

# Distribution Alignment

David I. Inouye



Elmore Family School of Electrical  
and Computer Engineering

# Standard ML assumes all data is *relevant*



ML  
User

But what if the data  
is biased or contains  
spurious correlations?

Give me all your data!

Don't worry,  
I'll figure it out.



Machine  
Learning  
System

# What if some data is assumed (or designed) to be *irrelevant*?

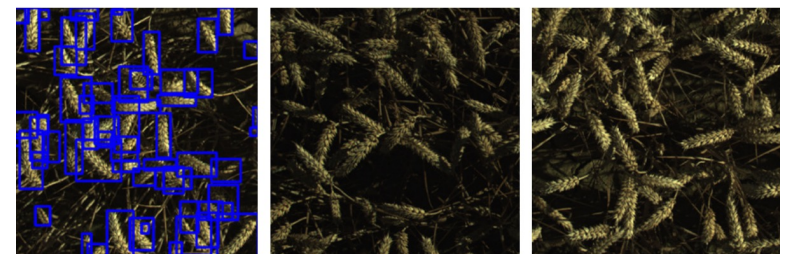
- Fair learning (e.g., FNF [Balunović et al., 2022])
  - Sensitive attributes (e.g., race) are *designed* to be irrelevant for social applications (e.g., loan approval)
- Robust learning (e.g., DANN [Ganin et al., 2016], IRM [Arjovsky et al., 2019])
  - The domain of images (e.g., photo vs sketch) is *assumed* to be irrelevant for object detection



Wheat images from Norway



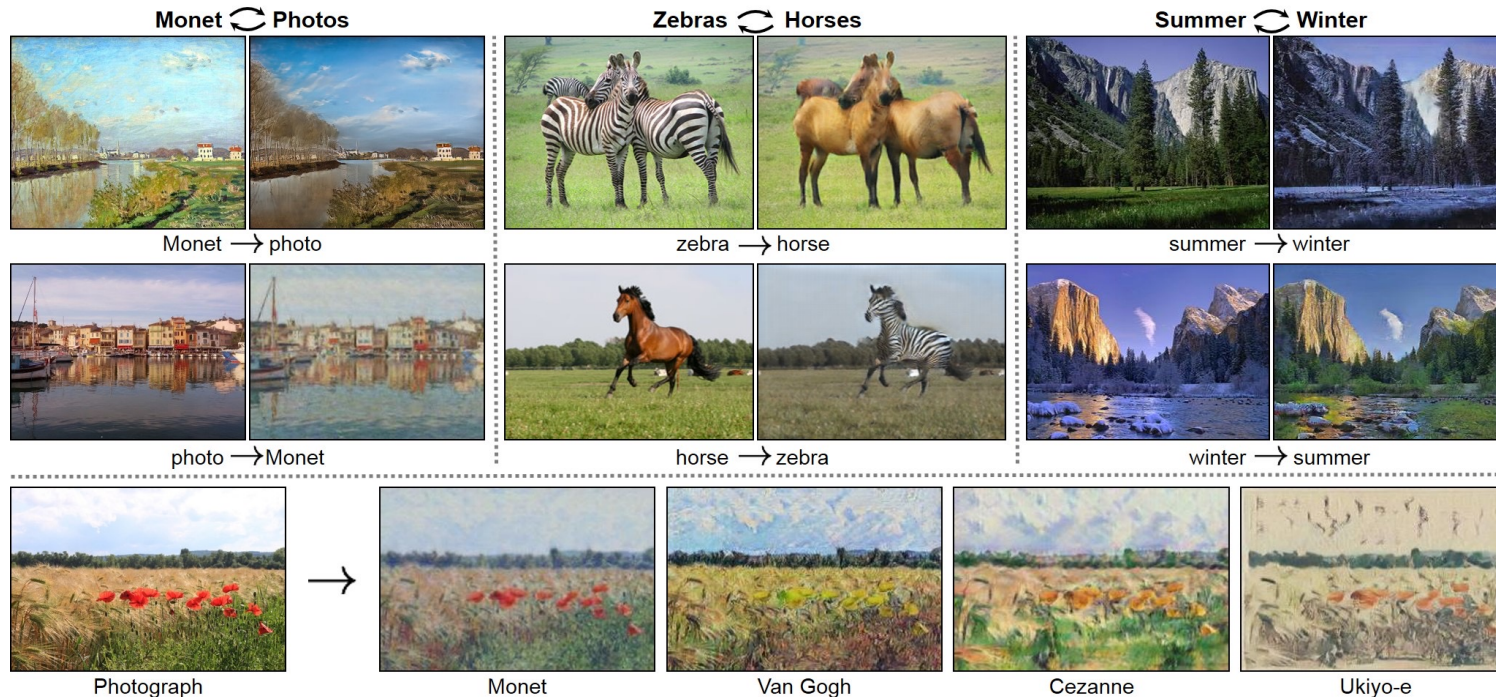
Wheat images from France



Wheat images from Belgium

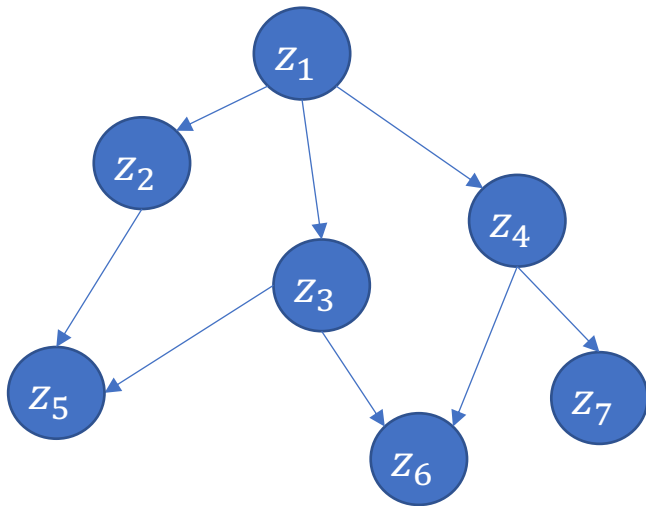
# What if some data is assumed (or designed) to be *irrelevant*?

- Unsupervised translation (e.g., CycleGAN [Zhu et al. 2017])
  - The source of images (i.e., real or generated) is *designed* to be irrelevant



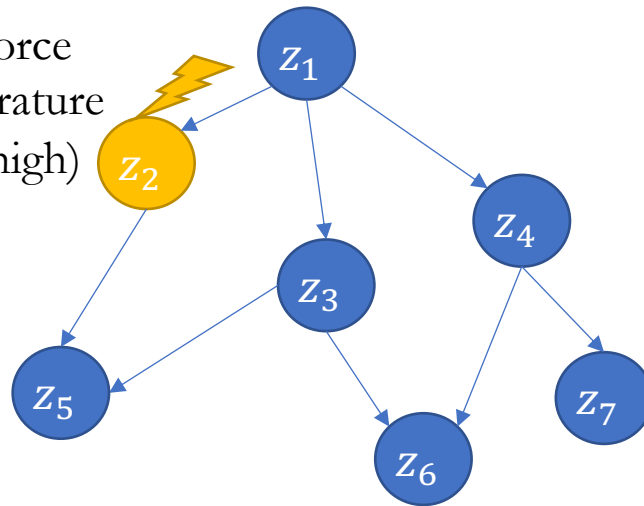
# What if some data is assumed (or designed) to be *irrelevant*?

- Causal discovery (e.g., ICP [Peters et al., 2016])
  - Interventions are *assumed* to be irrelevant for most causal mechanisms



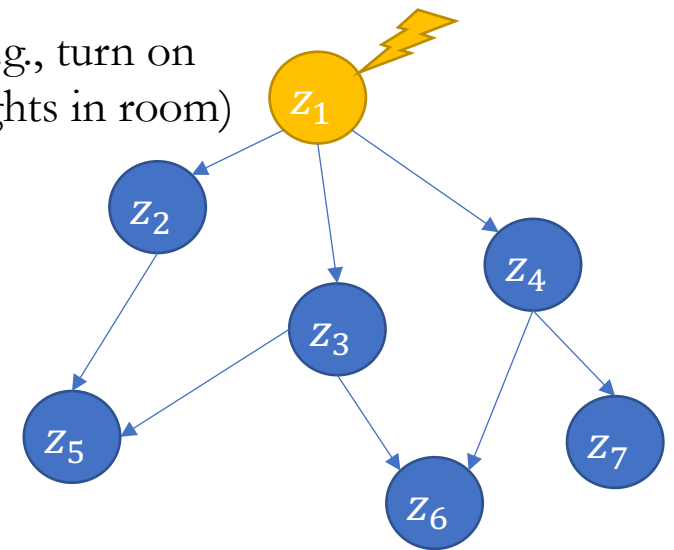
Observed distribution

(e.g., force temperature to be high)



Intervened distribution

(e.g., turn on lights in room)



Another intervened distribution

# How can **known** irrelevant information be used?

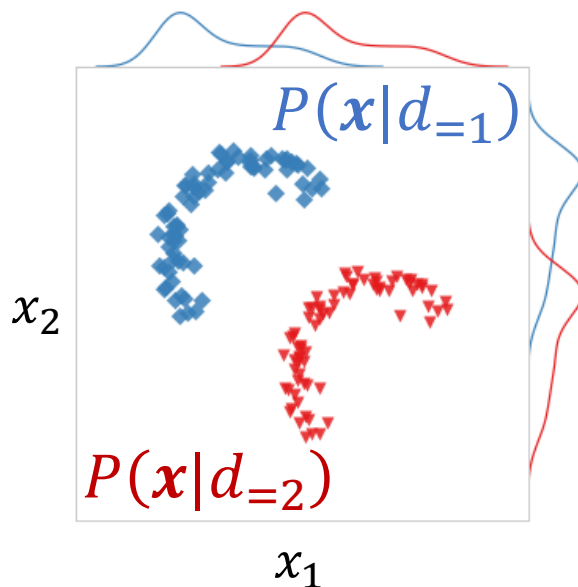
- Simply discard irrelevant features
  - However, other features may contain irrelevant information (e.g., while gender is removed, it can be predicted from an applicant's name)
  - Irrelevant features may be *unknown* or *entangled* with relevant features
- Model design
  - Hope model *implicitly* ignores irrelevant information (i.e., inductive bias)
  - Design model to *explicitly* ignore easy-to-formalize irrelevant information (e.g., graph models that are invariant to node permutations)
- Distribution alignment (this talk 😊 )
  - *Explicitly* minimize irrelevant information (even if infeasible formalize)

