](\mu_{X \leftarrow Z} h)}{=0} \\ &= [P_n^{(k-1)} - P](\mu_{X \leftarrow Z}^\perp h) + [P_n^{(k)} - P_n^{(k-1)}](\mu_{X \leftarrow Z}^\perp h) \end{aligned}$$

Theorem (Liu, M., Pal, Harchaoui)

$$\mathbb{E}_P \left[(P_n^{(k)}(h) - P(h))^2 \right] = \frac{\text{Var}(\overbrace{\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp}^{k \text{ times}} h)}{n} + \tilde{O} \left(\frac{k^6}{n^{3/2}} \right)$$

$$\begin{aligned} [P_n^{(k)} - P](h) &= [P_n^{(k)} - P](\mu_{X \leftarrow Z}^\perp h) + \overbrace{[P_n^{(k)} - P](\mu_{X \leftarrow Z} h)}{=0} \\ &= \underbrace{[P_n^{(k-1)} - P](\mu_{X \leftarrow Z}^\perp h)}_{\text{Like LHS, but with a conditional centering}} + \underbrace{[P_n^{(k)} - P_n^{(k-1)}](\mu_{X \leftarrow Z}^\perp h)}_{\text{error term}} \end{aligned}$$

Like LHS, but with a conditional centering

error term

Theorem (Liu, M., Pal, Harchaoui)

$$\mathbb{E}_P \left[\left(P_n^{(k)}(h) - P(h) \right)^2 \right] = \frac{\text{Var}(\overbrace{\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp}^{k \text{ times}} h)}{n} + \tilde{O} \left(\frac{k^6}{n^{3/2}} \right)$$

$$\begin{aligned} [P_n^{(k)} - P](h) &= [P_n^{(k)} - P](\mu_{X \leftarrow Z}^\perp h) + \overbrace{[P_n^{(k)} - P](\mu_{X \leftarrow Z} h)}{=0} \\ &= [P_n^{(k-1)} - P](\mu_{X \leftarrow Z}^\perp h) + [P_n^{(k)} - P_n^{(k-1)}](\mu_{X \leftarrow Z}^\perp h) \\ \text{unroll } k \text{ times } \curvearrowright &= [P_n^{(0)} - P](\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp h) + \sum_{\ell=1}^k [P_n^{(\ell)} - P_n^{(\ell-1)}](\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp h) \end{aligned}$$

Theorem (Liu, M., Pal, Harchaoui)

$$\mathbb{E}_P \left[\left(P_n^{(k)}(h) - P(h) \right)^2 \right] = \frac{\text{Var}(\overbrace{\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp}^{k \text{ times}} h)}{n} + \tilde{O} \left(\frac{k^6}{n^{3/2}} \right)$$

$$\begin{aligned} [P_n^{(k)} - P](h) &= [P_n^{(k)} - P](\mu_{X \leftarrow Z}^\perp h) + \overbrace{[P_n^{(k)} - P](\mu_{X \leftarrow Z} h)}{=0} \\ &= [P_n^{(k-1)} - P](\mu_{X \leftarrow Z}^\perp h) + [P_n^{(k)} - P_n^{(k-1)}](\mu_{X \leftarrow Z}^\perp h) \\ &= [P_n^{(0)} - P](\underbrace{\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp}_{k \text{ times}} h) + \sum_{\ell=1}^k [P_n^{(\ell)} - P_n^{(\ell-1)}](\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp h) \end{aligned}$$

Theorem (Liu, M., Pal, Harchaoui)

$$\mathbb{E}_P \left[(P_n^{(k)}(h) - P(h))^2 \right] = \frac{\text{Var}(\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp h)}{n} + \tilde{O} \left(\frac{k^6}{n^{3/2}} \right)$$

$\overbrace{\hspace{15em}}^{k \text{ times}}$

$$\begin{aligned} [P_n^{(k)} - P](h) &= [P_n^{(k)} - P](\mu_{X \leftarrow Z}^\perp h) + \overbrace{[P_n^{(k)} - P](\mu_{X \leftarrow Z} h)}{=0} \\ &= [P_n^{(k-1)} - P](\mu_{X \leftarrow Z}^\perp h) + [P_n^{(k)} - P_n^{(k-1)}](\mu_{X \leftarrow Z}^\perp h) \\ &= \underbrace{[P_n^{(0)} - P](\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp h)}_{k \text{ times}} + \sum_{\ell=1}^k [P_n^{(\ell)} - P_n^{(\ell-1)}](\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp h) \end{aligned}$$

Theorem (Liu, M., Pal, Harchaoui)

$$\mathbb{E}_P \left[(P_n^{(k)}(h) - P(h))^2 \right] = \frac{\text{Var}(\overbrace{\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp}^{k \text{ times}} h)}{n} + \tilde{O} \left(\frac{k^6}{n^{3/2}} \right)$$

$$\begin{aligned} [P_n^{(k)} - P](h) &= [P_n^{(k)} - P](\mu_{X \leftarrow Z}^\perp h) + \overbrace{[P_n^{(k)} - P](\mu_{X \leftarrow Z} h)}{=0} \\ &= [P_n^{(k-1)} - P](\mu_{X \leftarrow Z}^\perp h) + [P_n^{(k)} - P_n^{(k-1)}](\mu_{X \leftarrow Z}^\perp h) \\ &= [P_n^{(0)} - P](\underbrace{\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp}_{k \text{ times}} h) + \sum_{\ell=1}^k [P_n^{(\ell)} - P_n^{(\ell-1)}](\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp h) \end{aligned}$$

Theorem (Liu, M., Pal, Harchaoui)

$$\mathbb{E}_P \left[(P_n^{(k)}(h) - P(h))^2 \right] = \frac{\text{Var}(\overbrace{\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp}^{k \text{ times}} h)}{n} + \tilde{O} \left(\frac{k^6}{n^{3/2}} \right)$$

$$\begin{aligned} [P_n^{(k)} - P](h) &= [P_n^{(k)} - P](\mu_{X \leftarrow Z}^\perp h) + \overbrace{[P_n^{(k)} - P](\mu_{X \leftarrow Z} h)}{=0} \\ &= [P_n^{(k-1)} - P](\mu_{X \leftarrow Z}^\perp h) + [P_n^{(k)} - P_n^{(k-1)}](\mu_{X \leftarrow Z}^\perp h) \\ &= [P_n^{(0)} - P](\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp h) + \sum_{\ell=1}^k [P_n^{(\ell)} - P_n^{(\ell-1)}](\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp h) \end{aligned}$$

$$[P_n^{(\ell)} - P_n^{(\ell-1)}](\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp h) = \sum_{\mathbf{x}, \mathbf{z}} \underbrace{\left(\frac{P_X(\mathbf{x})}{P_{n,X}^{(\ell-1)}(\mathbf{x})} - 1 \right)}_{\substack{\text{unmatched marginal} \\ \text{error } O(n^{-1/2})}} [\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp h](\mathbf{x}, \mathbf{z}) P_n^{(\ell-1)}(\mathbf{x}, \mathbf{z})$$

Theorem (Liu, M., Pal, Harchaoui)

$$\mathbb{E}_P \left[(P_n^{(k)}(h) - P(h))^2 \right] = \frac{\text{Var}(\overbrace{\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp}^{k \text{ times}} h)}{n} + \tilde{O} \left(\frac{k^6}{n^{3/2}} \right)$$

$$\begin{aligned} [P_n^{(k)} - P](h) &= [P_n^{(k)} - P](\mu_{X \leftarrow Z}^\perp h) + \overbrace{[P_n^{(k)} - P](\mu_{X \leftarrow Z} h)}{=0} \\ &= [P_n^{(k-1)} - P](\mu_{X \leftarrow Z}^\perp h) + [P_n^{(k)} - P_n^{(k-1)}](\mu_{X \leftarrow Z}^\perp h) \\ &= [P_n^{(0)} - P](\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp h) + \sum_{\ell=1}^k [P_n^{(\ell)} - P_n^{(\ell-1)}](\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp h) \end{aligned}$$

$$[P_n^{(\ell)} - P_n^{(\ell-1)}](\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp h) = \sum_{\mathbf{x}, \mathbf{z}} \left(\frac{P_X(\mathbf{x})}{P_{n,X}^{(\ell-1)}(\mathbf{x})} - 1 \right) \underbrace{[\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp h](\mathbf{x}, \mathbf{z}) P_n^{(\ell-1)}(\mathbf{x}, \mathbf{z})}_{\text{control using similar recursion } O(n^{-1/2})}$$

control using similar
recursion $O(n^{-1/2})$

Theorem (Liu, M., Pal, Harchaoui)

$$\mathbb{E}_P \left[\left(P_n^{(k)}(h) - P(h) \right)^2 \right] = \frac{\text{Var}(\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp h)}{n} + \tilde{O} \left(\frac{k^6}{n^{3/2}} \right)$$

$$\begin{aligned} [P_n^{(k)} - P](h) &= [P_n^{(k)} - P](\mu_{X \leftarrow Z}^\perp h) + \overbrace{[P_n^{(k)} - P](\mu_{X \leftarrow Z}^\perp h)}^{=0} \\ &= [P_n^{(k-1)} - P](\mu_{X \leftarrow Z}^\perp h) + [P_n^{(k)} - P_n^{(k-1)}](\mu_{X \leftarrow Z}^\perp h) \\ &= [P_n^{(0)} - P](\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp h) + \sum_{\ell=1}^k [P_n^{(\ell)} - P_n^{(\ell-1)}](\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp h) \end{aligned}$$

What is the actual
reduction in variance?

Theorem (Liu, M., Pal, Harchaoui)

$$\mathbb{E}_P \left[\left(P_n^{(k)}(h) - P(h) \right)^2 \right] = \frac{\text{Var}(\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp h)}{n} + \tilde{O} \left(\frac{k^6}{n^{3/2}} \right)$$

$$\begin{aligned} [P_n^{(k)} - P](h) &= [P_n^{(k)} - P](\mu_{X \leftarrow Z}^\perp h) + \overbrace{[P_n^{(k)} - P](\mu_{X \leftarrow Z} h)}{=0} \\ &= [P_n^{(k-1)} - P](\mu_{X \leftarrow Z}^\perp h) + [P_n^{(k)} - P_n^{(k-1)}](\mu_{X \leftarrow Z}^\perp h) \\ &= [P_n^{(0)} - P](\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp h) + \sum_{\ell=1}^k [P_n^{(\ell)} - P_n^{(\ell-1)}](\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp h) \end{aligned}$$

$$\mu_{X \leftarrow Z}^\perp = I - \mu_{X \leftarrow Z}$$

$$\mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp = I - \mu_{X \leftarrow Z} - \mu_{Z \leftarrow X} + \mu_{Z \leftarrow X} \mu_{X \leftarrow Z}$$

$$\mu_{X \leftarrow Z}^\perp \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp = I - \mu_{X \leftarrow Z} - \mu_{Z \leftarrow X} + \mu_{Z \leftarrow X} \mu_{X \leftarrow Z} + \mu_{X \leftarrow Z} \mu_{Z \leftarrow X} - \mu_{X \leftarrow Z} \mu_{Z \leftarrow X} \mu_{X \leftarrow Z}$$

Theorem (Liu, M., Pal, Harchaoui)

$$\mathbb{E}_P \left[\left(P_n^{(k)}(h) - P(h) \right)^2 \right] = \frac{\text{Var}(\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp h)}{n} + \tilde{O} \left(\frac{k^6}{n^{3/2}} \right)$$

$$\begin{aligned} [P_n^{(k)} - P](h) &= [P_n^{(k)} - P](\mu_{X \leftarrow Z}^\perp h) + \overbrace{[P_n^{(k)} - P](\mu_{X \leftarrow Z} h)}{=0} \\ &= [P_n^{(k-1)} - P](\mu_{X \leftarrow Z}^\perp h) + [P_n^{(k)} - P_n^{(k-1)}](\mu_{X \leftarrow Z}^\perp h) \\ &= [P_n^{(0)} - P](\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp h) + \sum_{\ell=1}^k [P_n^{(\ell)} - P_n^{(\ell-1)}](\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp h) \end{aligned}$$

$$\mu_{X \leftarrow Z}^\perp = I - \mu_{X \leftarrow Z}$$

$$\mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp = I - \mu_{X \leftarrow Z} - \mu_{Z \leftarrow X} + \mu_{Z \leftarrow X} \mu_{X \leftarrow Z}$$

$$\mu_{X \leftarrow Z}^\perp \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp = I - \mu_{X \leftarrow Z} - \mu_{Z \leftarrow X} + \mu_{Z \leftarrow X} \mu_{X \leftarrow Z} + \mu_{X \leftarrow Z} \mu_{Z \leftarrow X} - \mu_{X \leftarrow Z} \mu_{Z \leftarrow X} \mu_{X \leftarrow Z}$$

Theorem (Liu, M., Pal, Harchaoui)

$$\mathbb{E}_P \left[\left(P_n^{(k)}(h) - P(h) \right)^2 \right] = \frac{\text{Var}(\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp h)}{n} + \tilde{O} \left(\frac{k^6}{n^{3/2}} \right)$$

$$\begin{aligned} [P_n^{(k)} - P](h) &= [P_n^{(k)} - P](\mu_{X \leftarrow Z}^\perp h) + \overbrace{[P_n^{(k)} - P](\mu_{X \leftarrow Z} h)}{=0} \\ &= [P_n^{(k-1)} - P](\mu_{X \leftarrow Z}^\perp h) + [P_n^{(k)} - P_n^{(k-1)}](\mu_{X \leftarrow Z}^\perp h) \\ &= [P_n^{(0)} - P](\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp h) + \sum_{\ell=1}^k [P_n^{(\ell)} - P_n^{(\ell-1)}](\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp h) \end{aligned}$$

$$\mu_{X \leftarrow Z}^\perp = I - \mu_{X \leftarrow Z}$$

$$\mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp = I - \mu_{X \leftarrow Z} - \mu_{Z \leftarrow X} + \mu_{Z \leftarrow X} \mu_{X \leftarrow Z}$$

$$\mu_{X \leftarrow Z}^\perp \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp = I - \mu_{X \leftarrow Z} - \mu_{Z \leftarrow X} + \mu_{Z \leftarrow X} \mu_{X \leftarrow Z} + \mu_{X \leftarrow Z} \mu_{Z \leftarrow X} - \mu_{X \leftarrow Z} \mu_{Z \leftarrow X} \mu_{X \leftarrow Z}$$

Theorem (Liu, M., Pal, Harchaoui)

$$\mathbb{E}_P \left[\left(P_n^{(k)}(h) - P(h) \right)^2 \right] = \frac{\text{Var}(\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp h)}{n} + \tilde{O} \left(\frac{k^6}{n^{3/2}} \right)$$

$$\begin{aligned} [P_n^{(k)} - P](h) &= [P_n^{(k)} - P](\mu_{X \leftarrow Z}^\perp h) + \overbrace{[P_n^{(k)} - P](\mu_{X \leftarrow Z} h)}^{=0} \\ &= [P_n^{(k-1)} - P](\mu_{X \leftarrow Z}^\perp h) + [P_n^{(k)} - P_n^{(k-1)}](\mu_{X \leftarrow Z}^\perp h) \\ &= [P_n^{(0)} - P](\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp h) + \sum_{\ell=1}^k [P_n^{(\ell)} - P_n^{(\ell-1)}](\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp h) \end{aligned}$$

$$\mu_{X \leftarrow Z}^\perp = I - \mu_{X \leftarrow Z}$$

$$\mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp = I - \mu_{X \leftarrow Z} - \mu_{Z \leftarrow X} + \mu_{Z \leftarrow X} \mu_{X \leftarrow Z}$$

$$\mu_{X \leftarrow Z}^\perp \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp = I - \mu_{X \leftarrow Z} - \mu_{Z \leftarrow X} + \mu_{Z \leftarrow X} \mu_{X \leftarrow Z} + \mu_{X \leftarrow Z} \mu_{Z \leftarrow X} - \mu_{X \leftarrow Z} \mu_{Z \leftarrow X} \mu_{X \leftarrow Z}$$

Theorem (Liu, M., Pal, Harchaoui)

$$\mathbb{E}_P \left[\left(P_n^{(k)}(h) - P(h) \right)^2 \right] = \frac{\text{Var}(\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp h)}{n} + \tilde{O} \left(\frac{k^6}{n^{3/2}} \right)$$

$$\begin{aligned} [P_n^{(k)} - P](h) &= [P_n^{(k)} - P](\mu_{X \leftarrow Z}^\perp h) + \overbrace{[P_n^{(k)} - P](\mu_{X \leftarrow Z} h)}{=0} \\ &= [P_n^{(k-1)} - P](\mu_{X \leftarrow Z}^\perp h) + [P_n^{(k)} - P_n^{(k-1)}](\mu_{X \leftarrow Z}^\perp h) \\ &= [P_n^{(0)} - P](\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp h) + \sum_{\ell=1}^k [P_n^{(\ell)} - P_n^{(\ell-1)}](\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp h) \end{aligned}$$

$$\mu_{X \leftarrow Z}^\perp = I - \mu_{X \leftarrow Z}$$

$$\mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp = I - \mu_{X \leftarrow Z} - \mu_{Z \leftarrow X} + \mu_{Z \leftarrow X} \mu_{X \leftarrow Z} \quad \mathbb{E}_{P_{X,Z}} [h(X, Z) | X] (\cdot)$$

$$\mu_{X \leftarrow Z}^\perp \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp = I - \mu_{X \leftarrow Z} - \mu_{Z \leftarrow X} + \mu_{Z \leftarrow X} \mu_{X \leftarrow Z} + \mu_{X \leftarrow Z} \mu_{Z \leftarrow X} - \underbrace{\mu_{X \leftarrow Z} \mu_{Z \leftarrow X} \mu_{X \leftarrow Z}}_h h$$



Theorem (Liu, M., Pal, Harchaoui)

\mathbb{E}

Old Friend: Singular Value Decomposition

Basis of $\mathbf{L}^2(P_X)$: $\alpha_1, \alpha_2, \dots$

$$\mu_{X \leftarrow Z} \beta_i = s_i \alpha_i$$

Basis of $\mathbf{L}^2(P_Z)$: β_1, β_2, \dots

$$\mu_{Z \leftarrow X} \alpha_i = s_i \beta_i$$

$[P_n^{(k)}$

$$= [P_n^{(0)} - P](\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp h) + \sum_{\ell=1}^k [P_n^{(\ell)} - P_n^{(\ell-1)}](\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp h)$$

$$\mu_{X \leftarrow Z}^\perp = I - \mu_{X \leftarrow Z}$$

$$\mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp = I - \mu_{X \leftarrow Z} - \mu_{Z \leftarrow X} + \mu_{Z \leftarrow X} \mu_{X \leftarrow Z}$$

$\mathbb{E}_{P_{X,Z}} [h(X, Z) | X] (\cdot)$

$$\mu_{X \leftarrow Z}^\perp \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp = I - \mu_{X \leftarrow Z} - \mu_{Z \leftarrow X} + \mu_{Z \leftarrow X} \mu_{X \leftarrow Z} + \mu_{X \leftarrow Z} \mu_{Z \leftarrow X} - \underbrace{\mu_{X \leftarrow Z} \mu_{Z \leftarrow X} \mu_{X \leftarrow Z}}_h$$



Theorem (Liu, M., Pal, Harchaoui)

\mathbb{E}

Old Friend: Singular Value Decomposition

Basis of $\mathbf{L}^2(P_X)$: $\alpha_1, \alpha_2, \dots$

$$\mu_{X \leftarrow Z} \beta_i = s_i \alpha_i$$

Basis of $\mathbf{L}^2(P_Z)$: β_1, β_2, \dots

$$\mu_{Z \leftarrow X} \alpha_i = s_i \beta_i$$

$[P_n^{(k)}$

$$= [P_n^{(0)} - P](\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp h) + \sum_{\ell=1}^k [P_n^{(\ell)} - P_n^{(\ell-1)}](\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp h)$$

$$\mu_{X \leftarrow Z}^\perp = I - \mu_{X \leftarrow Z}$$

$$\mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp = I - \mu_{X \leftarrow Z} - \mu_{Z \leftarrow X} + \mu_{Z \leftarrow X} \mu_{X \leftarrow Z}$$

$$\mu_{X \leftarrow Z}^\perp \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp = I - \mu_{X \leftarrow Z} - \mu_{Z \leftarrow X} + \mu_{Z \leftarrow X} \mu_{X \leftarrow Z} + \mu_{X \leftarrow Z} \mu_{Z \leftarrow X} - \underbrace{\mu_{X \leftarrow Z} \mu_{Z \leftarrow X} \mu_{X \leftarrow Z}}_h$$

$\mathbb{E}_{P_{X,Z}} [h(X, Z) | X] (\cdot)$

$$\mu_{X \leftarrow Z} h = \sum_i u_i \alpha_i$$



Theorem (Liu, M., Pal, Harchaoui)

\mathbb{E}

Old Friend: Singular Value Decomposition

Basis of $\mathbf{L}^2(P_X)$: $\alpha_1, \alpha_2, \dots$

$$\mu_{X \leftarrow Z} \beta_i = s_i \alpha_i$$

Basis of $\mathbf{L}^2(P_Z)$: β_1, β_2, \dots

$$\mu_{Z \leftarrow X} \alpha_i = s_i \beta_i$$

$[P_n^{(k)}$

$$= [P_n^{(0)} - P](\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp h) + \sum_{\ell=1}^k [P_n^{(\ell)} - P_n^{(\ell-1)}](\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp h)$$

$$\mu_{X \leftarrow Z}^\perp = I - \mu_{X \leftarrow Z}$$

$$\mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp = I - \mu_{X \leftarrow Z} - \mu_{Z \leftarrow X} + \mu_{Z \leftarrow X} \mu_{X \leftarrow Z}$$

$$\mu_{X \leftarrow Z}^\perp \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp = I - \mu_{X \leftarrow Z} - \mu_{Z \leftarrow X} + \mu_{Z \leftarrow X} \mu_{X \leftarrow Z} + \mu_{X \leftarrow Z} \mu_{Z \leftarrow X} - \underbrace{\mu_{X \leftarrow Z} \mu_{Z \leftarrow X} \mu_{X \leftarrow Z}}_h$$

$$\mu_{X \leftarrow Z} h = \sum_i s_i u_i \beta_i$$



Theorem (Liu, M., Pal, Harchaoui)

\mathbb{E}

Old Friend: Singular Value Decomposition

Basis of $\mathbf{L}^2(P_X)$: $\alpha_1, \alpha_2, \dots$

$$\mu_{X \leftarrow Z} \beta_i = s_i \alpha_i$$

Basis of $\mathbf{L}^2(P_Z)$: β_1, β_2, \dots

$$\mu_{Z \leftarrow X} \alpha_i = s_i \beta_i$$

$[P_n^{(k)}$

$$= [P_n^{(0)} - P](\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp h) + \sum_{\ell=1}^k [P_n^{(\ell)} - P_n^{(\ell-1)}](\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp h)$$

$$\mu_{X \leftarrow Z}^\perp = I - \mu_{X \leftarrow Z}$$

$$\mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp = I - \mu_{X \leftarrow Z} - \mu_{Z \leftarrow X} + \mu_{Z \leftarrow X} \mu_{X \leftarrow Z}$$

$$\mu_{X \leftarrow Z}^\perp \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp = I - \mu_{X \leftarrow Z} - \mu_{Z \leftarrow X} + \mu_{Z \leftarrow X} \mu_{X \leftarrow Z} + \mu_{X \leftarrow Z} \mu_{Z \leftarrow X} - \underbrace{\mu_{X \leftarrow Z} \mu_{Z \leftarrow X} \mu_{X \leftarrow Z}}_h$$

$$\mu_{X \leftarrow Z} h = \sum_i s_i^2 u_i \alpha_i$$



Theorem (Liu, M., Pal, Harchaoui)

\mathbb{E}

Old Friend: Singular Value Decomposition

Basis of $\mathbf{L}^2(P_X)$: $\alpha_1, \alpha_2, \dots$

$$\mu_{X \leftarrow Z} \beta_i = s_i \alpha_i$$

Basis of $\mathbf{L}^2(P_Z)$: β_1, β_2, \dots

$$\mu_{Z \leftarrow X} \alpha_i = s_i \beta_i$$

$[P_n^{(k)}$

$$= [P_n^{(0)} - P](\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp h) + \sum_{\ell=1}^k [P_n^{(\ell)} - P_n^{(\ell-1)}](\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp h)$$

$$\mu_{X \leftarrow Z}^\perp = I - \mu_{X \leftarrow Z}$$

$$\mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp = I - \mu_{X \leftarrow Z} - \mu_{Z \leftarrow X} + \mu_{Z \leftarrow X} \mu_{X \leftarrow Z}$$

$$\mu_{X \leftarrow Z}^\perp \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp = I - \mu_{X \leftarrow Z} - \mu_{Z \leftarrow X} + \mu_{Z \leftarrow X} \mu_{X \leftarrow Z} + \mu_{X \leftarrow Z} \mu_{Z \leftarrow X} - \mu_{X \leftarrow Z} \mu_{Z \leftarrow X} \mu_{X \leftarrow Z} h$$

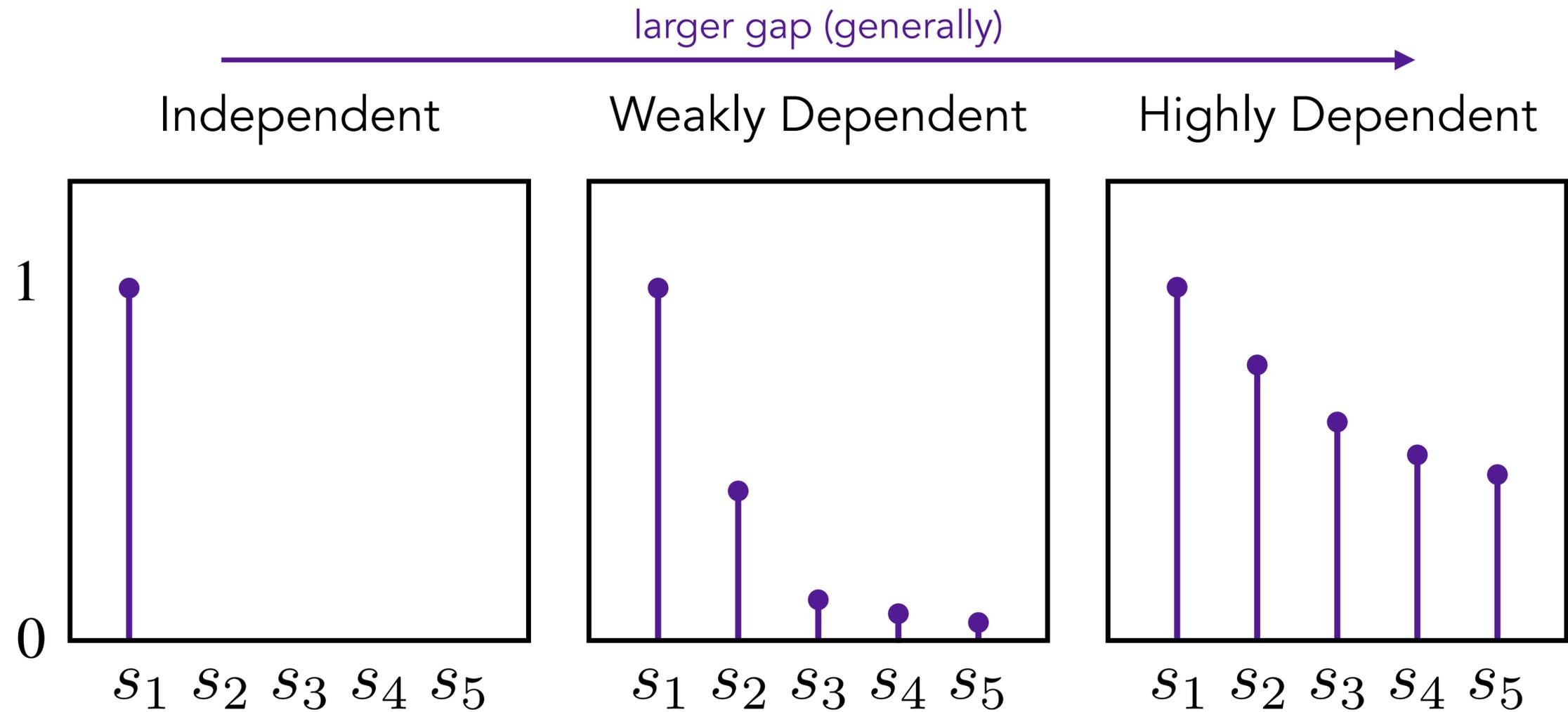
$$\mu_{X \leftarrow Z} h = \sum_i s_i^{10} u_i \alpha_i$$

$$\text{Var}(h) - \text{Var}(\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp h) = \sum_{i=2,3,\dots} \frac{u_i^2 + v_i^2 - 2s_i u_i v_i}{1 - s_i^2}$$

$u_i = \langle \alpha_i, \mu_{X \leftarrow Z} h \rangle$
 $v_i = \langle \beta_i, \mu_{Z \leftarrow X} h \rangle$

$$\text{Var}(h) - \text{Var}(\dots \mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp h) = \sum_{i=2,3,\dots} \frac{u_i^2 + v_i^2 - 2s_i u_i v_i}{1 - s_i^2}$$

$u_i = \langle \alpha_i, \mu_{X \leftarrow Z} h \rangle$
 $v_i = \langle \beta_i, \mu_{Z \leftarrow X} h \rangle$



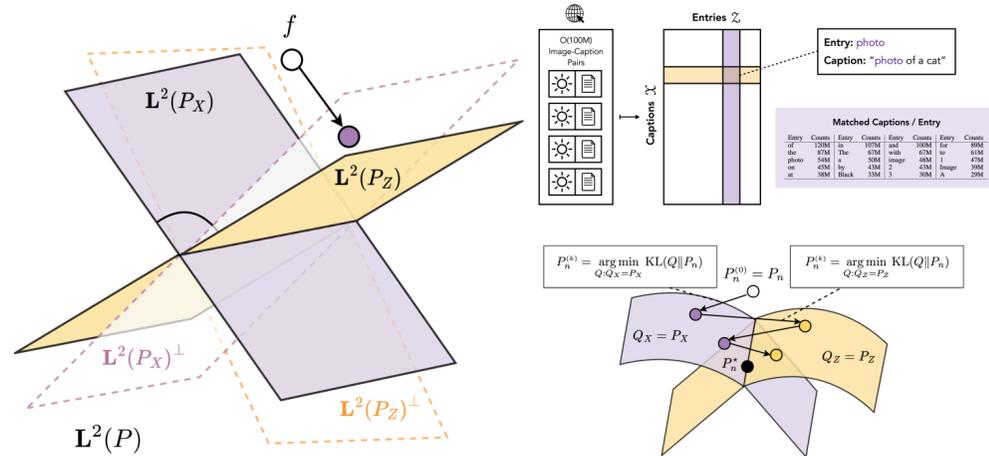
Research Contributions



Exact variance reduction analysis of balancing-based mean estimation!