# Alignment Concepts



# Distribution alignment is the *opposite* objective of classification

Original Space



Optimization Objective

## Classification

$$\max_g \phi(P(g(x)|d_{=1}), P(g(x)|d_{=2}))$$

where  $g: \mathbb{R}^2 \rightarrow \mathbb{R}$  and  $\phi$  is a distribution divergence (e.g., KL, JSD,  $W_2$ )

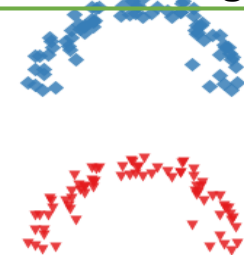
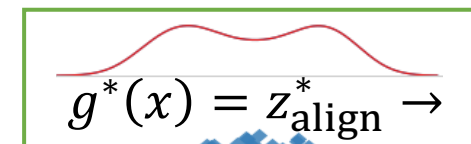
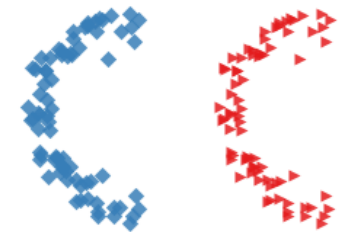
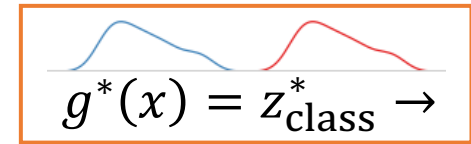
## Distribution alignment

$$\min_g \phi(P(g(x)|d_{=1}), P(g(x)|d_{=2}))$$

Optimal solution

$$P(g^*(x)|d_{=1}) = P(g^*(x)|d_{=2})$$

Latent Space





Alignment can be with respect to the marginal, conditional, or joint distribution

**Marginal alignment**

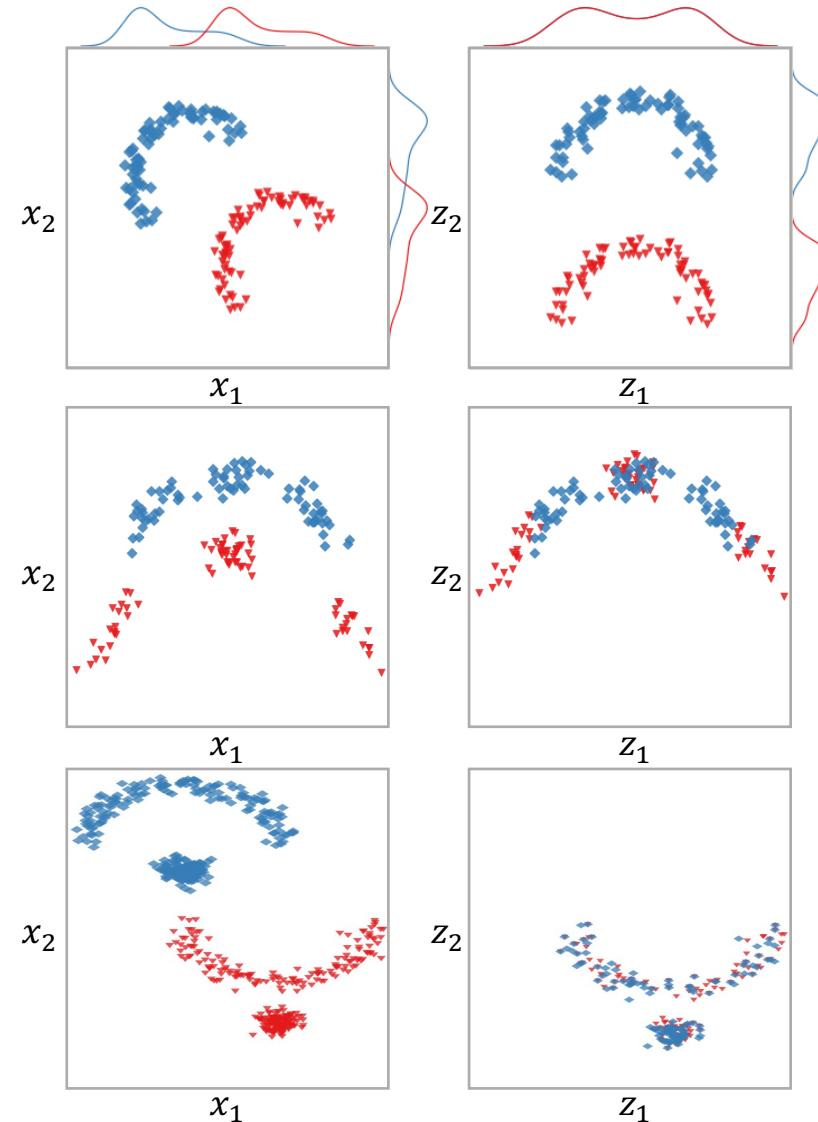
$$P(z_1 | d_{=1}) = P(z_1 | d_{=2})$$

**Conditional alignment**

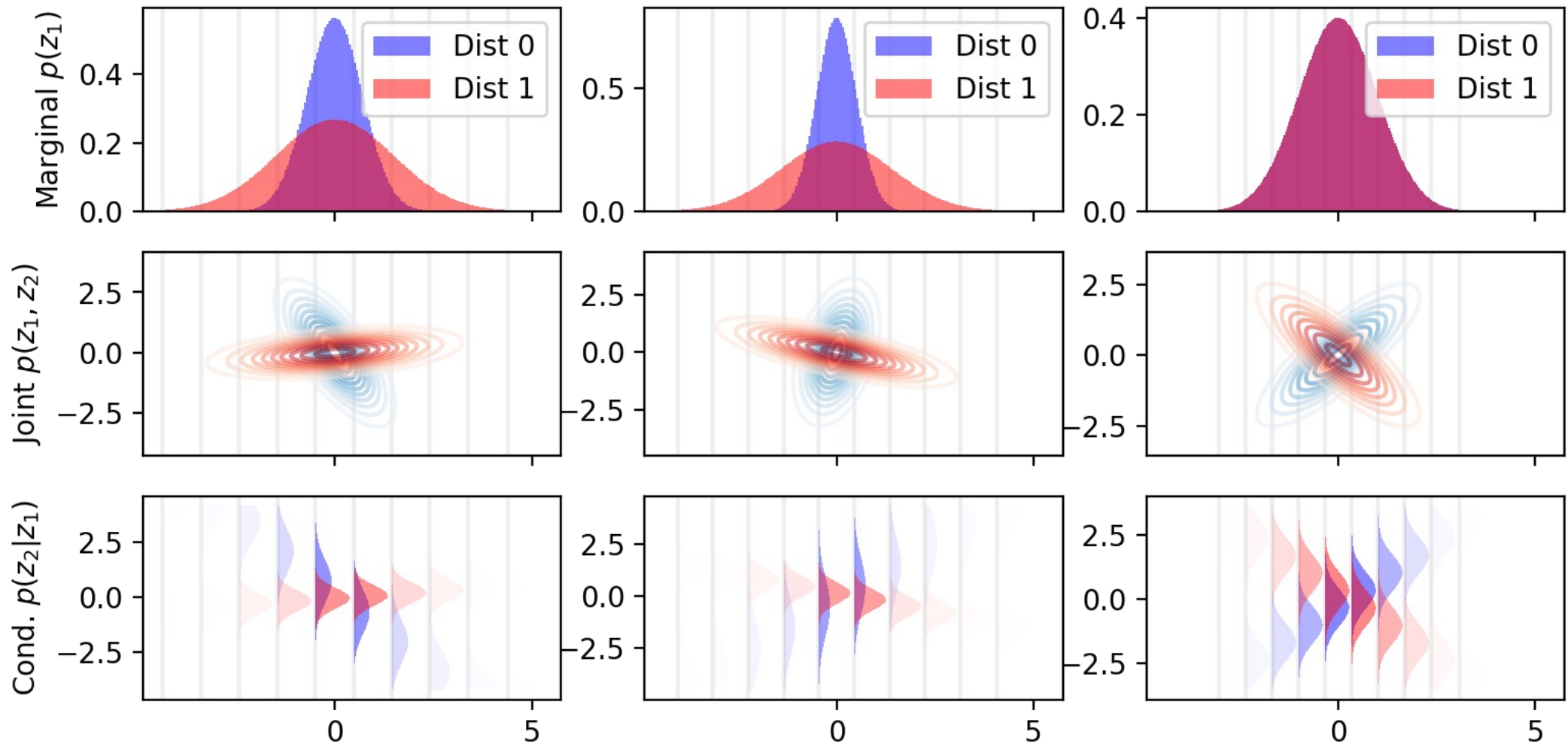
$$P(z_2 | z_1, d_{=1}) = P(z_2 | z_1, d_{=2})$$