Theorem (Liu, M., Pal, Harchaoui)

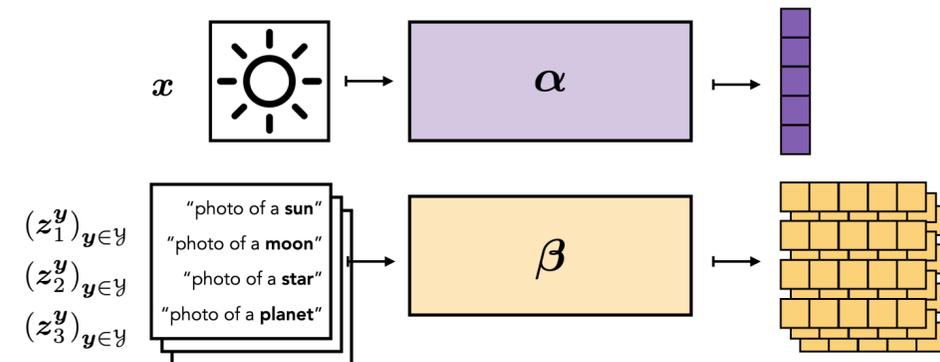
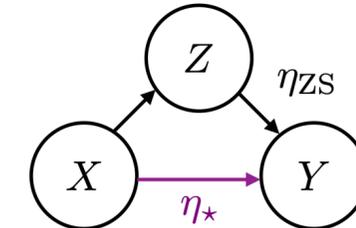
$$\mathbb{E}_P [(P_n^{(k)}(h) - P(h))^2] = \frac{\text{Var}(\dots \overset{k \text{ times}}{\mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp} h)}{n} + \tilde{O}\left(\frac{k^6}{n^{3/2}}\right)$$



Theoretical framework for obtaining generalization guarantees for two broad classes of zero-shot prediction models!

Thm. 1 (M., Harchaoui)

$$\|\eta_* - \eta_{ZS}\|_{\mathbf{L}^2(Q_X)}^2 \lesssim \mathbb{E}_{Q_Z} [I(X, Y|Z)] + \|g_Q - g_R\|_{\mathbf{L}^2(Q_Z)}^2$$



Preliminaries

Balanced Pre-Training

Zero-Shot Prediction

Conclusion



Zaid Harchaoui
University of
Washington

The Prediction Path to Zero-Shot Generalization

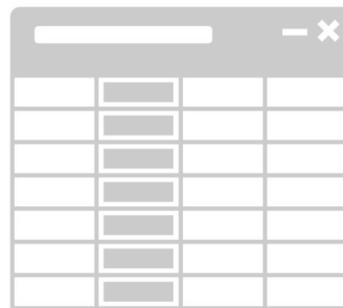
Ronak Mehta¹ Zaid Harchaoui¹

ICML 2025 Spotlight (Top
2.6% of Submissions)

Motivating Questions

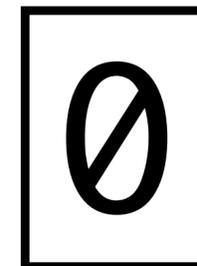
Foundation Modeling

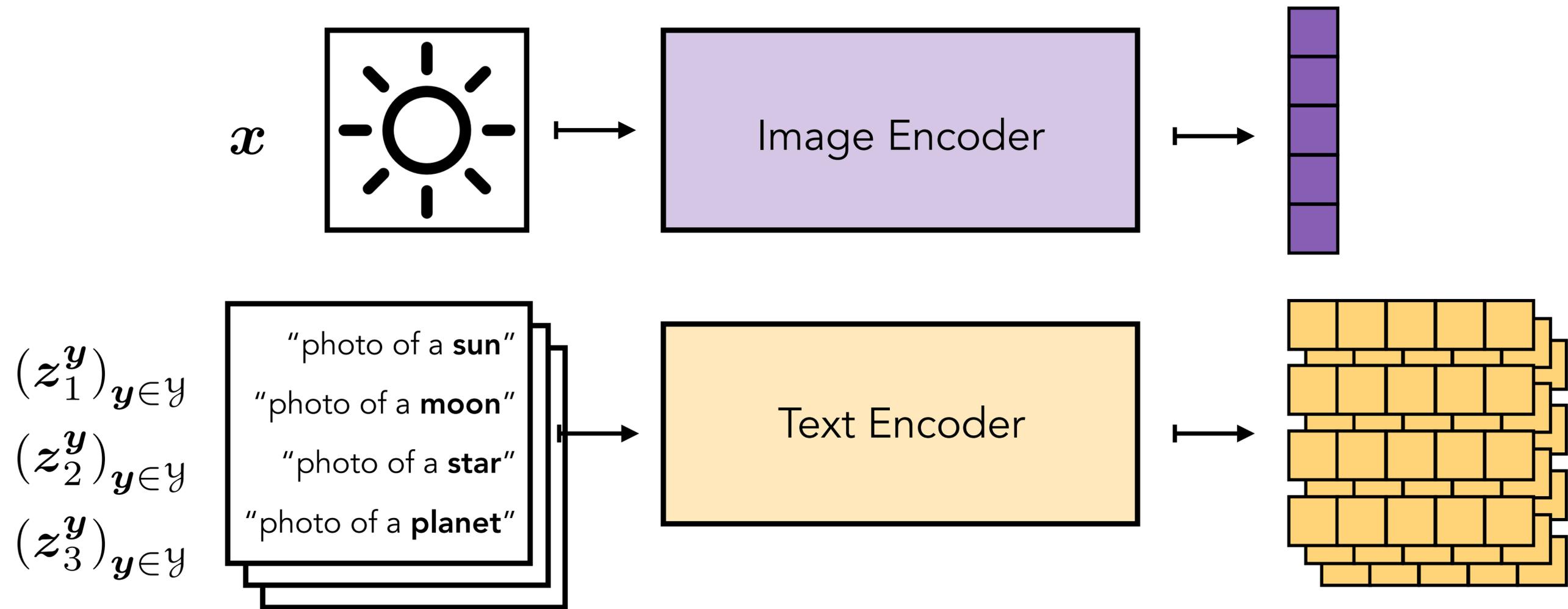
What is the statistical effect of **balancing** methods on the pre-training and the resulting foundation model?

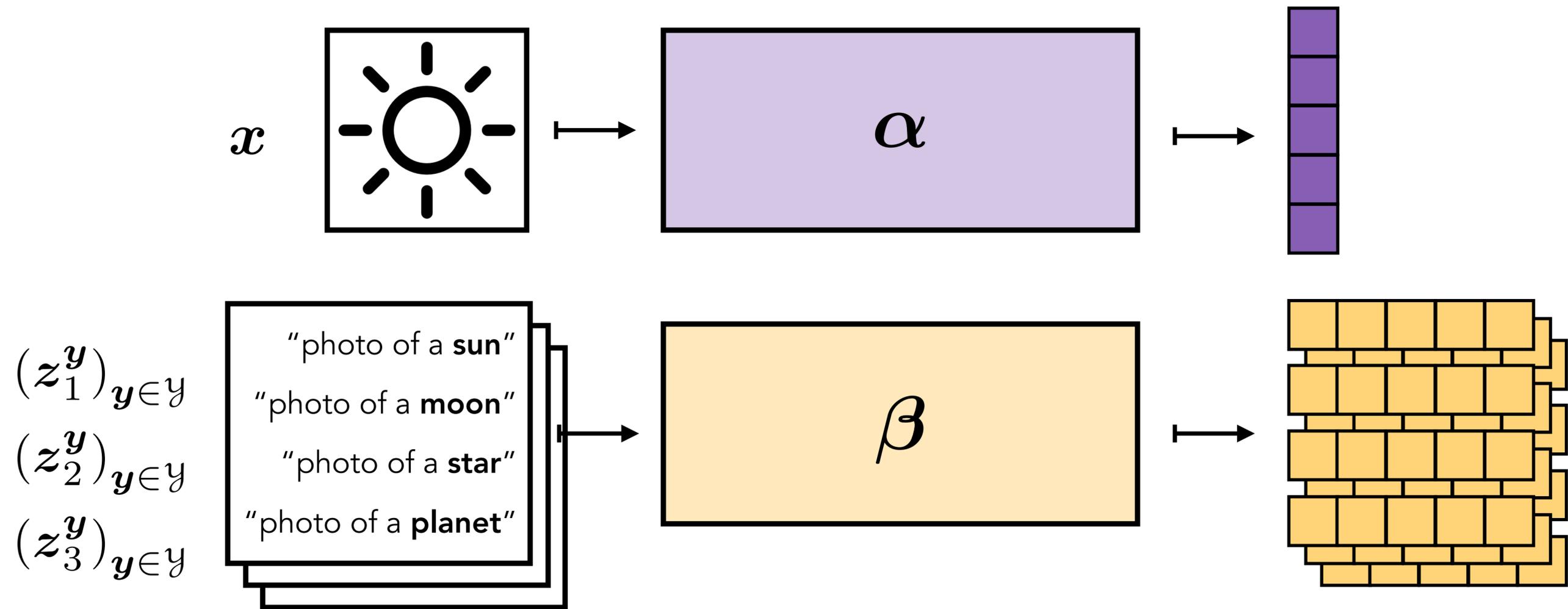


Zero-Shot Prediction

How do we analyze prompt-based **zero-shot prediction** as a statistical estimator, such as its comparison to direct supervision?







$$\mathbf{x} \mapsto \arg \max_{\mathbf{y} \in \mathcal{Y}} \frac{1}{m} \sum_{k=1}^m \langle \alpha(\mathbf{x}), \beta(z_k^{\mathbf{y}}) \rangle$$

Zero-Shot Prediction: Compatibility of Three Distributions

How should we interpret this procedure from the perspective of statistical learning theory?

Zero-Shot Prediction: Compatibility of Three Distributions

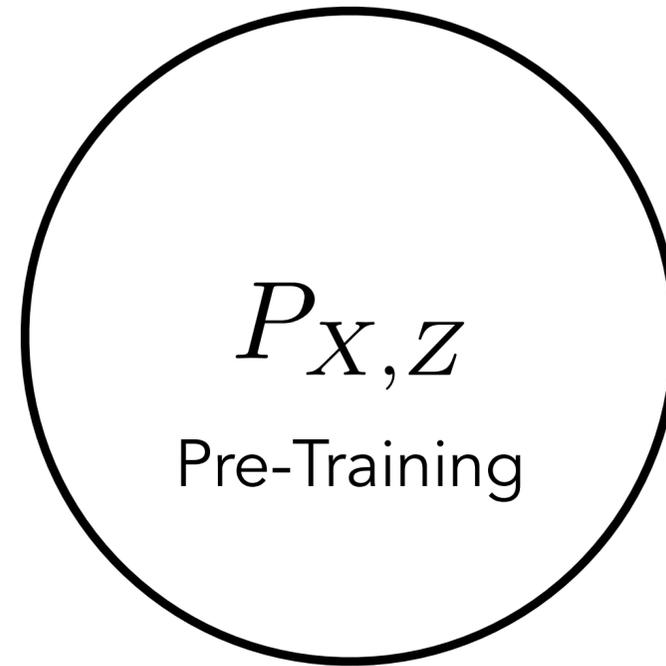
How should we interpret this procedure from the perspective of statistical learning theory?

Step 1: Define the downstream task we wish to solve (no different from supervised learning).

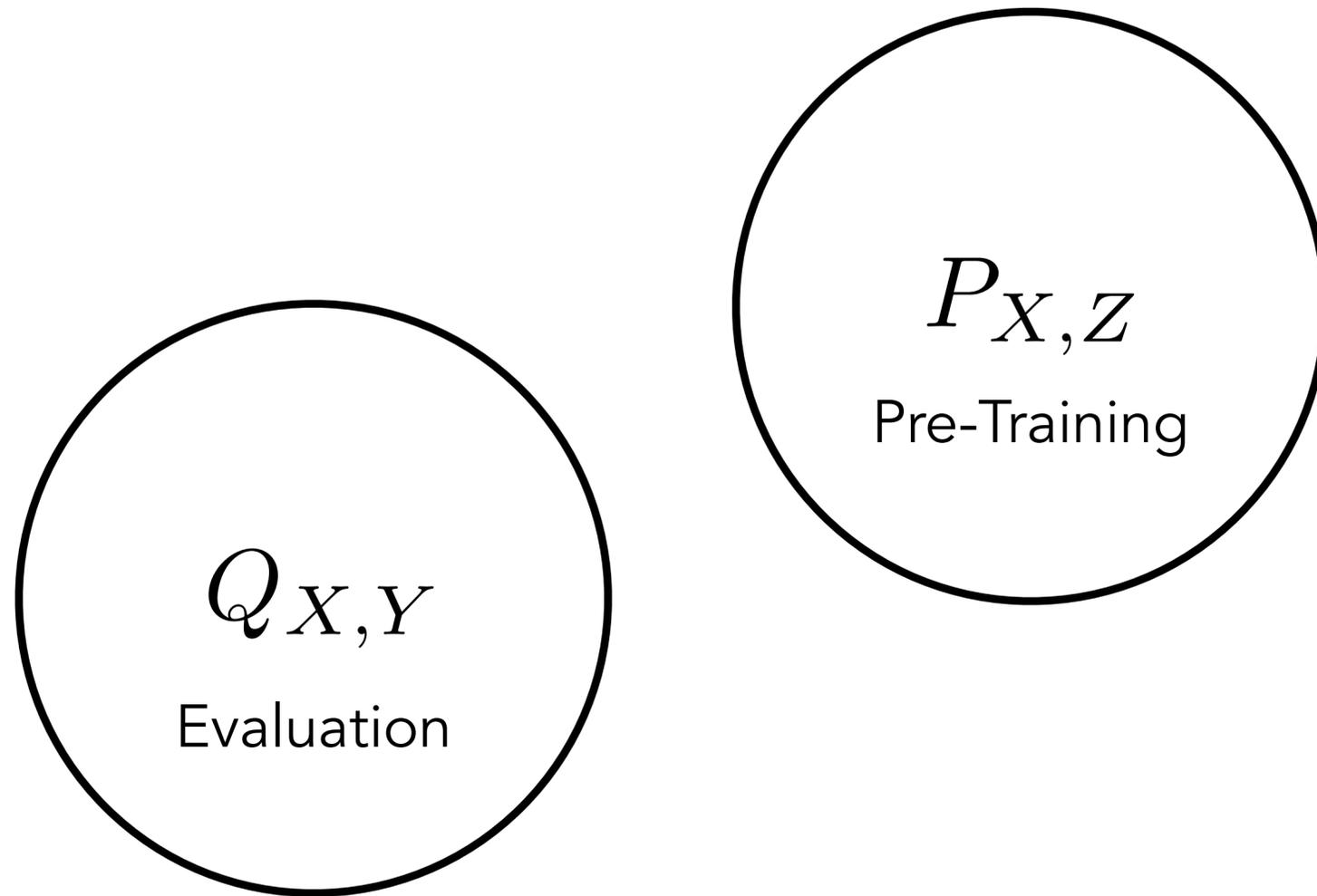
Step 2: Define the population version of the zero-shot prediction procedure and justify it.

Step 3: Provide two estimation frameworks and generic error decompositions that can be used to prove generalization risk bounds.

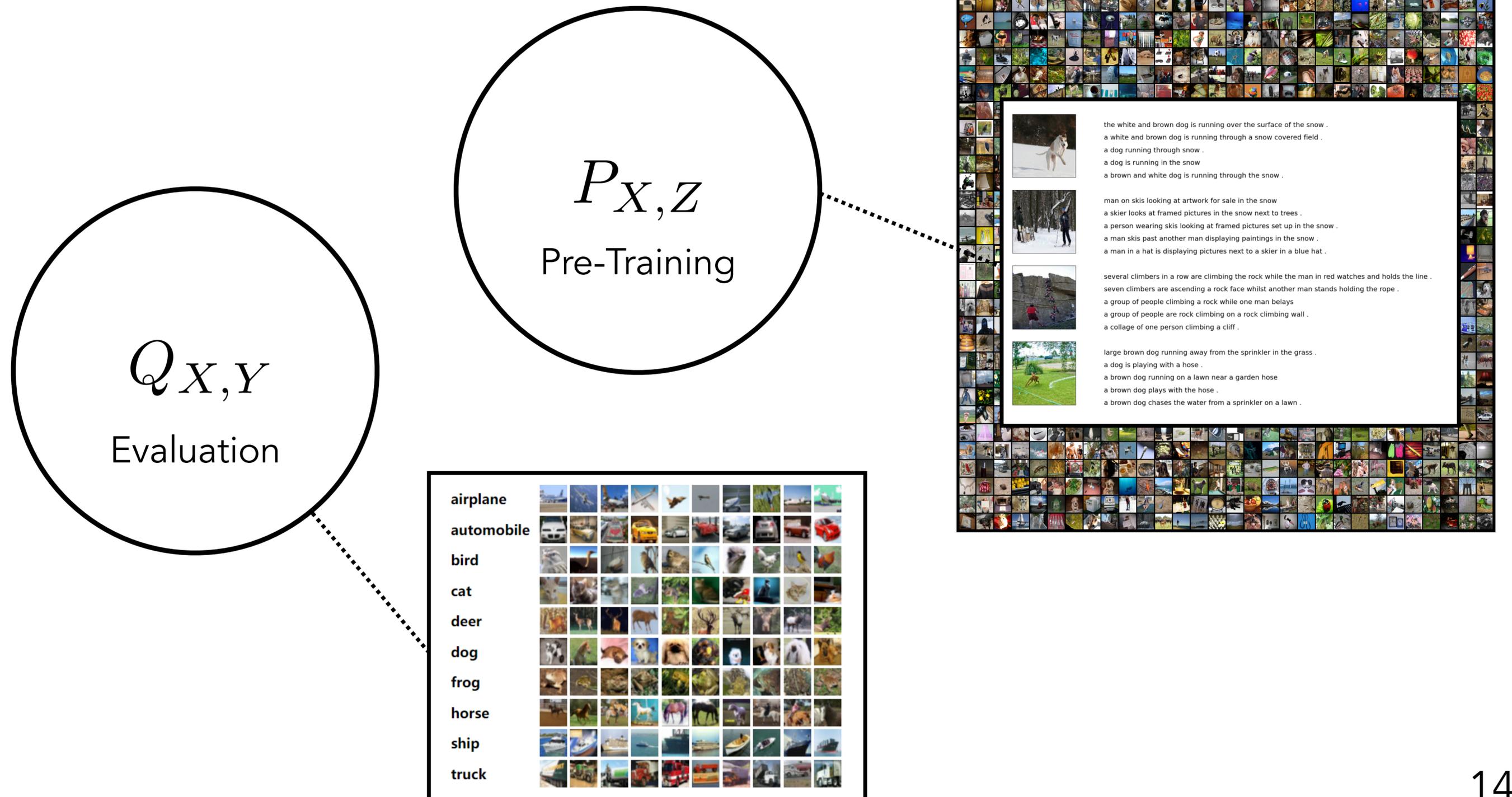
Zero-Shot Prediction: Compatibility of Three Distributions



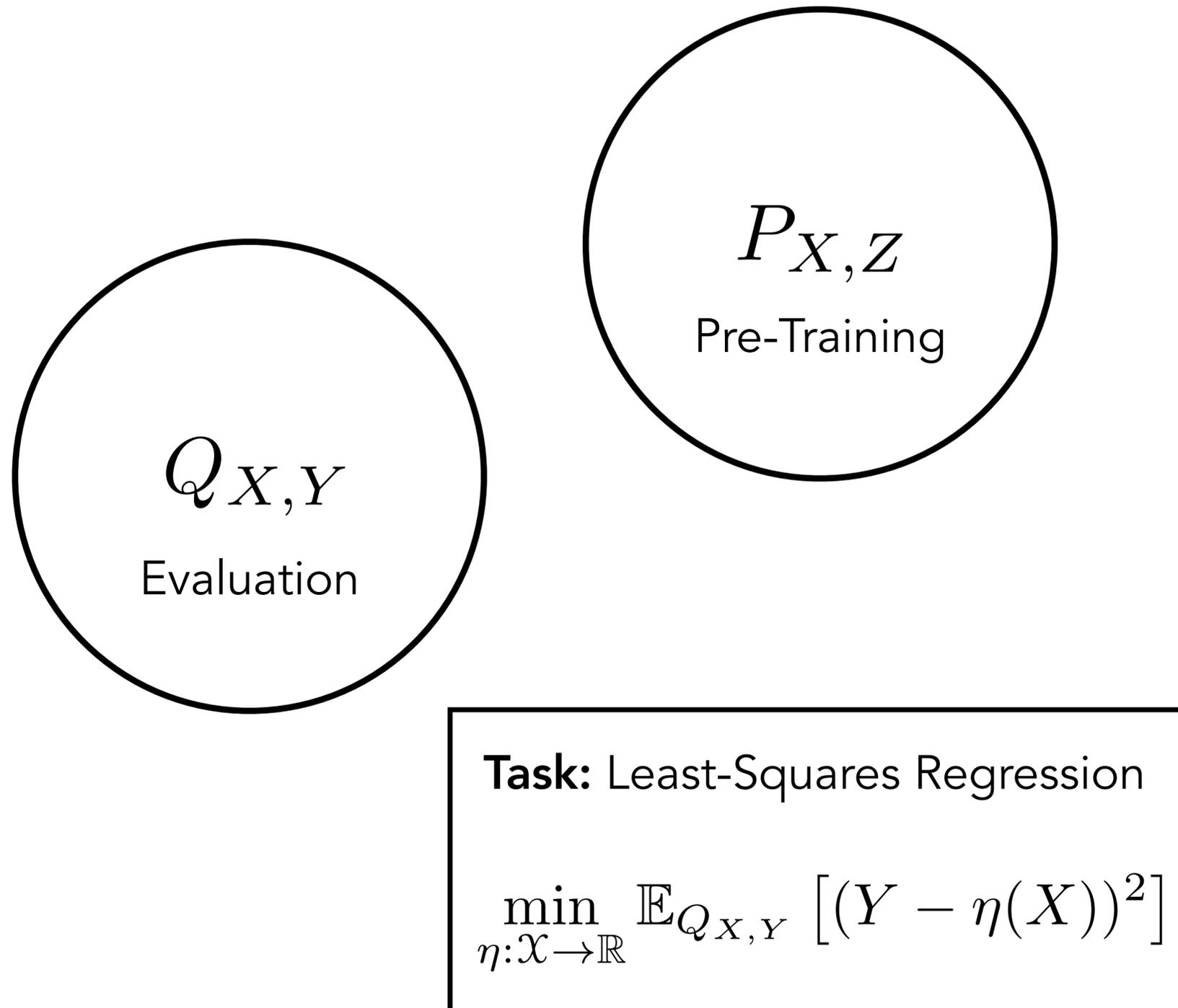
Zero-Shot Prediction: Compatibility of Three Distributions



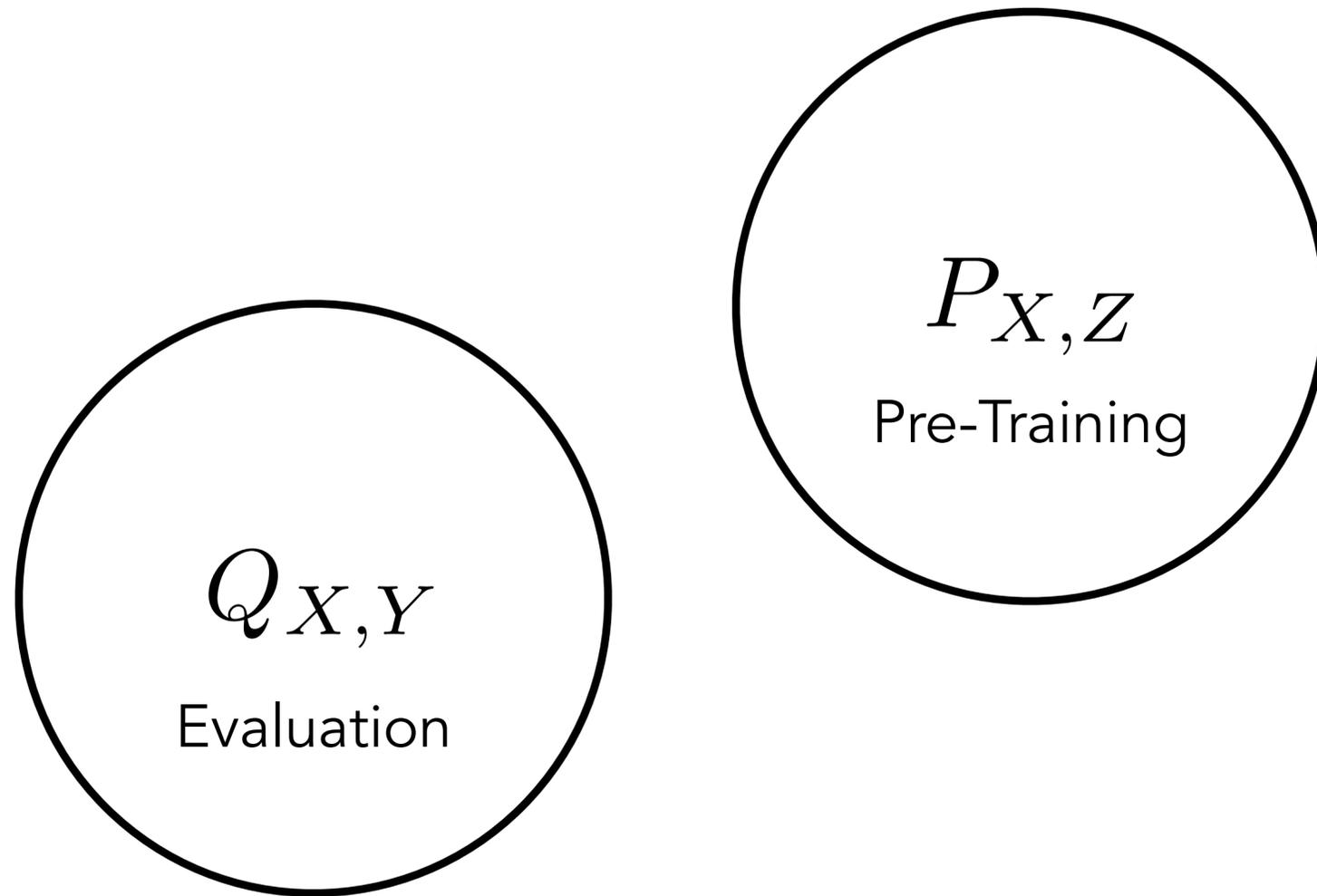
Zero-Shot Prediction: Compatibility of Three Distributions