**Joint alignment**

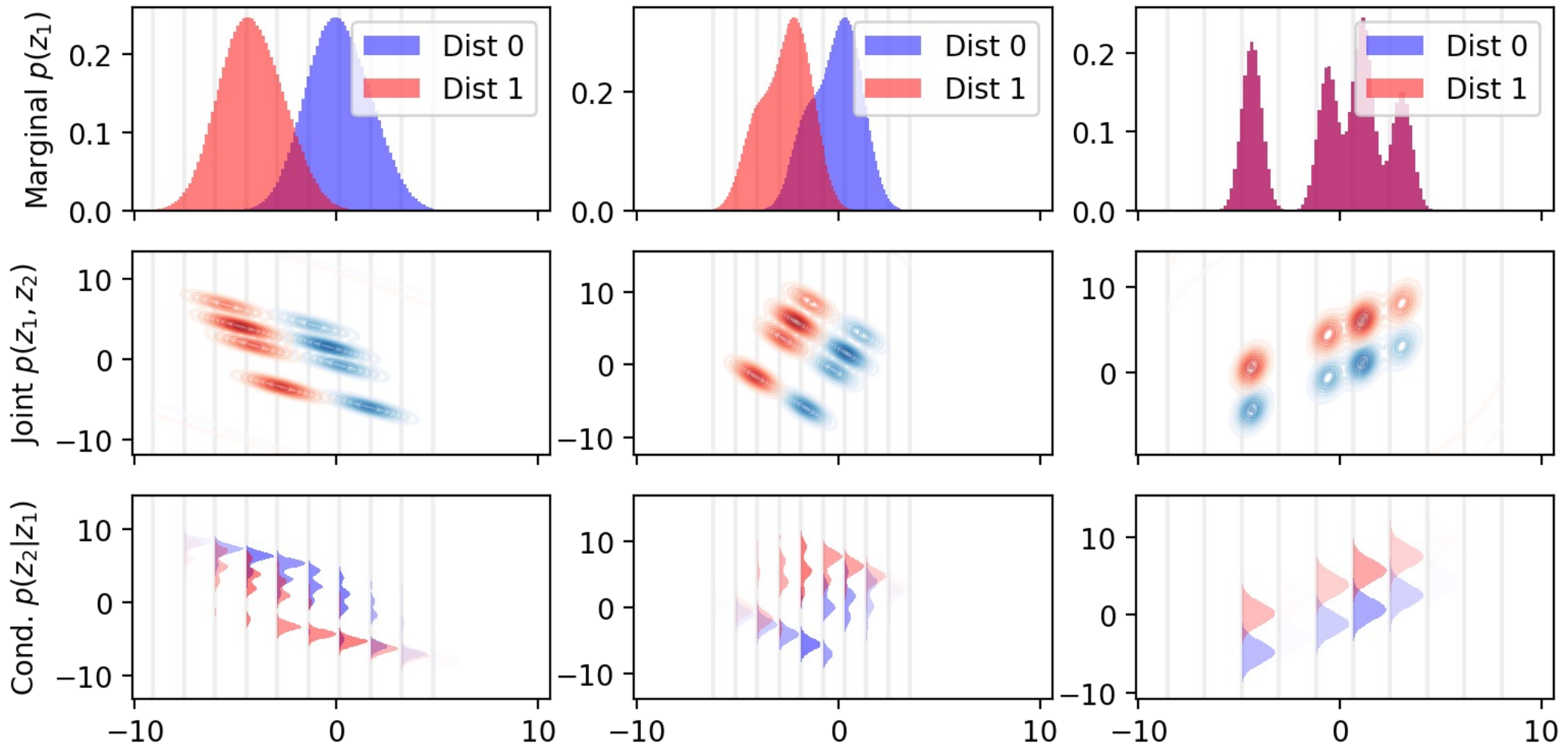
$$P(z_1, z_2 | d_{=1}) = P(z_1, z_2 | d_{=2})$$



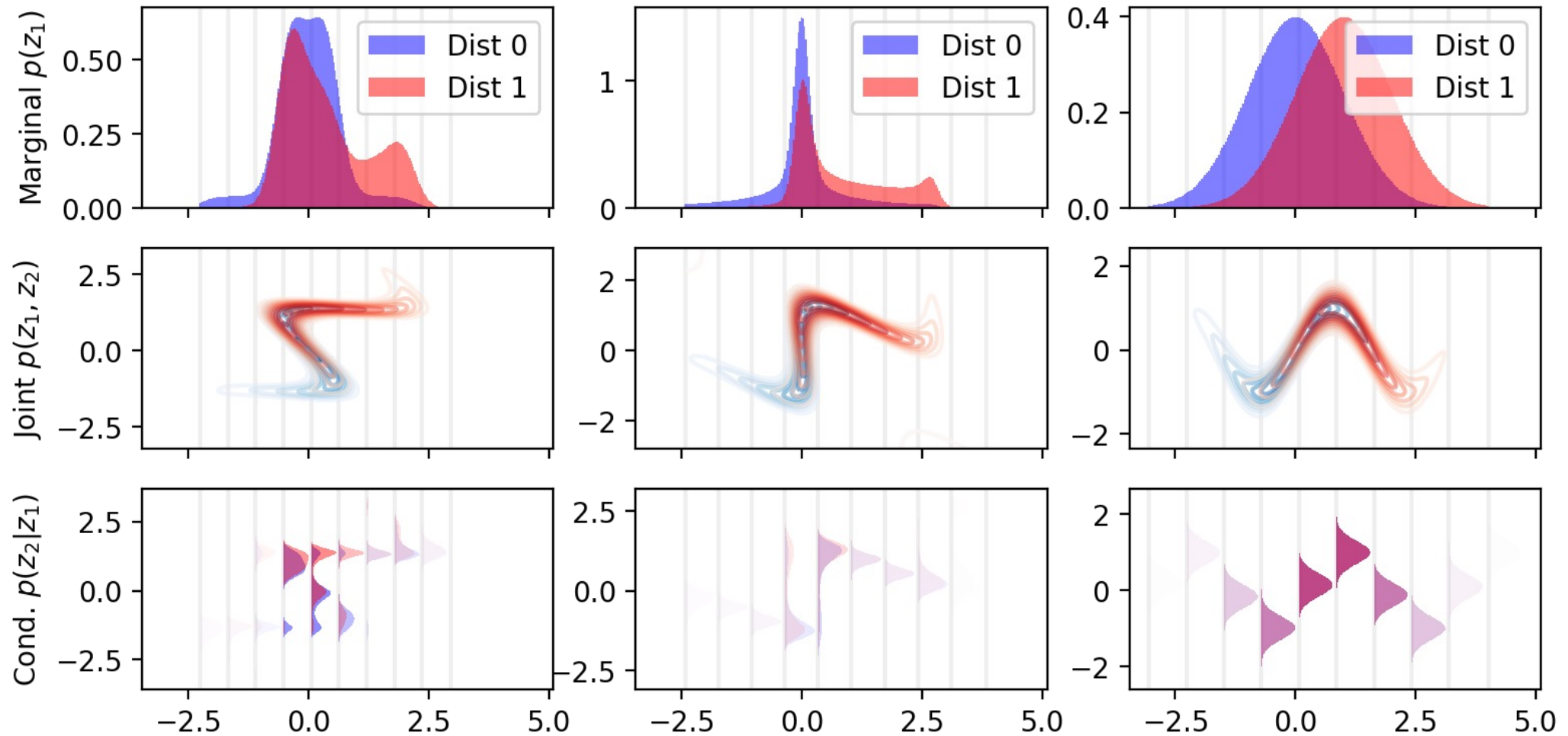
# Example: Marginal alignment without conditional alignment



# Example: Marginal alignment without conditional alignment



# Example: Conditional alignment without marginal alignment



# Distribution alignment minimizes the divergence between two distributions

## Definition 1: Joint Distribution Alignment

Given samples from the joint distribution  $P(\mathbf{x}, d)$ , *distribution alignment* is the problem of finding an *aligner*  $g: \mathcal{X} \times \mathcal{D} \rightarrow \mathcal{Z}$  that minimizes a distribution divergence  $\phi: \mathcal{P} \times \mathcal{P} \rightarrow \mathbb{R}_+$  between the domain-conditional distributions:

$$\min_{g \in \mathcal{G}} \phi(P(\mathbf{z} | d_{=1}), P(\mathbf{z} | d_{=2})), \quad \text{where } \mathbf{z} \equiv g(\mathbf{x}, d).$$

Any distribution divergence that satisfies non-negativity and  $\phi(P, Q) = 0$  if and only if  $P = Q$  (e.g., KL, JSD,  $W_2$ ).

Aligner can depend on domain label  $d$

## Definition 2: Conditional Distribution Alignment

Given two variable index sets  $\mathcal{A}, \mathcal{B} \in \{1, 2, \dots, m\}$ , *conditional alignment* minimizes an aggregation, defined by an aggregator  $\Omega_{\mathcal{Z}_{\mathcal{B}}}[\cdot]$ , over all conditional divergences:

$$\min_{g \in \mathcal{G}} \Omega_{\mathcal{Z}_{\mathcal{B}}} [ \phi(P(\mathbf{z}_{\mathcal{A}} | \mathbf{z}_{\mathcal{B}}, d_{=1}), P(\mathbf{z}_{\mathcal{A}} | \mathbf{z}_{\mathcal{B}}, d_{=2})) ], \quad \text{where } \mathbf{z} \equiv g(\mathbf{x}, d).$$

Usually this is merely the expectation over  $\mathbf{z}_{\mathcal{B}}$ , i.e.,  $\mathbb{E}_{P(\mathbf{z}_{\mathcal{B}})}[\cdot]$

# Constraints on aligners can be explicit or implicit

- Explicit constraints

- *Translation* aligner, i.e.,  $g(\mathbf{x}, d) = \begin{cases} \mathbf{x}, & \text{if } d = 1 \\ \tilde{g}(\mathbf{x}), & \text{otherwise} \end{cases}$

- *Shared* aligner between domains, i.e.,  $g(\mathbf{x}, d) = \tilde{g}(\mathbf{x})$

- *Invertible* aligner, i.e.,  $\exists g^{-1}$  s.t.  $\forall \mathbf{x}, g^{-1}(g(\mathbf{x}, d), d) = \mathbf{x}$

- *Approximately invertible* via cycle consistency  $\exists f$  s.t.  $\forall \mathbf{x}, f(g(\mathbf{x}, d), d) \approx \mathbf{x}$