Zero-Shot Prediction: Compatibility of Three Distributions

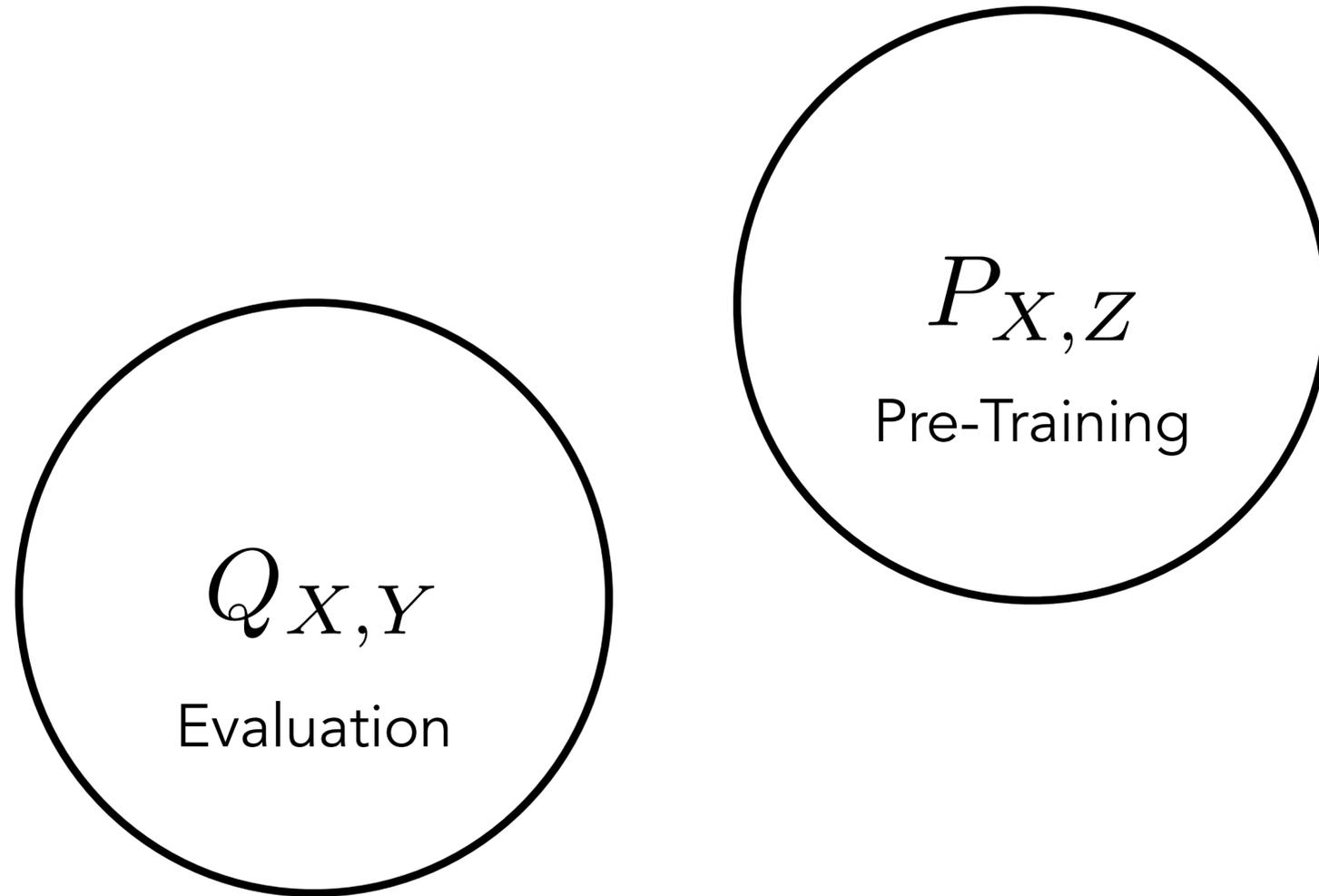


Zero-Shot Prediction: Compatibility of Three Distributions



$$\eta_{\star}(\mathbf{x}) = \mathbb{E}_{Q_{X,Y}} [Y|X](\mathbf{x})$$

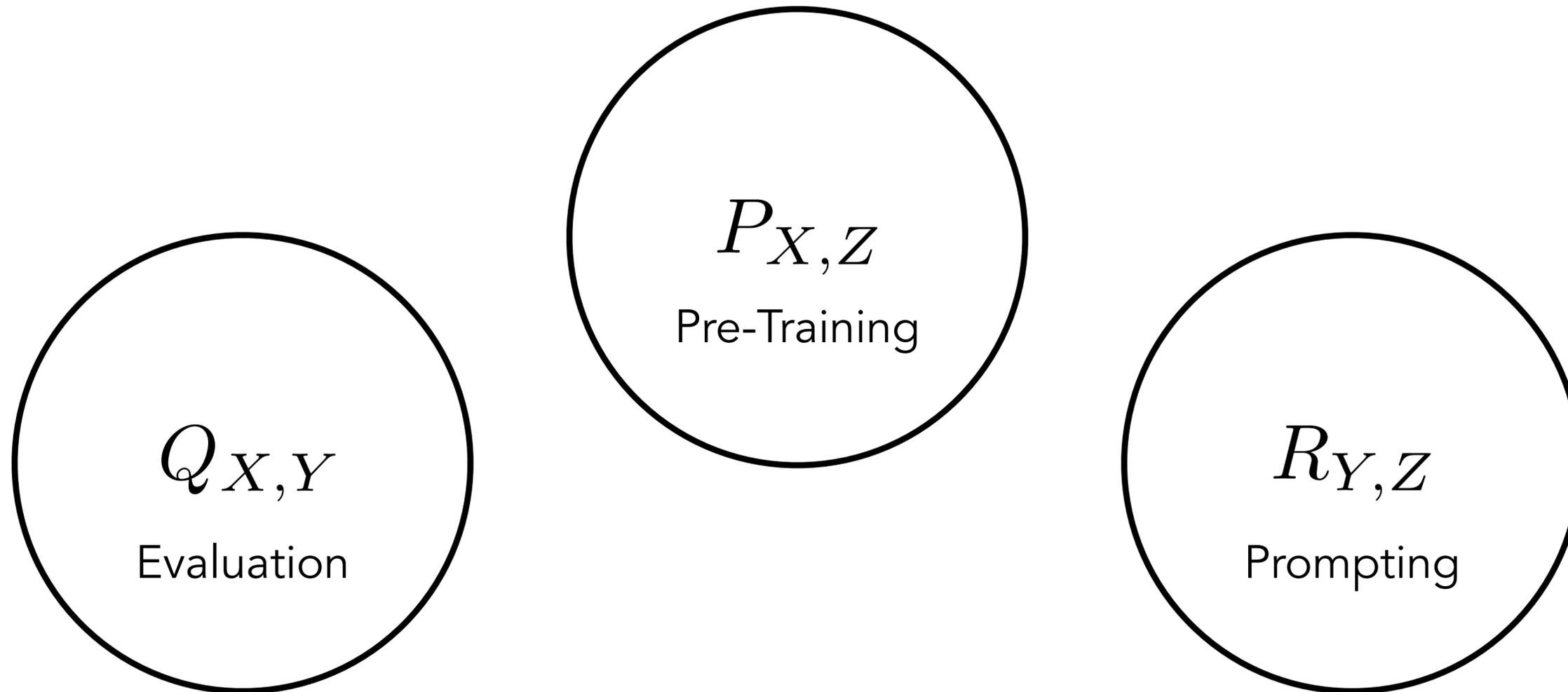
Zero-Shot Prediction: Compatibility of Three Distributions



$$\eta_{\star}(\mathbf{x}) = \mathbb{E}_{Q_{X,Y}} [Y|X] (\mathbf{x})$$

$$\eta_{ZS}(\mathbf{x}) =$$

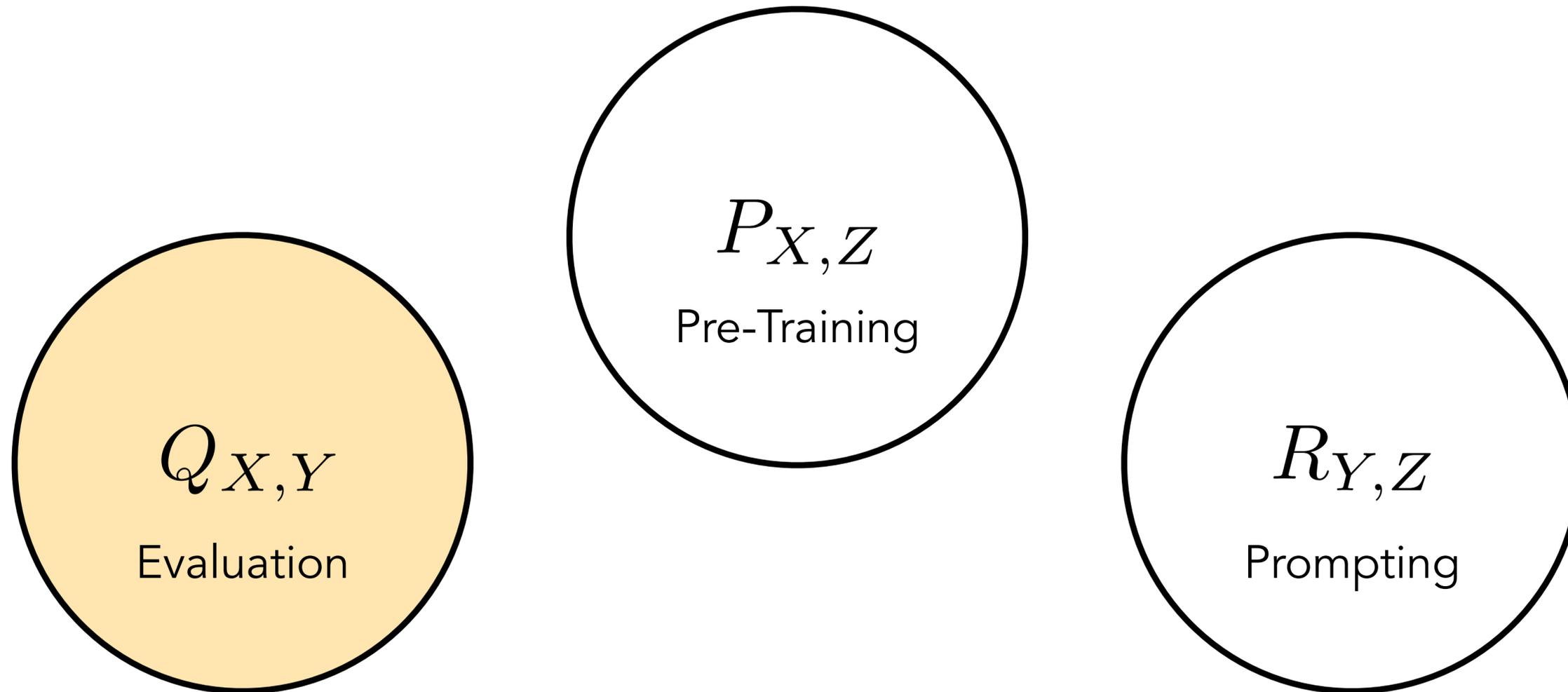
Zero-Shot Prediction: Compatibility of Three Distributions



$$\eta_{\star}(\mathbf{x}) = \mathbb{E}_{Q_{X,Y}} [Y|X] (\mathbf{x})$$

$$\eta_{\text{ZS}}(\mathbf{x}) = \mathbb{E}_{P_{X,Z}} [\mathbb{E}_{R_{Y,Z}} [Y|Z] |X] (\mathbf{x})$$

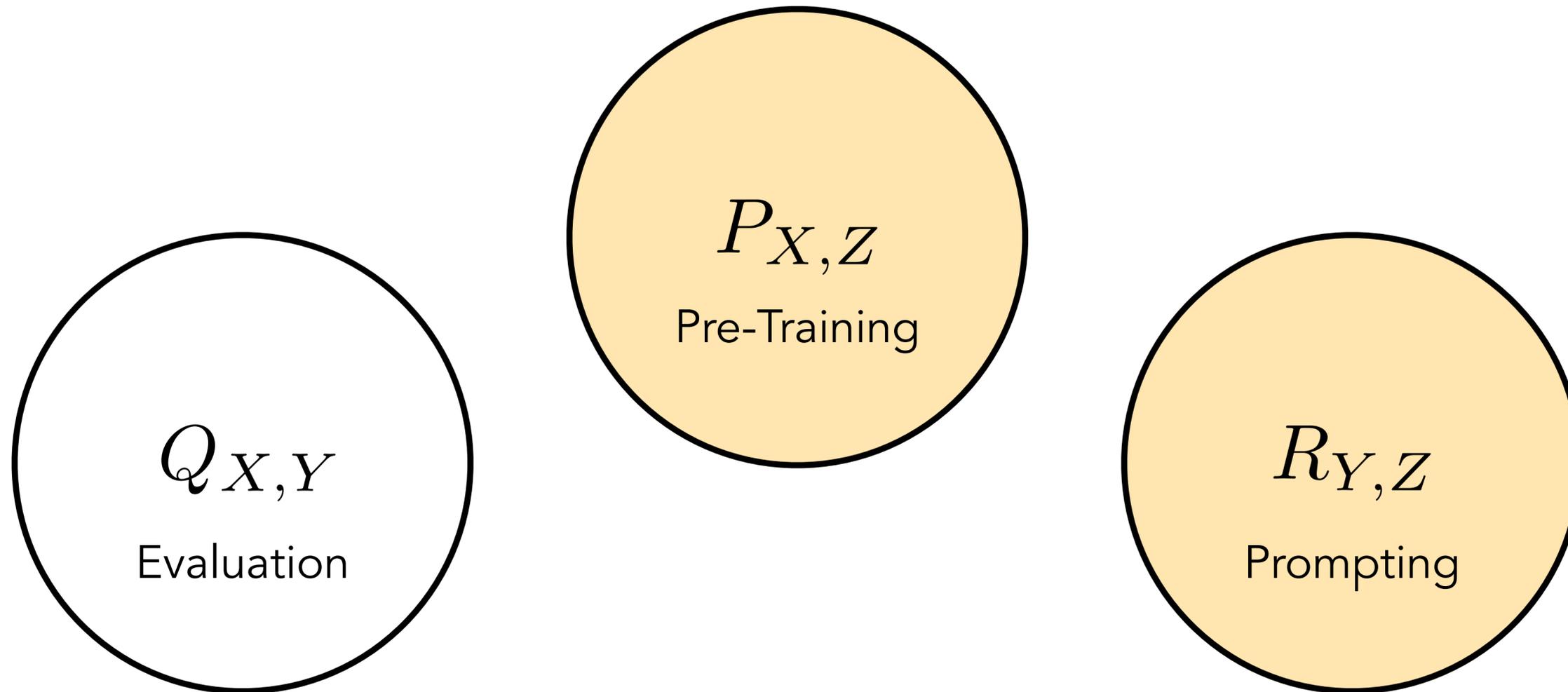
Zero-Shot Prediction: Compatibility of Three Distributions



$$\eta_{\star}(\mathbf{x}) = \mathbb{E}_{Q_{X,Y}} [Y|X] (\mathbf{x})$$

$$\eta_{ZS}(\mathbf{x}) = \mathbb{E}_{P_{X,Z}} [\mathbb{E}_{R_{Y,Z}} [Y|Z] |X] (\mathbf{x})$$

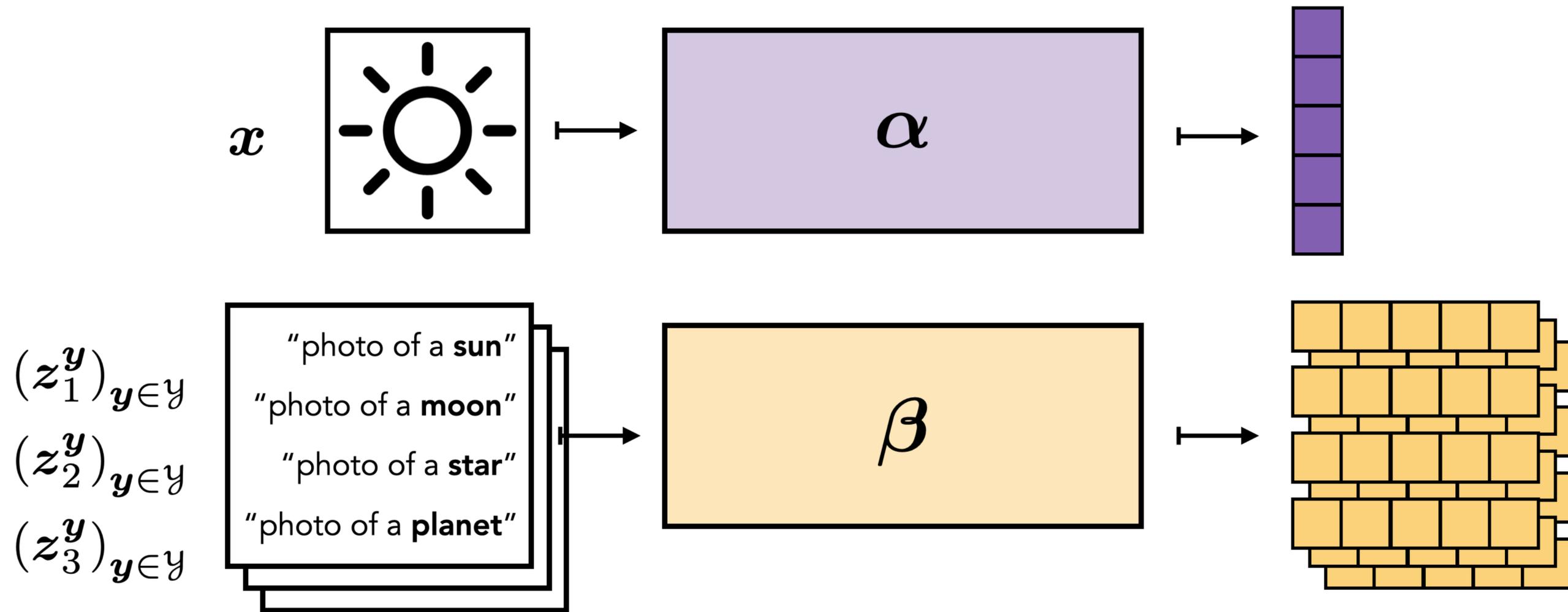
Zero-Shot Prediction: Compatibility of Three Distributions



$$\eta_{\star}(\mathbf{x}) = \mathbb{E}_{Q_{X,Y}} [Y|X] (\mathbf{x})$$

$$\eta_{\text{ZS}}(\mathbf{x}) = \mathbb{E}_{P_{X,Z}} [\mathbb{E}_{R_{Y,Z}} [Y|Z] |X] (\mathbf{x})$$