- Implicit (soft-)constraints via other optimization terms

- We will get to this in **alignment applications**

# These definitions encompass all alignment types under a unified framework

## Marginal alignment

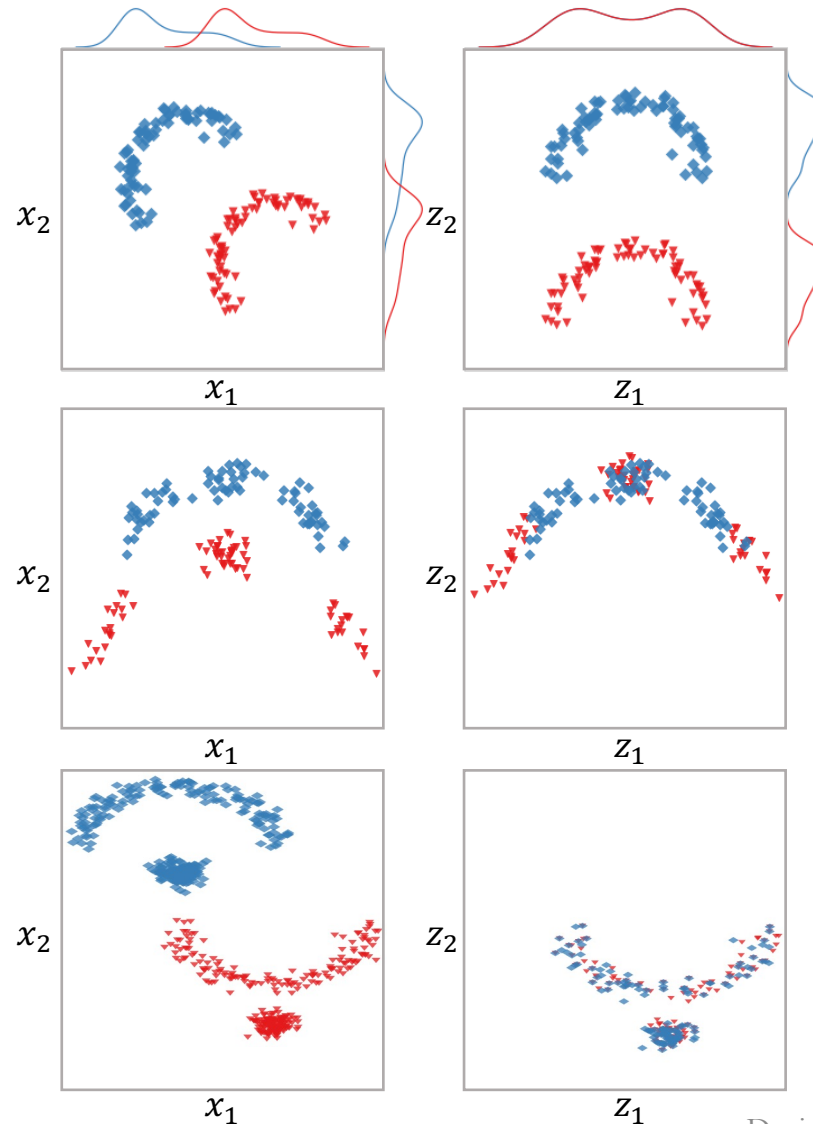
$$P(z_1 | d_{=1}) = P(z_1 | d_{=2})$$

## Conditional alignment

$$P(z_2 | z_1, d_{=1}) = P(z_2 | z_1, d_{=2})$$

## Joint alignment

$$P(z_1, z_2 | d_{=1}) = P(z_1, z_2 | d_{=2})$$



Shared aligner  
 $g(x, d) = Qx$

$$z_{\mathcal{A}} = z_1, z_{\mathcal{B}} = \emptyset$$

Shift only on y-axis

$$g(x, d_{=1}) = x$$

$$g(x, d_{=2}) = x + [0, a]^T$$

$$z_{\mathcal{A}} = z_2, z_{\mathcal{B}} = z_1$$

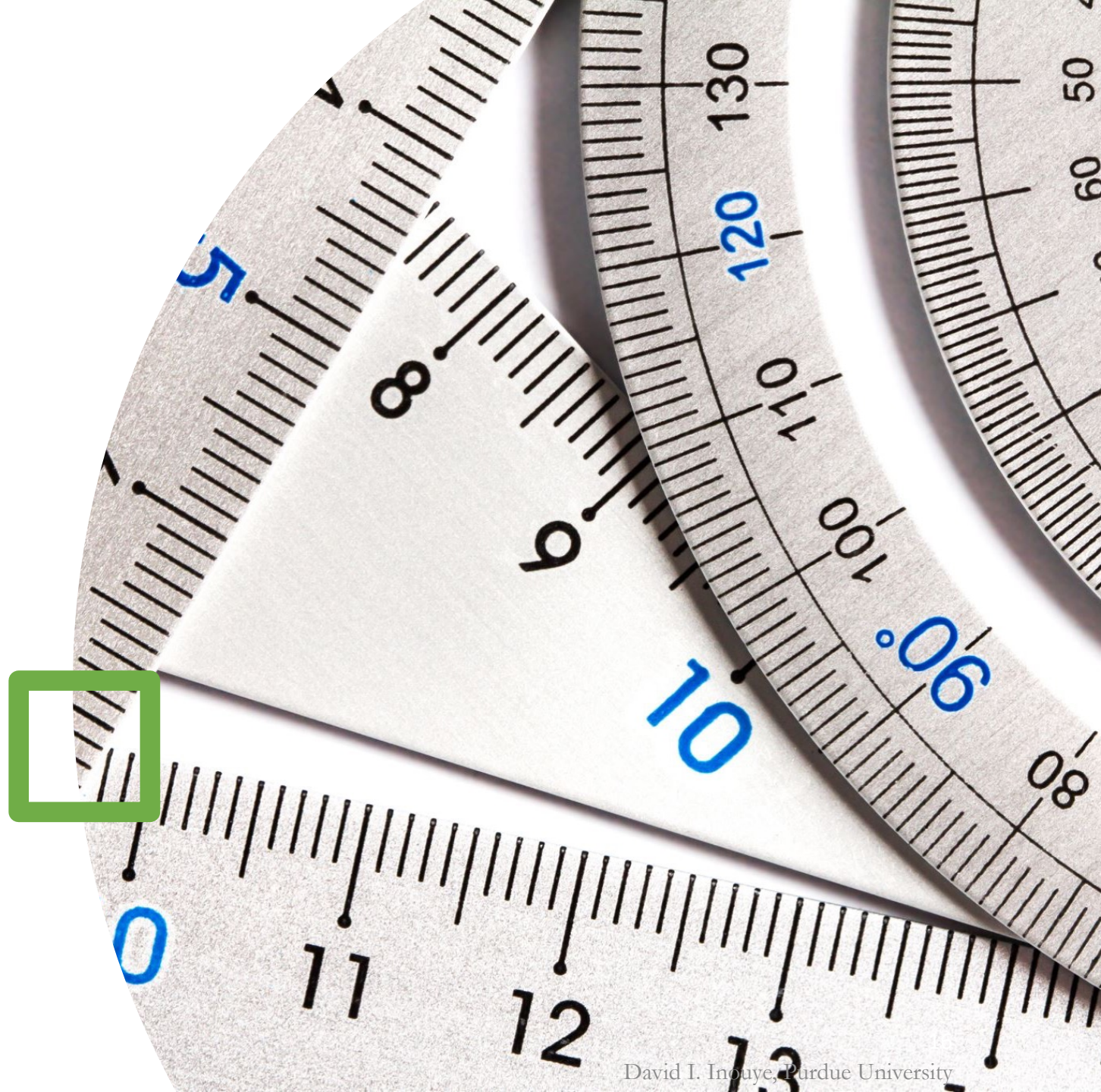
Translation

$$g(x, d_{=1}) = \tilde{g}(x)$$

$$g(x, d_{=2}) = x$$

$$z_{\mathcal{A}} = z, z_{\mathcal{B}} = \emptyset$$

# Tractable Alignment Measures





# *Tractable* alignment measures are needed for evaluation and alignment algorithms

- Two primary uses for alignment measures:
  1. Evaluating or comparing alignment methods
  2. Designing objectives for alignment algorithms (i.e., directly minimize alignment measure)
- While theoretic divergences are elegant (e.g., KL, JSD, TV), most of them are intractable to estimate given only samples
  - Thus, this talk focuses on *tractable* alignment measures

# Extrinsic alignment measures have been used for evaluation (but usually not training)

- External task metric
  - Classification accuracy under fair (alignment) constraints
  - Generalization performance on unseen domain (for domain generalization methods that use feature alignment)
  - **Does not measure alignment explicitly**
- Frechet Inception Distance (FID) or Inception Score (IS)
  - Evaluates quality of images from deep generative models based on latent space of Inception v3 network
  - Perceptual measure of image quality and diversity
  - **Inapplicable for applications with limited data or without a well-established semantic latent space**

# Intrinsic measures are used for training but not for evaluation

- Adversarial measures are variational **lower** bounds of divergences

$$\phi_{GAN}(g) = \max_f \mathbb{E}_{P(\mathbf{x}|d=1)}[\log f(g(\mathbf{x}, 1))] + \mathbb{E}_{P(\mathbf{x}|d=2)}[\log(1 - f(g(\mathbf{x}, 2)))]$$

- If solved perfectly, then  $\phi_{GAN}(g) = JSD(P(g(\mathbf{x}, 1)|d=1), P(g(\mathbf{x}, 2)|d=2)) + const$
  - If non-optimal, then it is **lower** bound.
  - **Difficult for training (min-max/adversarial) and rarely used for evaluation**
- Other intrinsic measure based on Wasserstein distance
    - Empirical optimal transport algorithms – **scales quadratically in number of samples**
    - Sliced Wasserstein distance – **closed-form solution in 1D via sorting**

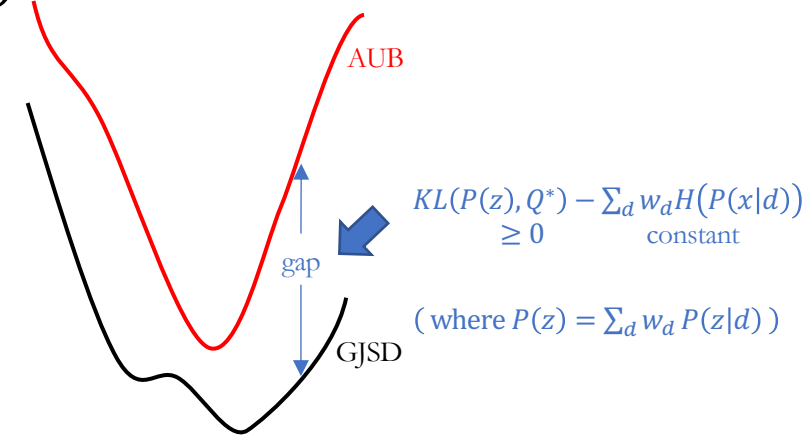
# Alignment Upper Bound (AUB) generalizes alignment measures based on *invertible* models