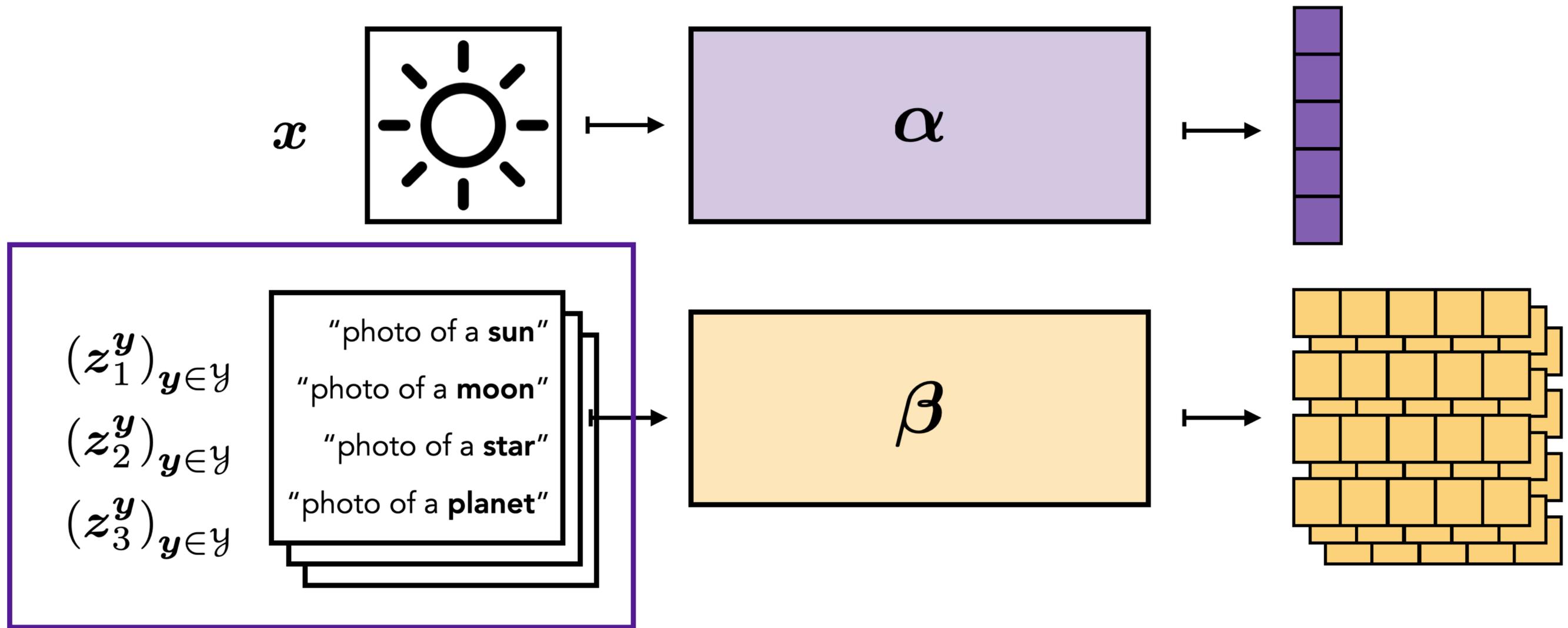
Zero-Shot Prediction: Compatibility of Three Distributions



$$\eta_{\star}(\mathbf{x}) = \mathbb{E}_{Q_{X,Y}} [Y | X] (\mathbf{x})$$

$$\eta_{ZS}(\mathbf{x}) = \mathbb{E}_{P_{X,Z}} [\mathbb{E}_{R_{Y,Z}} [Y | Z] | X] (\mathbf{x})$$

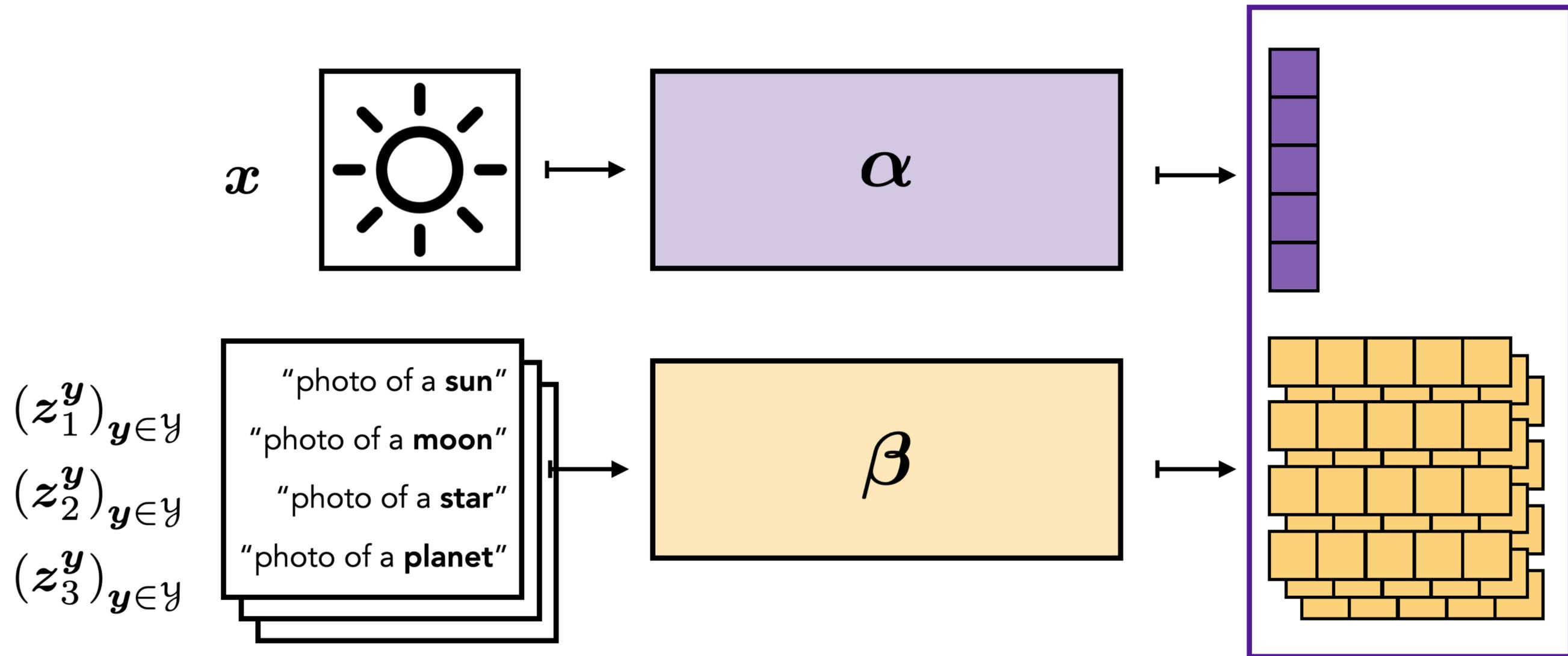
Zero-Shot Prediction: Compatibility of Three Distributions



$$\eta_{\star}(\mathbf{x}) = \mathbb{E}_{Q_{X,Y}} [Y|X](\mathbf{x})$$

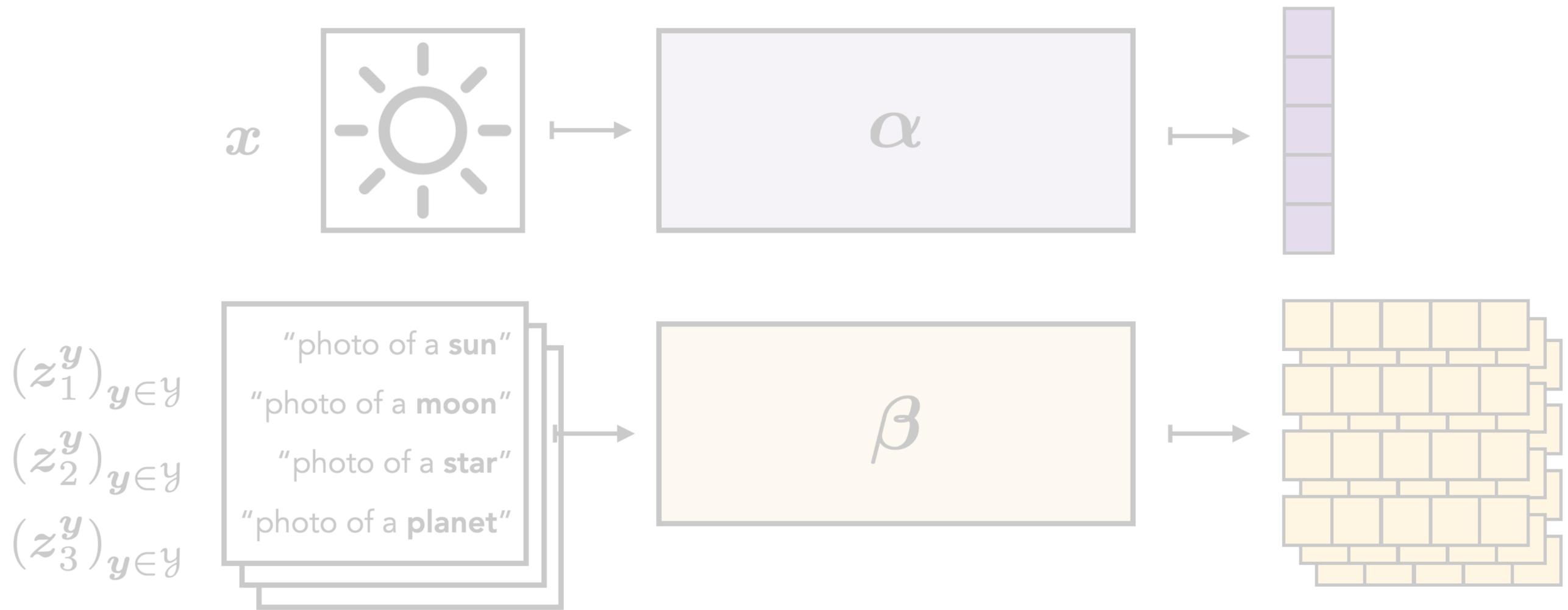
$$\eta_{ZS}(\mathbf{x}) = \mathbb{E}_{P_{X,Z}} [\mathbb{E}_{R_{Y,Z}} [Y|Z]|X](\mathbf{x})$$

Zero-Shot Prediction: Compatibility of Three Distributions



$$\eta_{\star}(\mathbf{x}) = \mathbb{E}_{Q_{X,Y}} [Y | X] (\mathbf{x})$$

$$\eta_{ZS}(\mathbf{x}) = \mathbb{E}_{P_{X,Z}} [\mathbb{E}_{R_{Y,Z}} [Y | Z] | X] (\mathbf{x})$$



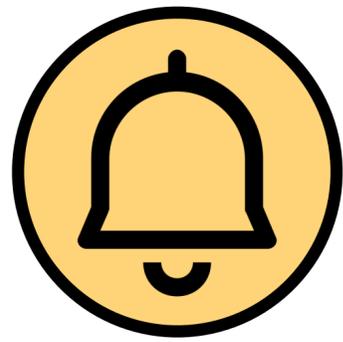
$$\eta_{\star}(\mathbf{x}) = \mathbb{E}_{Q_{X,Y}} [Y|X](\mathbf{x})$$

$$\eta_{ZS}(\mathbf{x}) = \mathbb{E}_{P_{X,Z}} [\mathbb{E}_{R_{Y,Z}} [Y|Z]|X](\mathbf{x})$$

$$\eta_{ZS}(\mathbf{x}) = \mathbb{E}_{P_{X,Z}} [g_R(Z)|X](\mathbf{x})$$

Conditional Mean

Information Density



Conditional Mean

$$\mu_{X \leftarrow Z} : \mathbf{L}^2(P_Z) \rightarrow \mathbf{L}^2(P_X)$$

$$[\mu_{X \leftarrow Z} g](\mathbf{x}) = \mathbb{E}_P [g(Z) | X](\mathbf{x})$$

Information Density

$$R(\mathbf{x}, \mathbf{z}) := \frac{dP_{X,Z}}{dP_X P_Z}(\mathbf{x}, \mathbf{z}) = \frac{dP_{Z|X=\mathbf{x}}}{dP_Z}(\mathbf{z})$$

Conditional Mean

$$\eta_{ZS}(\mathbf{x}) = [\mu_{X \leftarrow Z} g_R](\mathbf{x})$$

Information Density

$$\eta_{ZS}(\mathbf{x}) = \mathbb{E}_{R_{Y,Z}} [Y \cdot R(\mathbf{x}, Z)] + \text{err}(P_Z, R_Z)$$

Conditional Mean

$$\eta_{ZS}(\mathbf{x}) = [\mu_{X \leftarrow Z} g_R](\mathbf{x})$$

Conditional Mean Approach

$$\hat{\eta}_{ZS}(\mathbf{x}) = [\hat{\mu}_{X \leftarrow Z} \hat{g}_R](\mathbf{x})$$

- 1** $(X_1, Z_1), \dots, (X_n, Z_n) \sim P_{X,Z}$
Estimate $\mu_{X \leftarrow Z}$ (operator regression)
- 2** $(Y_1, Z_1), \dots, (Y_M, Z_M) \sim R_{Y,Z}$
Estimate g_R (nonparametric regression)

Information Density

$$\eta_{ZS}(\mathbf{x}) = \mathbb{E}_{R_{Y,Z}} [Y \cdot R(\mathbf{x}, Z)] + \text{err}(P_Z, R_Z)$$

Conditional Mean

$$\eta_{ZS}(\mathbf{x}) = [\mu_{X \leftarrow Z} g_R](\mathbf{x})$$

Conditional Mean Approach

$$\hat{\eta}_{ZS}(\mathbf{x}) = [\hat{\mu}_{X \leftarrow Z} \hat{g}_R](\mathbf{x})$$

- 1 $(X_1, Z_1), \dots, (X_n, Z_n) \sim P_{X,Z}$
Estimate $\mu_{X \leftarrow Z}$ (operator regression)
- 2 $(Y_1, Z_1), \dots, (Y_M, Z_M) \sim R_{Y,Z}$
Estimate g_R (nonparametric regression)

Information Density

$$\eta_{ZS}(\mathbf{x}) = \mathbb{E}_{R_{Y,Z}} [Y \cdot R(\mathbf{x}, Z)] + \text{err}(P_Z, R_Z)$$

Information Density Approach

$$\hat{\eta}_{ZS}(\mathbf{x}) = \mathbb{E}_{\hat{R}_{Y,Z}} [Y \cdot \hat{R}(\mathbf{x}, Z)]$$

- 1 $(X_1, Z_1), \dots, (X_n, Z_n) \sim P_{X,Z}$
 $(X'_1, Z'_1), \dots, (X'_n, Z'_n) \sim P_X P_Z$
Estimate R (RND estimation)
- 2 $(Y_1, Z_1), \dots, (Y_M, Z_M) \sim R_{Y,Z}$
Estimate $\mathbb{E}_{R_{Y,Z}} [\cdot]$ (infinite-dimensional mean estimation)