- A variational **upper** bound of JSD:

$$\phi_{AUB}(g) = \min_{Q \in \mathcal{Q}} \sum_{d=1}^k \mathbb{E}_{P(\mathbf{x}|d)} [-\log |J_{g_d}| Q(g(\mathbf{x}, d))]$$

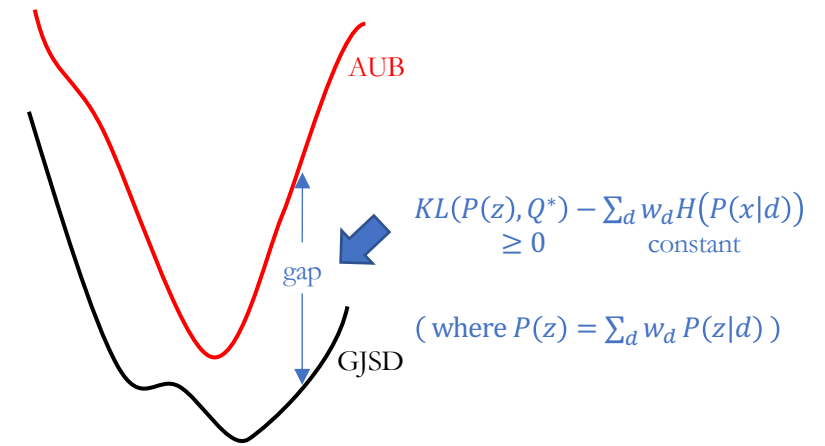
- $Q$  is a density model *shared* among domains
- $g$  is *invertible* and  $|J_{g_d}|$  is the determinant Jacobian of  $g(\cdot, d)$

- **Bound gap** is exactly  $KL(\sum_d w_d P(z|d), Q(z))$
- **Any**  $Q$  provides an **upper** bound on JSD + const
- Alignment is **cooperative**:  $\min_g \phi_{AUB}(g) = \min_g \min_Q \dots$ 
  - The optimal solution aligns the distributions **regardless of**  $Q$



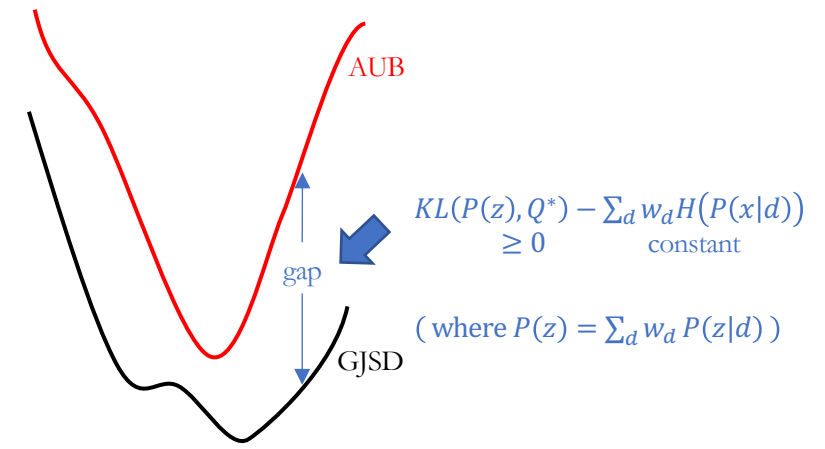
# AUB(1): JSD as entropy of mixture minus mixture of entropies

- $JSD(P(z|d=1), P(z|d=2))$
- $= \sum_d \frac{1}{2} KL(P(z|d), P(z))$  (Let  $P(z) = \sum_d \frac{1}{2} P(z|d)$ , i. e., a mixture)
- $= \sum_d \frac{1}{2} \mathbb{E}_{P(z|d)} \left[ \log \frac{P(z|d)}{P(z)} \right]$
- $= \sum_d \frac{1}{2} \mathbb{E}_{P(z|d)} [-\log P(z)] - \sum_d \frac{1}{2} \mathbb{E}_{P(z|d)} [-\log P(z|d)]$
- $= \sum_d \frac{1}{2} \int_{\mathcal{Z}} P(z|d) (-\log P(z)) dz - \sum_d \frac{1}{2} H(P(z|d))$
- $= \int_{\mathcal{Z}} \sum_d \frac{1}{2} P(z|d) (-\log P(z)) dz - \sum_d \frac{1}{2} H(P(z|d))$
- $= \int_{\mathcal{Z}} P(z) (-\log P(z)) dz - \sum_d \frac{1}{2} H(P(z|d))$
- $= H(P(z)) - \sum_d \frac{1}{2} H(P(z|d))$

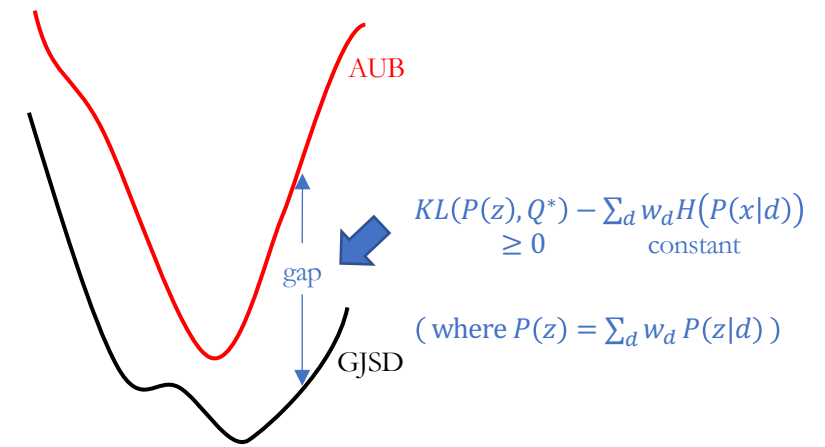


AUB(2): Latent entropy is observed entropy + log determinant term

- $H(P(z|d))$
- $= \mathbb{E}_{P(z|d)} [-\log P(z|d)]$
- $= \mathbb{E}_{P(x|d)} [-\log P(z = g(x, d)|d)]$
- $= \mathbb{E}_{P(x|d)} \left[ -\log P(x|d) |J_{g_d}(x)|^{-1} \right]$
- $= \mathbb{E}_{P(x|d)} [-\log P(x|d)] + \mathbb{E}_{P(x|d)} \left[ -\log |J_{g_d}(x)|^{-1} \right]$
- $= H(P(x|d)) + \mathbb{E}_{P(x|d)} \left[ -\log |J_{g_d}(x)|^{-1} \right]$



AUB(3): Latent cross entropy is weighted observed cross entropy



- $H_c(P(z), Q(z)) \equiv \mathbb{E}_{P(z)}[-\log(Q(z))]$   
 (Note that:  $KL(P, Q) = H_c(P, Q) - H(P)$ )

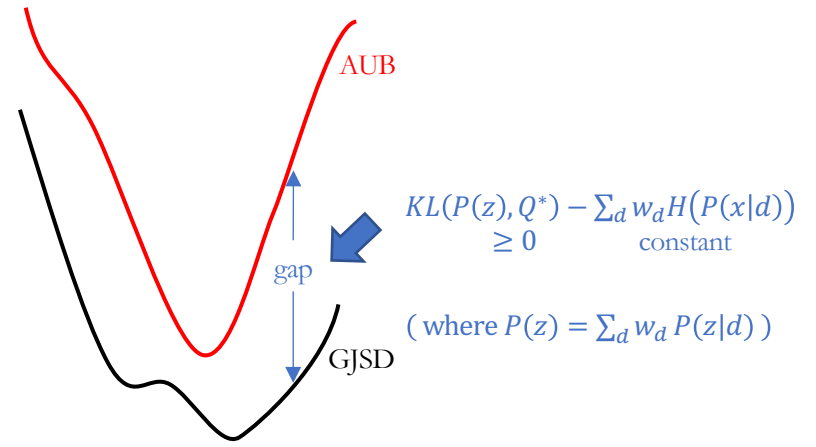
- $= - \int_{\mathcal{Z}} \sum_d \frac{1}{2} P(z|d) \log(Q(z)) dz$

- $= - \sum_d \frac{1}{2} \int_{\mathcal{Z}} P(z|d) \log(Q(z)) dz$

- $= \sum_d \frac{1}{2} \mathbb{E}_{P(z|d)}[-\log(Q(z))]$

- $= \sum_d \frac{1}{2} \mathbb{E}_{P(x|d)}[-\log(Q(g(x, d)))]$

# AUB(4): AUB is upper bound on JSD + const



- $JSD(P(z|d=1), P(z|d=2))$
- $= H(P(z)) - \sum_d \frac{1}{2} H(P(z|d))$
- $= H_c(P(z), Q(z)) - H_c(P(z), Q(z)) + H(P(z)) - \sum_d \frac{1}{2} H(P(z|d))$
- $= H_c(P(z), Q(z)) - KL(P(z), Q(z)) - \sum_d \frac{1}{2} H(P(z|d))$
- $\leq H_c(P(z), Q(z)) - \sum_d \frac{1}{2} H(P(z|d))$
- $= \sum_d \frac{1}{2} \mathbb{E}_{P(x|d)} \left[ -\log(Q(g(x, d))) \right] - \sum_d \frac{1}{2} \left( \mathbb{E}_{P(x|d)} [\log |J_{g_d}(x)|] + H(P(x|d)) \right)$
- $= \sum_d \frac{1}{2} \mathbb{E}_{P(x|d)} \left[ -\log(|J_{g_d}(x)| Q(g(x, d))) \right] - \sum_d \frac{1}{2} H(P(x|d))$