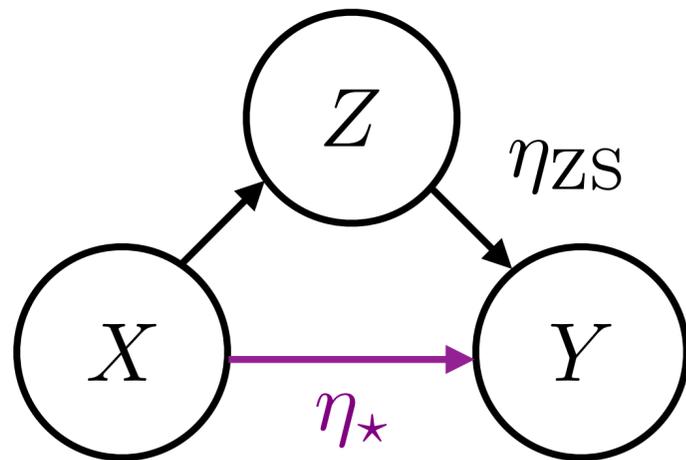
$$\|\eta_\star - \hat{\eta}_{\text{ZS}}\|_{\mathbf{L}^2(Q_X)}^2 \leq 2\|\eta_\star - \eta_{\text{ZS}}\|_{\mathbf{L}^2(Q_X)}^2 + 2\|\eta_{\text{ZS}} - \hat{\eta}_{\text{ZS}}\|_{\mathbf{L}^2(Q_X)}^2$$

$$\|\eta_{\star} - \hat{\eta}_{ZS}\|_{\mathbf{L}^2(Q_X)}^2 \leq 2\|\eta_{\star} - \eta_{ZS}\|_{\mathbf{L}^2(Q_X)}^2 + 2\|\eta_{ZS} - \hat{\eta}_{ZS}\|_{\mathbf{L}^2(Q_X)}^2$$

Thm. 1 (M., Harchaoui)

$$\|\eta_{\star} - \eta_{ZS}\|_{\mathbf{L}^2(Q_X)}^2 \lesssim$$

$Q_{X,Y,Z}$ joint evaluation distribution
with **latent caption**

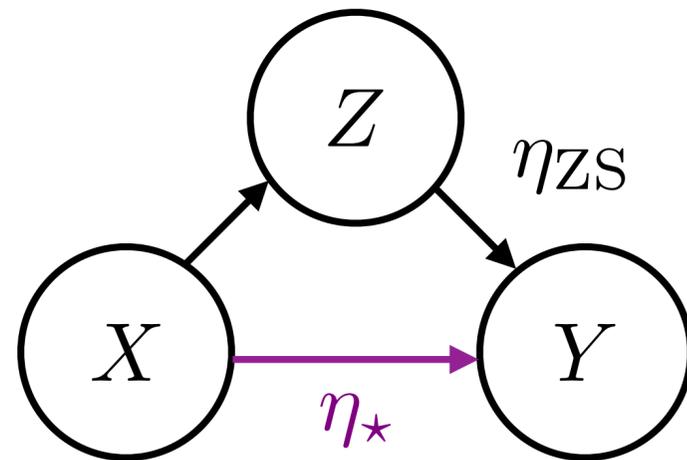


$$\|\eta_{\star} - \hat{\eta}_{ZS}\|_{\mathbf{L}^2(Q_X)}^2 \leq 2\|\eta_{\star} - \eta_{ZS}\|_{\mathbf{L}^2(Q_X)}^2 + 2\|\eta_{ZS} - \hat{\eta}_{ZS}\|_{\mathbf{L}^2(Q_X)}^2$$

Thm. 1 (M., Harchaoui)

$$\|\eta_{\star} - \eta_{ZS}\|_{\mathbf{L}^2(Q_X)}^2 \lesssim \mathbb{E}_{Q_Z} [I(X, Y | Z)]$$

$Q_{X,Y,Z}$ joint evaluation distribution with **latent caption**



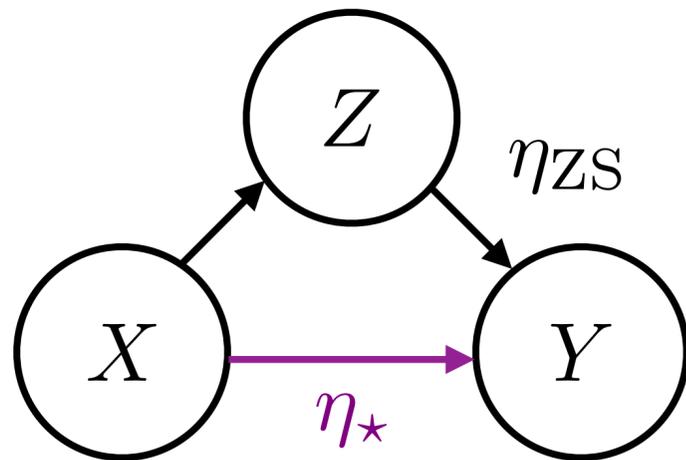
Direct and indirect paths differ when X and Y are **conditionally dependent** given Z .

$$\|\eta_{\star} - \hat{\eta}_{ZS}\|_{\mathbf{L}^2(Q_X)}^2 \leq 2\|\eta_{\star} - \eta_{ZS}\|_{\mathbf{L}^2(Q_X)}^2 + 2\|\eta_{ZS} - \hat{\eta}_{ZS}\|_{\mathbf{L}^2(Q_X)}^2$$

Thm. 1 (M., Harchaoui)

$$\|\eta_{\star} - \eta_{ZS}\|_{\mathbf{L}^2(Q_X)}^2 \lesssim \mathbb{E}_{Q_Z} [I(X, Y | Z)] + \|g_Q - g_R\|_{\mathbf{L}^2(Q_Z)}^2$$

$Q_{X,Y,Z}$ joint evaluation distribution with **latent caption**



Direct and indirect paths differ when X and Y are **conditionally dependent** given Z .

$g_R(\mathbf{z}) := \mathbb{E}_{R_{Y,Z}} [Y | Z](\mathbf{z})$ prompt distribution

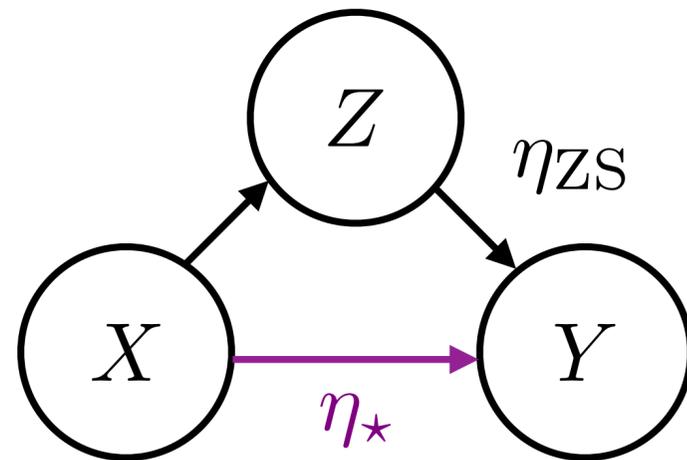
$g_Q(\mathbf{z}) := \mathbb{E}_{Q_{Y,Z}} [Y | Z](\mathbf{z})$ true relationship

$$\|\eta_{\star} - \hat{\eta}_{ZS}\|_{\mathbf{L}^2(Q_X)}^2 \leq 2\|\eta_{\star} - \eta_{ZS}\|_{\mathbf{L}^2(Q_X)}^2 + 2\|\eta_{ZS} - \hat{\eta}_{ZS}\|_{\mathbf{L}^2(Q_X)}^2$$

Thm. 1 (M., Harchaoui)

$$\|\eta_{\star} - \eta_{ZS}\|_{\mathbf{L}^2(Q_X)}^2 \lesssim \mathbb{E}_{Q_Z} [I(X, Y | Z)] + \|g_Q - g_R\|_{\mathbf{L}^2(Q_Z)}^2$$

$Q_{X,Y,Z}$ joint evaluation distribution with **latent caption**



Direct and indirect paths differ when X and Y are **conditionally dependent** given Z .

$g_R(\mathbf{z}) := \mathbb{E}_{R_{Y,Z}} [Y | Z](\mathbf{z})$ prompt distribution

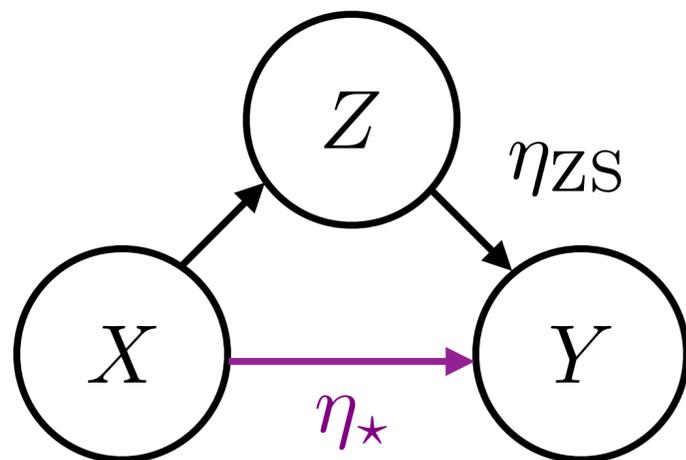
$g_Q(\mathbf{z}) := \mathbb{E}_{Q_{Y,Z}} [Y | Z](\mathbf{z})$ true relationship

$$\|\eta_{\star} - \hat{\eta}_{ZS}\|_{\mathbf{L}^2(Q_X)}^2 \leq 2\|\eta_{\star} - \eta_{ZS}\|_{\mathbf{L}^2(Q_X)}^2 + 2\|\eta_{ZS} - \hat{\eta}_{ZS}\|_{\mathbf{L}^2(Q_X)}^2$$

Thm. 1 (M., Harchaoui)

$$\|\eta_{\star} - \eta_{ZS}\|_{\mathbf{L}^2(Q_X)}^2 \lesssim \mathbb{E}_{Q_Z} [I(X, Y|Z)] + \|g_Q - g_R\|_{\mathbf{L}^2(Q_Z)}^2$$

$Q_{X,Y,Z}$ joint evaluation distribution with **latent caption**



Direct and indirect paths differ when X and Y are **conditionally dependent** given Z .

Thm. 2 (M., Harchaoui)

Conditional Mean

$$\|\eta_{ZS} - \hat{\eta}_{ZS}\|_{\mathbf{L}^2(P_X)}^2 \lesssim \text{plog}(1/\delta) \left[n^{-\frac{q}{q+1}} + M^{-\frac{2p}{2p-1}} \right] \text{ w.p. } 1 - \delta$$

$$\hat{\eta}_{ZS}(\mathbf{x}) = [\hat{\mu}_{X \leftarrow Z} \hat{g}_R](\mathbf{x})$$

$g_R(\mathbf{z}) := \mathbb{E}_{R_{Y,Z}} [Y|Z](\mathbf{z})$ prompt distribution

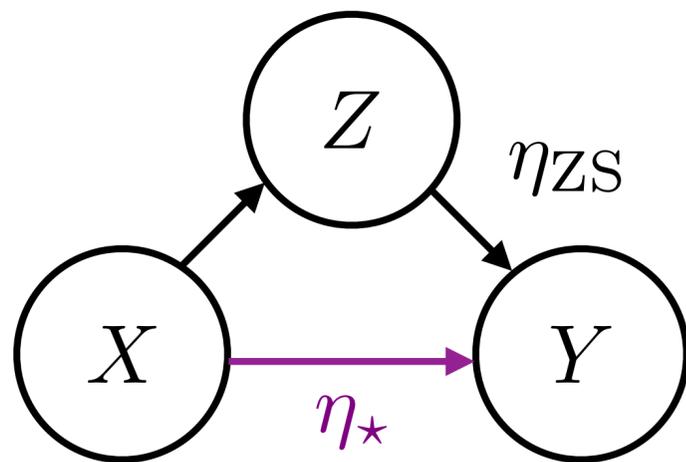
$g_Q(\mathbf{z}) := \mathbb{E}_{Q_{Y,Z}} [Y|Z](\mathbf{z})$ true relationship

$$\|\eta_{\star} - \hat{\eta}_{ZS}\|_{\mathbf{L}^2(Q_X)}^2 \leq 2\|\eta_{\star} - \eta_{ZS}\|_{\mathbf{L}^2(Q_X)}^2 + 2\|\eta_{ZS} - \hat{\eta}_{ZS}\|_{\mathbf{L}^2(Q_X)}^2$$

Thm. 1 (M., Harchaoui)

$$\|\eta_{\star} - \eta_{ZS}\|_{\mathbf{L}^2(Q_X)}^2 \lesssim \mathbb{E}_{Q_Z} [I(X, Y|Z)] + \|g_Q - g_R\|_{\mathbf{L}^2(Q_Z)}^2$$

$Q_{X,Y,Z}$ joint evaluation distribution with **latent caption**



Direct and indirect paths differ when X and Y are **conditionally dependent** given Z .

$g_R(\mathbf{z}) := \mathbb{E}_{R_{Y,Z}} [Y|Z](\mathbf{z})$ prompt distribution

$g_Q(\mathbf{z}) := \mathbb{E}_{Q_{Y,Z}} [Y|Z](\mathbf{z})$ true relationship

Thm. 2 (M., Harchaoui)

Conditional Mean

$$\|\eta_{ZS} - \hat{\eta}_{ZS}\|_{\mathbf{L}^2(P_X)}^2 \lesssim \text{plog}(1/\delta) \left[n^{-\frac{q}{q+1}} + M^{-\frac{2p}{2p-1}} \right] \text{ w.p. } 1 - \delta$$

$$\hat{\eta}_{ZS}(\mathbf{x}) = [\hat{\mu}_{X \leftarrow Z} \hat{g}_R](\mathbf{x})$$

$n^{-\frac{q}{q+1}}$ (estimating $\mu_{X \leftarrow Z}$ using $\hat{\mu}_{X \leftarrow Z}$)

$1 \leftarrow q \rightarrow \infty$ effective dimension of Z is small relative to effective dimension of X
 X and Z are (approximately) independent

$M^{-\frac{2p}{2p-1}}$ (estimating g_R using \hat{g}_R)

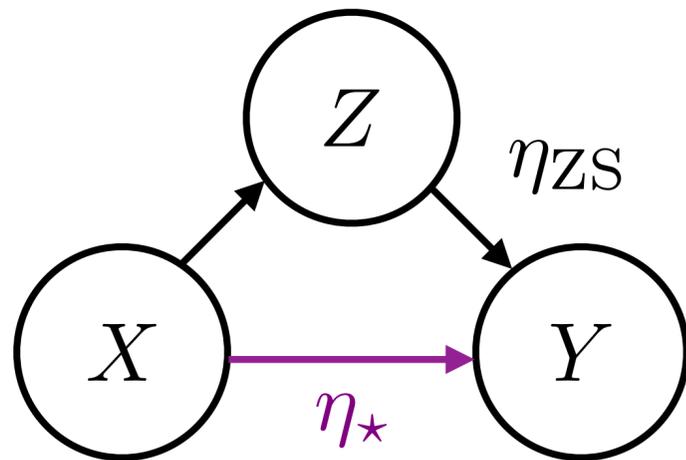
$1 \leftarrow p \rightarrow \infty$ effective dimension of Z is small

$$\|\eta_{\star} - \hat{\eta}_{ZS}\|_{\mathbf{L}^2(Q_X)}^2 \leq 2\|\eta_{\star} - \eta_{ZS}\|_{\mathbf{L}^2(Q_X)}^2 + 2\|\eta_{ZS} - \hat{\eta}_{ZS}\|_{\mathbf{L}^2(Q_X)}^2$$

Thm. 1 (M., Harchaoui)

$$\|\eta_{\star} - \eta_{ZS}\|_{\mathbf{L}^2(Q_X)}^2 \lesssim \mathbb{E}_{Q_Z} [I(X, Y|Z)] + \|g_Q - g_R\|_{\mathbf{L}^2(Q_Z)}^2$$

$Q_{X,Y,Z}$ joint evaluation distribution with **latent caption**



Direct and indirect paths differ when X and Y are **conditionally dependent** given Z .