Constant w.r.t  $g$



# Alignment Algorithms

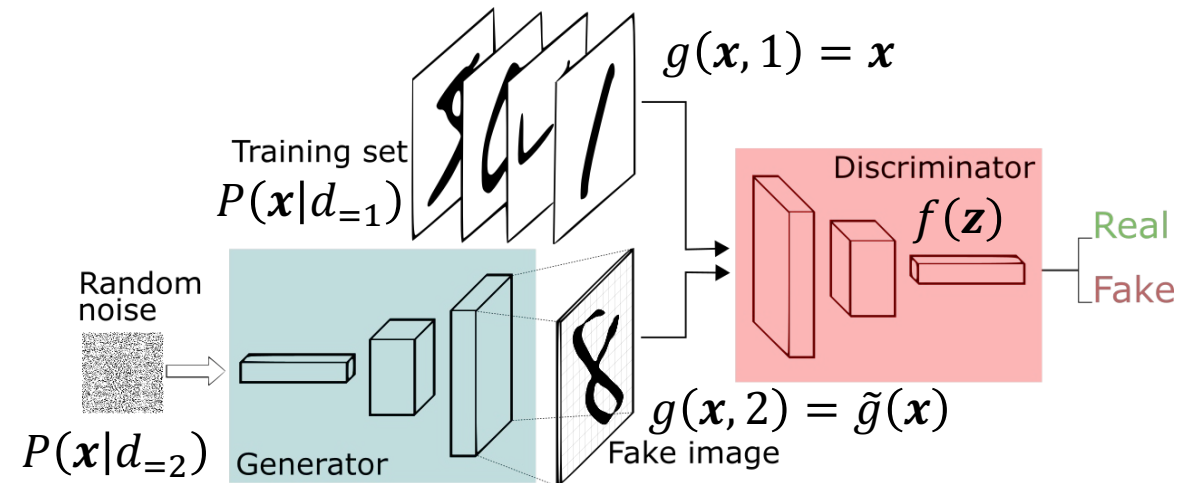
---

# Adversarial optimization (GAN-inspired) is the standard approach to alignment

- Intuition – Competitive game
  - Counterfeiter is trying to avoid getting caught
  - Police is trying to catch counterfeiter
- Algorithm – Usually alternating optimization between min and max
- Benefits
  - No constraints on generator and discriminator models
- Drawbacks
  - Lacks domain-agnostic evaluation metrics (e.g., unable to check for overfitting)
  - Unstable or poorly conditioned optimization

## Adversarial alignment problem

$$\min_g \max_f \mathbb{E}_{P(\mathbf{x}|d=1)} [\log f(g(\mathbf{x}, 1))] + \mathbb{E}_{P(\mathbf{x}|d=2)} [\log (1 - f(g(\mathbf{x}, 2)))]$$

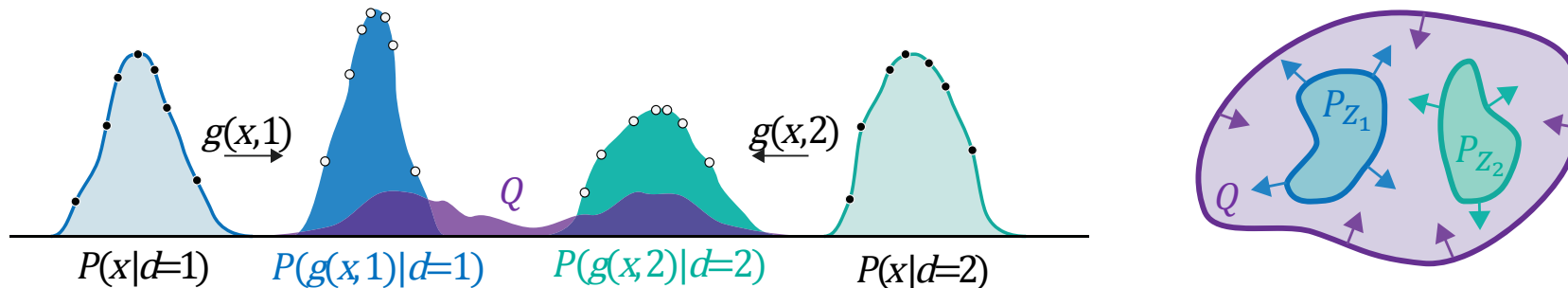


<https://www.freecodecamp.org/news/an-intuitive-introduction-to-generative-adversarial-networks-gans-7a2264a81394/>

# AUB optimization provides a **cooperative** alternative to adversarial alignment

**AUB cooperative alignment problem**

$$\min_g \min_{Q \in \mathcal{Q}} \sum_{j=1}^k \mathbb{E}_{P(\mathbf{x}|d)} [\log |J_{g_d}| Q(g(\mathbf{x}, d))]$$



- Minimizing  $g$  makes distributions closer to current  $Q$  (left)
- Minimizing  $Q$  tightens bound by getting closer to the latent mixture, i.e.,  $\sum_d P(g(x, d)|d)$  (right)

# AUB can perform alignment on tabular data and between multiple domains

	MINIBOONE (42)	GAS (7)	HEPMASS (20)	POWER (5)
LRMF	12.79	-6.17	18.49	-0.93
AF (MLE)	14.08	-6.52	19.37	-0.77
AF (Adv. only)	18.18	-3.15	21.70	-0.39
AF (hybrid)	19.49	-3.76	21.42	-0.43
Ours	<b>12.11</b>	<b>-7.09</b>	<b>18.26</b>	<b>-1.19</b>



AlignFlow (MLE)

Ours

These results on 4 benchmark tabular datasets demonstrate that our algorithm can improve the AUB alignment measure on test data.

Our AUB algorithm can translate between 10 domains (MNIST digits here) better than the closest competitor (AlignFlow) for invertible models. (Original real digits are far left and grid is translations to all other digits.)

# Iterative alignment flows iteratively solve 1D alignment problems to create deep aligner

1. Find 1D projection that is **maximally misaligned** (i.e., max sliced Wasserstein distance)

$$\max_{\theta} W_2(P(\theta^T \mathbf{x}|d_{=1}), P(\theta^T \mathbf{x}|d_{=2}))$$

2. Align along this 1D projection by mapping to barycenter distribution

$$\min \mathbb{E}_{P(\tilde{x}, d)} [\|g(\tilde{x}, d) - \tilde{x}\|^2]$$

$$\text{s.t. } \tilde{x} = \theta^T \mathbf{x}, \quad P(g(\tilde{x}, 1)|d_{=1}) = P(g(\tilde{x}, 2)|d_{=2})$$

3. Update aligner (add one layer) and repeat

$$\tilde{g}(\mathbf{x}) = g(\theta^T \mathbf{x}, d)\theta + \mathbf{x}_{\perp}$$

$$\tilde{g}_{\text{global}}^{\text{new}} = \tilde{g} \circ \tilde{g}_{\text{global}}^{\text{old}}$$

$$\mathbf{x}^{\text{new}} = \tilde{g}(\mathbf{x})$$

# INB is significantly faster than the closest invertible model baselines

	Model	WD	FID	TC	Time(s)
Ours	NB	60.010 $\pm$ 0.000	229.551 $\pm$ 0.000	<b>28.115 <math>\pm</math> 0.000</b>	<b>25</b>
	INB ( $L = 20$ )	23.481 $\pm$ 0.161	43.196 $\pm$ 0.588	31.671 $\pm$ 0.056	430
	INB ( $L = 250$ )	<b>23.183 <math>\pm</math> 0.095</b>	<b>37.480 <math>\pm</math> 0.008</b>	32.841 $\pm$ 0.097	2200
Iterative Baselines	DD	39.079 $\pm$ 0.000	166.320 $\pm$ 0.000	235.164 $\pm$ 0.000	360
	SINF-Align( $0 \Rightarrow 1$ )	50.151 $\pm$ 0.950	247.142 $\pm$ 0.972	—	50
	SINF-Align( $1 \Rightarrow 0$ )	42.658 $\pm$ 1.253	202.058 $\pm$ 1.716	—	50
Deep Baselines	AlignFlow( $\lambda = 1e-4$ )	56.386	158.654	392.578	220000
	AlignFlow( $\lambda = 1e-5$ )	60.452	191.983	412.531	220000

# Alignment Applications

---

Alignment applications can be unified as a task objective + *(soft) alignment constraints*

## Task objective

- “What we want”
- Relevant information

## Alignment constraints

- “What we don’t want”
- Irrelevant information



# Fair classification aims to classify correctly while controlling for sensitive attributes

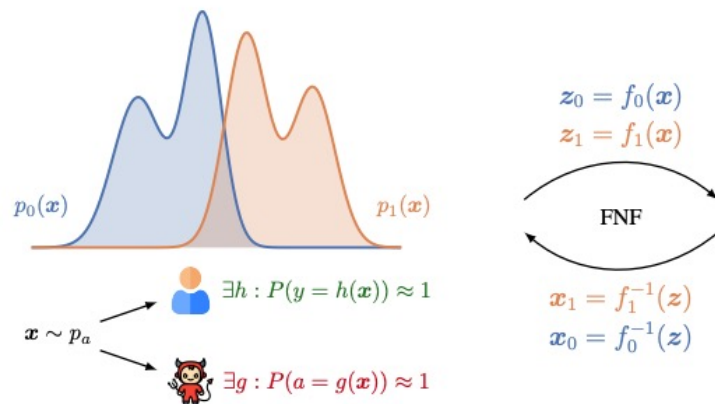
Task objective:

“What we want” / “relevant”

- Accurately predict whether a loan application should be approved
- Standard classification loss

$$\mathbb{E}_{P(\mathbf{x}, d)} [\ell(f(g(\mathbf{x}, d)), y)]$$

Raw representation is good for task but sensitive attribute can be determined

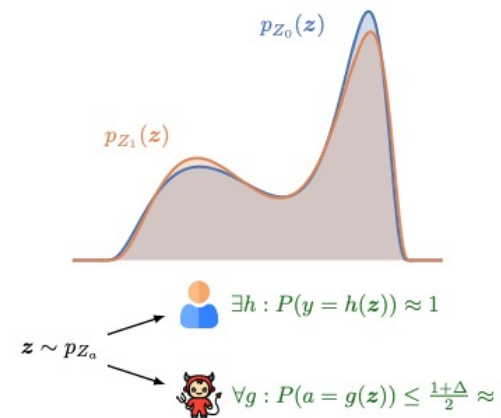


(Soft) alignment constraints:

“What we don’t want” / “irrelevant”

- The prediction must be *independent* of sensitive attribute  $d$
- Alignment constraint loss

$$\phi(P(g(\mathbf{x}, d)|d=1), P(g(\mathbf{x}, d)|d=2))$$



Aligned representation is good for task but sensitive attribute **cannot** be determined

# Unsupervised image-to-image translation aims to preserve content while changing domains

Task objective:

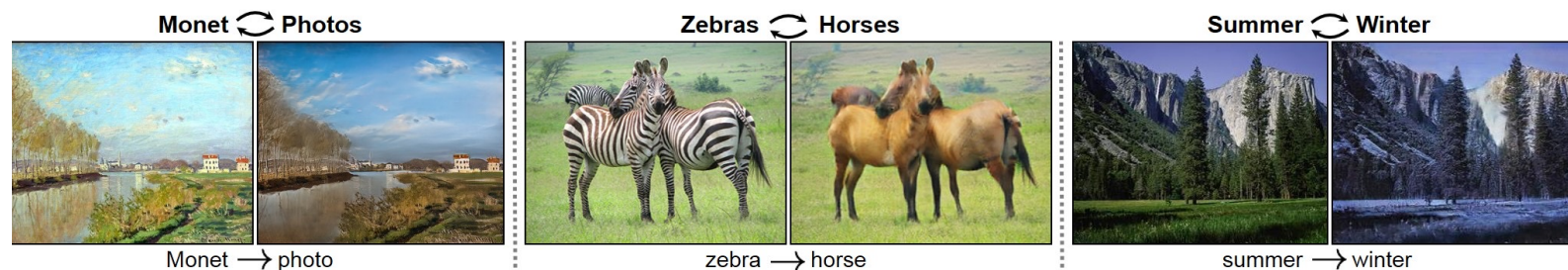
“What we want” / “relevant”

- Preserve semantic image content
- Both explicit and implicit methods (e.g., CycleGAN)
  - Cycle consistency loss (explicit)
  - Identity regularization (explicit)
  - CNN architecture (implicit)

(Soft) alignment constraints:

“What we don’t want” / “irrelevant”

- Change the style (or domain) of the image
  - Translated image should “look like” images from the other domain
- Alignment constraint loss
$$\phi(P(g(\mathbf{x}, d)|d_1), P(\mathbf{x}|d_2))$$



# Background: Causal probabilistic models *implicitly* encode the effect of **interventions**



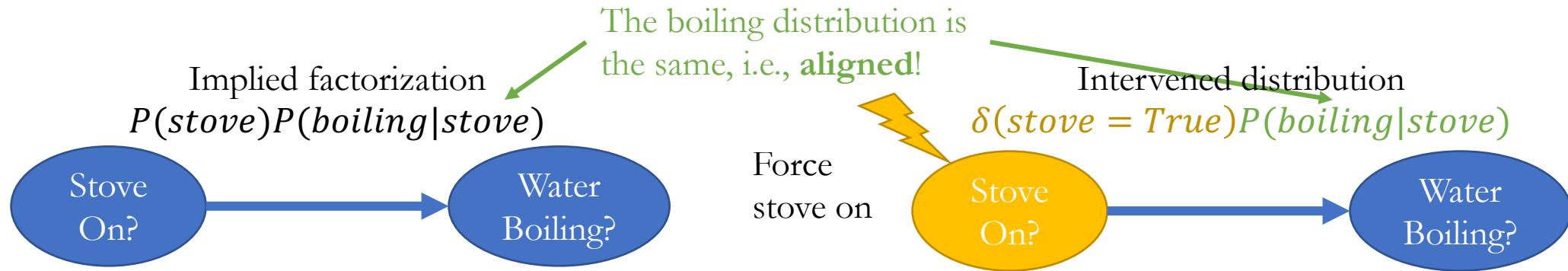
Both are valid factorizations.  
 But which factorization is *causal*?

One idea: The factorization that changes the least under an intervention.



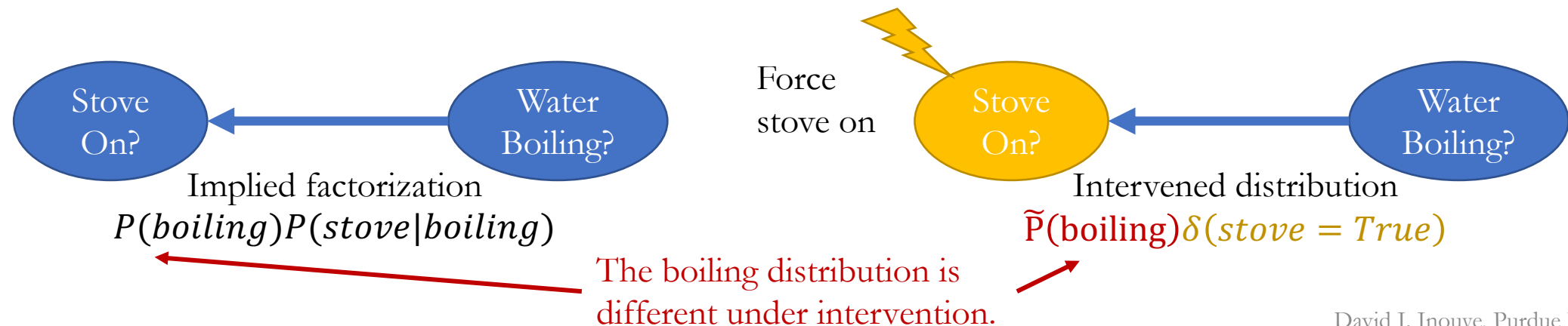
The stove distribution is different under intervention.  $\tilde{P}(\text{stove}|\text{boiling}) \equiv P(\text{stove})$

# Background: Causal probabilistic models *implicitly* encode the effect of **interventions**

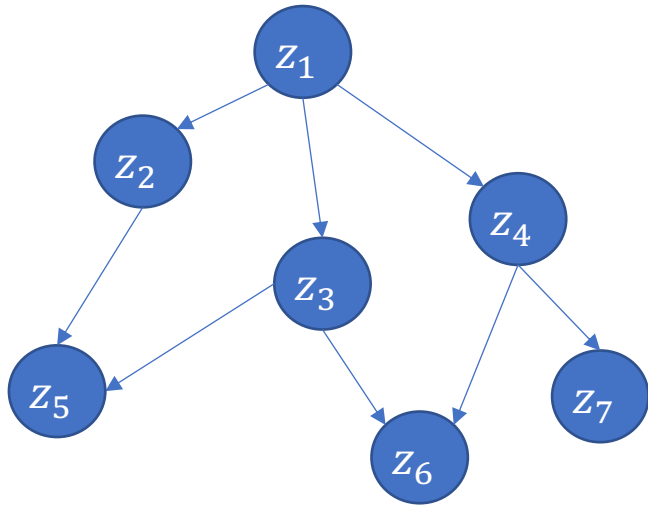


Both are valid factorizations.  
But which factorization is *causal*?

One idea: The factorization that changes the least under an intervention.

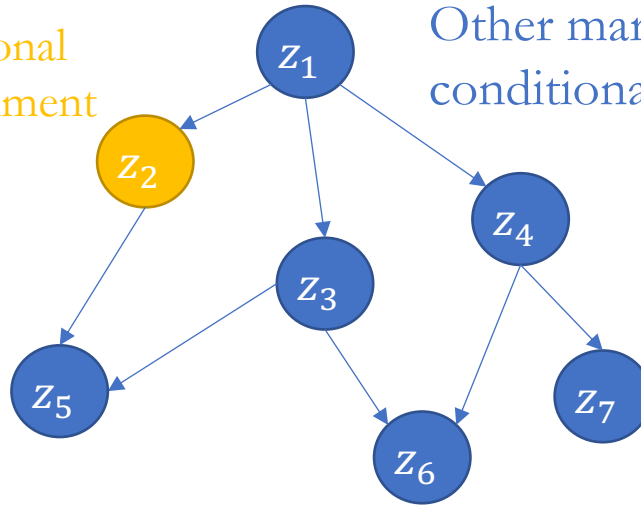


# Different domains can be viewed as *unknown* interventions in a *latent* causal space



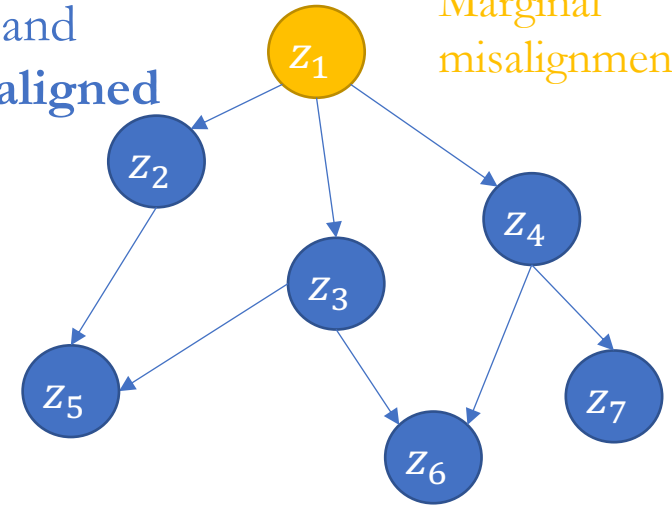
Latent space  $\mathbf{z}|d_{=0} \sim \text{CausalModel}$