Thm. 3 (M., Harchaoui)

Information Density

$$\|\eta_{ZS} - \hat{\eta}_{ZS}\|_{\mathbf{L}^2(P_X)}^2 \lesssim \text{err}(P_Z, R_Z) \text{plog}(1/\delta) \left[n^{-\frac{q}{q+1}} + M^{-1} \right] \text{ w.p. } 1 - \delta$$

$$\hat{\eta}_{ZS}(\mathbf{x}) = \mathbb{E}_{\hat{R}_{Y,Z}} [Y \cdot \hat{R}(\mathbf{x}, Z)]$$

$g_R(\mathbf{z}) := \mathbb{E}_{R_{Y,Z}} [Y|Z](\mathbf{z})$ prompt distribution

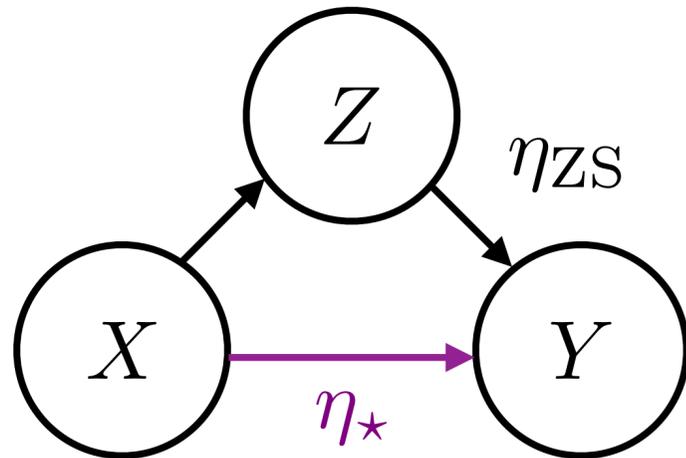
$g_Q(\mathbf{z}) := \mathbb{E}_{Q_{Y,Z}} [Y|Z](\mathbf{z})$ true relationship

$$\|\eta_{\star} - \hat{\eta}_{ZS}\|_{\mathbf{L}^2(Q_X)}^2 \leq 2\|\eta_{\star} - \eta_{ZS}\|_{\mathbf{L}^2(Q_X)}^2 + 2\|\eta_{ZS} - \hat{\eta}_{ZS}\|_{\mathbf{L}^2(Q_X)}^2$$

Thm. 1 (M., Harchaoui)

$$\|\eta_{\star} - \eta_{ZS}\|_{\mathbf{L}^2(Q_X)}^2 \lesssim \mathbb{E}_{Q_Z} [I(X, Y|Z)] + \|g_Q - g_R\|_{\mathbf{L}^2(Q_Z)}^2$$

$Q_{X,Y,Z}$ joint evaluation distribution with **latent caption**



Direct and indirect paths differ when X and Y are **conditionally dependent** given Z .

$g_R(\mathbf{z}) := \mathbb{E}_{R_{Y,Z}} [Y|Z](\mathbf{z})$ prompt distribution

$g_Q(\mathbf{z}) := \mathbb{E}_{Q_{Y,Z}} [Y|Z](\mathbf{z})$ true relationship

Thm. 3 (M., Harchaoui)

Information Density

$$\|\eta_{ZS} - \hat{\eta}_{ZS}\|_{\mathbf{L}^2(P_X)}^2 \lesssim \text{err}(P_Z, R_Z) \text{plog}(1/\delta) \left[n^{-\frac{q}{q+1}} + M^{-1} \right] \text{ w.p. } 1 - \delta$$

$$\hat{\eta}_{ZS}(\mathbf{x}) = \mathbb{E}_{\hat{R}_{Y,Z}} [Y \cdot \hat{R}(\mathbf{x}, Z)]$$

$n^{-\frac{q}{q+1}}$ (estimating R using \hat{R})

$1 \leftarrow q \rightarrow \infty$ X and Z are (approximately) independent

M^{-1} (estimating $R_{Y,Z}$ using $\hat{R}_{Y,Z}$)

$\text{err}(P_Z, R_Z)$ (error term from $\eta_{ZS}(\mathbf{x}) = \mathbb{E}_{R_{Y,Z}} [Y \cdot R(\mathbf{x}, Z)] + \text{err}(P_Z, R_Z)$)

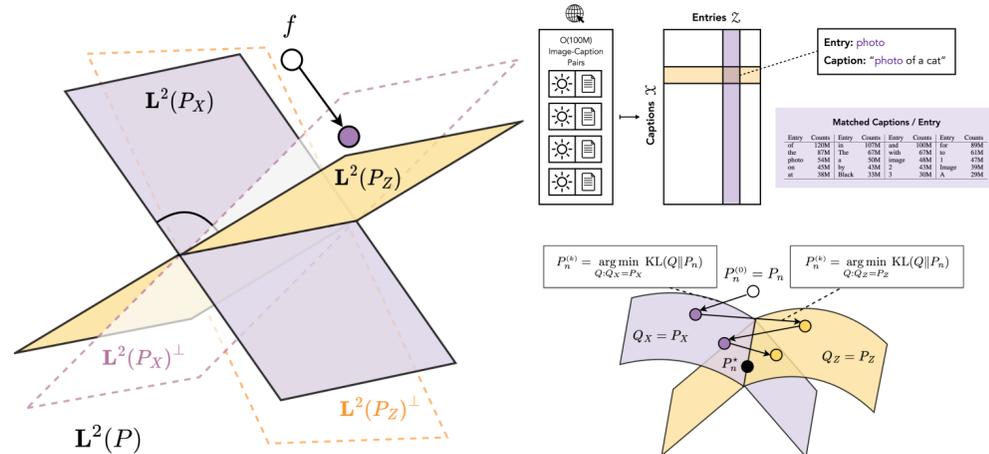
Research Contributions



Exact variance reduction analysis of balancing-based mean estimation!

Theorem (Liu, M., Pal, Harchaoui)

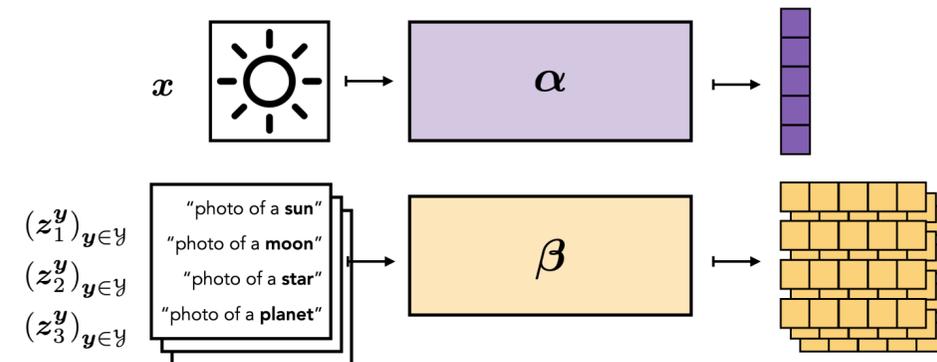
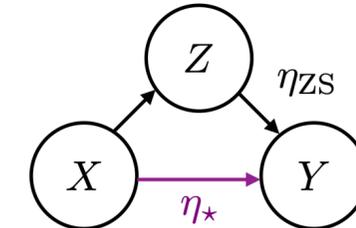
$$\mathbb{E}_P [(P_n^{(k)}(h) - P(h))^2] = \frac{\text{Var}(\dots \overset{k \text{ times}}{\mu_{Z \leftarrow X}^\perp \mu_{X \leftarrow Z}^\perp} h)}{n} + \tilde{O}\left(\frac{k^6}{n^{3/2}}\right)$$



Theoretical framework for obtaining generalization guarantees for two broad classes of zero-shot prediction models!

Thm. 1 (M., Harchaoui)

$$\|\eta_\star - \eta_{ZS}\|_{\mathbf{L}^2(Q_X)}^2 \lesssim \mathbb{E}_{Q_Z} [I(X, Y | Z)] + \|g_Q - g_R\|_{\mathbf{L}^2(Q_Z)}^2$$



Preliminaries

Balanced Pre-Training

Zero-Shot Prediction

Conclusion

Perspectives: Balanced Pre-Training

- Non-asymptotic mean squared error bounds for **nonlinear** functionals.

$$\Psi(P_n^{(k)}) + P_n(\phi_n^{(k)}) - \Psi(P) = [P_n - P](\phi_0) + \mathcal{R}_n^{(k)} + \mathcal{D}_n^{(k)}$$



Perspectives: Balanced Pre-Training

- Non-asymptotic mean squared error bounds for **nonlinear** functionals.

$$\Psi(P_n^{(k)}) + P_n(\phi_n^{(k)}) - \Psi(P) = [P_n - P](\phi_0) + \mathcal{R}_n^{(k)} + \mathcal{D}_n^{(k)}$$

• tangent space for “known marginal” model known (Bickel et al, '90)

$$\underbrace{\Psi(P_n^{(k)}) + P(\phi_n^{(k)}) - \Psi(P)}_{o_p(n^{-1/2}) \text{ remainder?}}$$

$$\underbrace{[P_n - P](\phi_n^{(k)} - \phi_0)}_{o_p(n^{-1/2}) \text{ drift?}}$$

Perspectives: Balanced Pre-Training

- Non-asymptotic mean squared error bounds for **nonlinear** functionals.

$$\Psi(P_n^{(k)}) + P_n(\phi_n^{(k)}) - \Psi(P) = [P_n - P](\phi_0) + \mathcal{R}_n^{(k)} + \mathcal{D}_n^{(k)}$$

$\underbrace{\Psi(P_n^{(k)}) + P(\phi_n^{(k)}) - \Psi(P)}_{o_p(n^{-1/2}) \text{ remainder?}} \quad \underbrace{[P_n - P](\phi_n^{(k)} - \phi_0)}_{o_p(n^{-1/2}) \text{ drift?}}$

tangent space for "known marginal" model known
 (Bickel et al, '90)

- Simultaneously optimizing (P_X, P_Z) for downstream performance (bilevel).

$$\hat{\theta}_n = \arg \min_{\theta \in \mathbb{R}^d} \mathbb{E}_{P_n^*} [h_\theta(X, Z)] \longrightarrow \mathbb{E}_P [h_{\hat{\theta}_n}(X, Z)] - \inf_{\theta \in \mathbb{R}^d} \mathbb{E}_P [h_\theta(X, Z)]$$

Perspectives: Balanced Pre-Training

- Non-asymptotic mean squared error bounds for **nonlinear** functionals.

$$\Psi(P_n^{(k)}) + P_n(\phi_n^{(k)}) - \Psi(P) = [P_n - P](\phi_0) + \mathcal{R}_n^{(k)} + \mathcal{D}_n^{(k)}$$

$\underbrace{\Psi(P_n^{(k)}) + P(\phi_n^{(k)}) - \Psi(P)}_{o(n^{-1/2}) \text{ remainder}} \quad \underbrace{[P_n - P](\phi_n^{(k)} - \phi_0)}_{o(n^{-1/2}) \text{ drift}}$

tangent space for "known marginal" model known
 (Bickel et al, '90)

- Simultaneously optimizing (P_X, P_Z) for downstream performance (bilevel).

$$\hat{\theta}_n = \arg \min_{\theta \in \mathbb{R}^d} \mathbb{E}_{P_n^*} [h_\theta(X, Z)] \longrightarrow \mathbb{E}_P [h_{\hat{\theta}_n}(X, Z)] - \inf_{\theta \in \mathbb{R}^d} \mathbb{E}_P [h_\theta(X, Z)]$$

\uparrow
 (P_X, P_Z)

Perspectives: Zero-Shot Prediction

- Incorporating ERM-based statistical results for models based optimizing foundation modeling objectives.

$$\hat{R} = \min_{R} L_n(R)$$

Perspectives: Zero-Shot Prediction

- Incorporating ERM-based statistical results for models based optimizing foundation modeling objectives.

$$\hat{R} = \min_{R} L_n(R)$$

$$\|\hat{R}_{\text{opt}} - R\|_{L^2(P_X P_Z)}^2 \leq$$

Perspectives: Zero-Shot Prediction

- Incorporating ERM-based statistical results for models based optimizing foundation modeling objectives.

$$\hat{R} = \min_{R} L_n(R)$$

$$\|\hat{R}_{\text{opt}} - R\|_{\mathbf{L}^2(P_X P_Z)}^2 \leq \underbrace{2\|\hat{R}_{\text{opt}} - \hat{R}\|_{\mathbf{L}^2(P_X P_Z)}^2}_{\text{optimization error (non-convex)}} + 2\|\hat{R} - R\|_{\mathbf{L}^2(P_X P_Z)}^2$$

optimization error
(non-convex)

Perspectives: Zero-Shot Prediction

- Incorporating ERM-based statistical results for models based optimizing foundation modeling objectives.

$$\hat{R} = \min_{R} L_n(R)$$

$$\|\hat{R}_{\text{opt}} - R\|_{\mathbf{L}^2(P_X P_Z)}^2 \leq \underbrace{2\|\hat{R}_{\text{opt}} - \hat{R}\|_{\mathbf{L}^2(P_X P_Z)}^2}_{\substack{\text{optimization error} \\ \text{(non-convex)}}} + \underbrace{2\|\hat{R} - R\|_{\mathbf{L}^2(P_X P_Z)}^2}_{\text{statistical error}}$$

$$\|\hat{R} - R\|_{\mathbf{L}^2(P_X P_Z)}^2 \leq C(L(\hat{R}) - L(R))$$

Perspectives: Zero-Shot Prediction

- Incorporating ERM-based statistical results for models based optimizing foundation modeling objectives.

$$\hat{R} = \min_{R} L_n(R)$$

$$\|\hat{R}_{\text{opt}} - R\|_{\mathbf{L}^2(P_X P_Z)}^2 \leq \underbrace{2\|\hat{R}_{\text{opt}} - \hat{R}\|_{\mathbf{L}^2(P_X P_Z)}^2}_{\substack{\text{optimization error} \\ \text{(non-convex)}}} + \underbrace{2\|\hat{R} - R\|_{\mathbf{L}^2(P_X P_Z)}^2}_{\text{statistical error}}$$

$$\begin{aligned} \|\hat{R} - R\|_{\mathbf{L}^2(P_X P_Z)}^2 &\leq C(L(\hat{R}) - L(R)) \\ &\leq 2C \sup_{R} (L_n(R) - L(R)) \end{aligned}$$

Preliminaries

Balanced Pre-Training

Zero-Shot Prediction

Conclusion

Thank you!

To my co-authors, mentors,



and to all my friends and family.