Conditional misalignment



$\mathbf{z}|d_{=1} \sim \text{IntervenedCausalModel}$

Other marginals and conditionals are aligned



$\mathbf{z}|d_{=2} \sim \text{IntervenedCausalModel}$

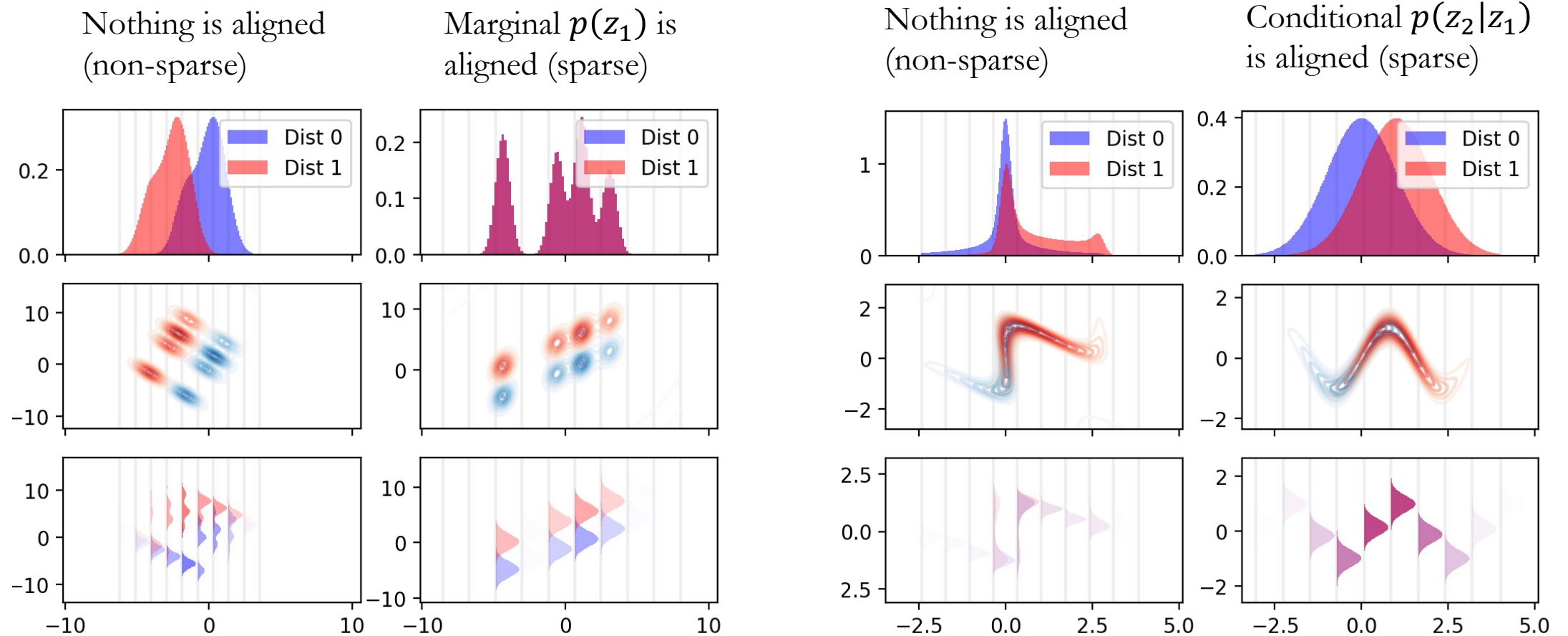
Marginal misalignment

Observed space  $\mathbf{x} = g^{-1}(\mathbf{z})$



# Sparse intervention assumption $\Rightarrow$ misalignment sparsity (Only a few conditionals are misaligned)

In 2D this means that either the marginal or conditionals are misaligned but **not both**.



# Future Vision: Alignment as (soft) constraints to combat underspecification in deep learning

Underspecification Presents Challenges for Credibility in Modern Machine Learning

Alexander D'Amour\*  
 Katherine Heller\*  
 Dan Moldovan\*  
 Ben Adlam  
 Babak Alipanahi  
 Alex Beutel  
 Christina Chen  
 Jonathan Deaton  
 Jacob Eisenstein  
 Matthew D. Hoffman  
 Farhad Hormozdiari  
 Neil Houlsby  
 Shaobo Hou  
 Ghassen Jerfel  
 Alan Karthikesalingam  
 Mario Lucic  
 Yian Ma  
 Cory McLean  
 Diana Mincu  
 Akinori Mitani  
 Andrea Montanari  
 Zachary Nado  
 Vivek Natarajan  
 Christopher Nielson†  
 Thomas F. Osborne†  
 Rajiv Raman  
 Kim Ramasamy  
 Rory Sayres  
 Jessica Schrouff  
 Martin Seneviratne  
 Shannon Sequeira  
 Harini Suresh  
 Victor Veitch  
 Max Vladymyrov  
 Xuezhi Wang  
 Kellie Webster  
 Steve Yadlowsky  
 Taedong Yun  
 Xiaohua Zhai  
 D. Sculley

ALEXDAMOUR@GOOGLE.COM  
 KHELLER@GOOGLE.COM  
 MDAN@GOOGLE.COM  
 ADLAM@GOOGLE.COM  
 BABAKA@GOOGLE.COM  
 ALEXBEUTEL@GOOGLE.COM  
 CHRISTINIUM@GOOGLE.COM  
 JDEATON@GOOGLE.COM  
 JEISENSTEIN@GOOGLE.COM  
 MHOFFMAN@GOOGLE.COM  
 FHORMOZ@GOOGLE.COM  
 NEILHOULSBY@GOOGLE.COM  
 SHAOBOHOU@GOOGLE.COM  
 GHASSEN@GOOGLE.COM  
 ALANKARTHI@GOOGLE.COM  
 LUCIC@GOOGLE.COM  
 YIANMA@UCSD.EDU  
 CYM@GOOGLE.COM  
 DMINCU@GOOGLE.COM  
 AMITANI@GOOGLE.COM  
 MONTANARI@STANFORD.EDU  
 ZNADO@GOOGLE.COM  
 NATVIV@GOOGLE.COM  
 CHRISTOPHER.NIELSON@VA.GOV  
 THOMAS.OSBORNE@VA.GOV  
 DRRRN@SNMAIL.ORG  
 KIM@ARAVIND.ORG  
 SAYRES@GOOGLE.COM  
 SCHROUFF@GOOGLE.COM  
 MARTSEN@GOOGLE.COM  
 SHNNN@GOOGLE.COM  
 HSURESH@MIT.EDU  
 VICTORVEITCH@GOOGLE.COM  
 MXV@GOOGLE.COM  
 XUEZHIV@GOOGLE.COM  
 WEBSTERK@GOOGLE.COM  
 YADLOWSKY@GOOGLE.COM  
 TEDYUN@GOOGLE.COM  
 XZHAI@GOOGLE.COM  
 DSCULLEY@GOOGLE.COM

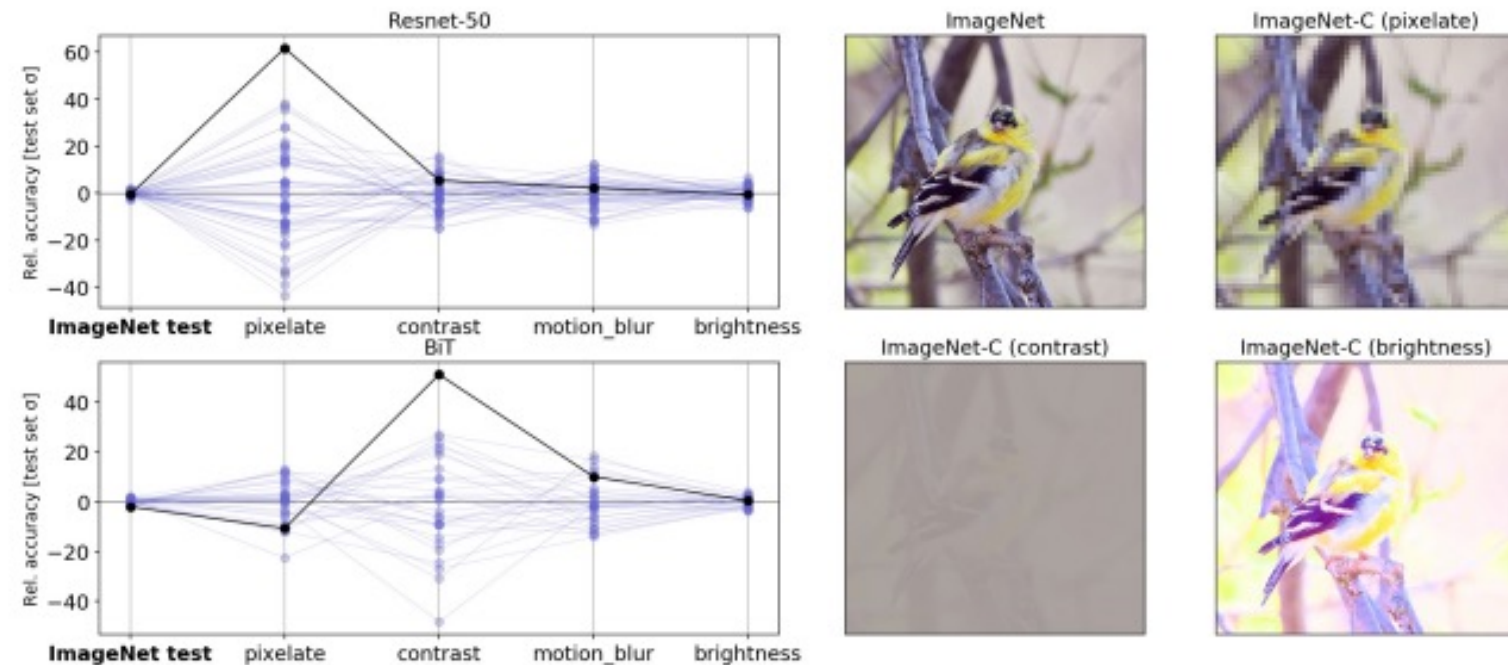


Figure 4: Image classification model performance on stress tests is sensitive to random initialization in ways that are not apparent in iid evaluation. (Top Left)

# Future research opportunities in all areas of distribution alignment

- Alignment concepts
  - Conditional alignment in particular
- Alignment measures
  - More application-agnostic measures
  - Rigorous evaluation protocols
- Alignment algorithms
  - Beyond adversarial
  - More stable optimization
- Alignment applications
  - What robustness can we achieve?
  - Can we make this more general?