



Adversarial Support Alignment

Shangyuan Tong¹ Timur Garipov¹

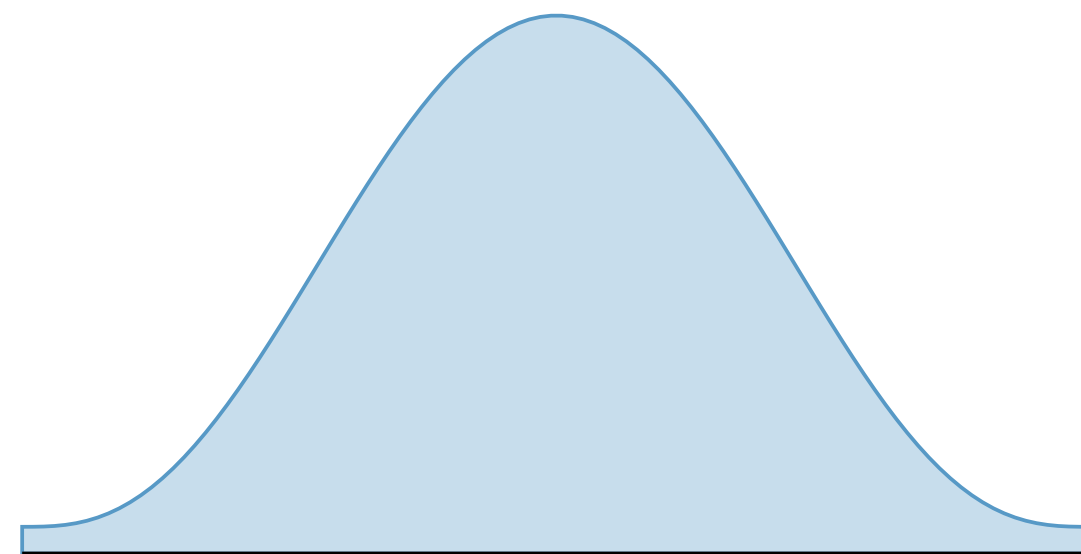
Yang Zhang² Shiyu Chang³ Tommi Jaakkola¹

¹MIT CSAIL ²MIT-IBM Watson AI Lab ³UC Santa Barbara

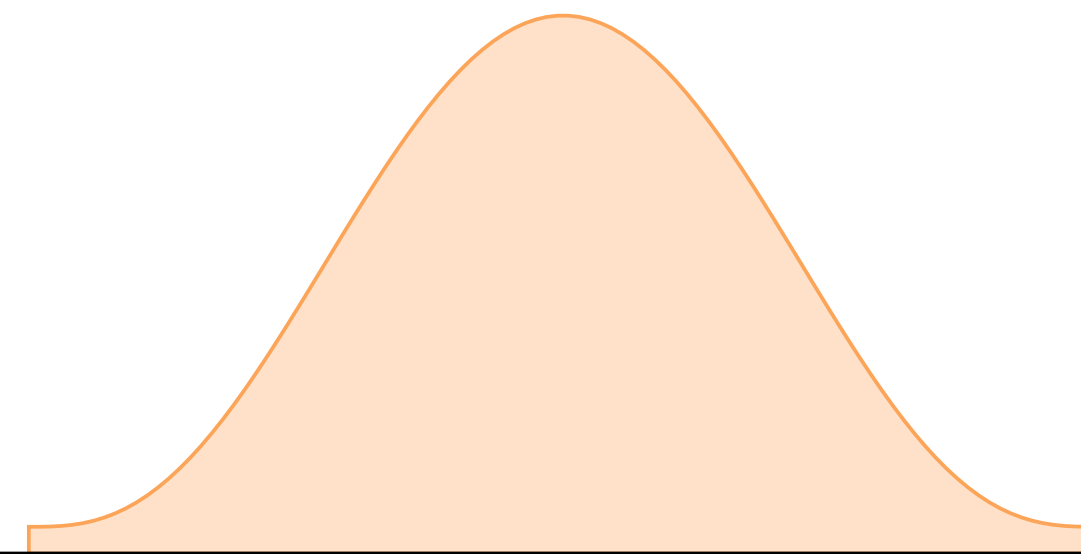
Background: distribution alignment

Given $\mathcal{P} = \{p^\theta \mid \theta \in \Theta\}$ $\mathcal{Q} = \{q^\theta \mid \theta \in \Theta\}$

$p^\theta(x)$



$q^\theta(x)$

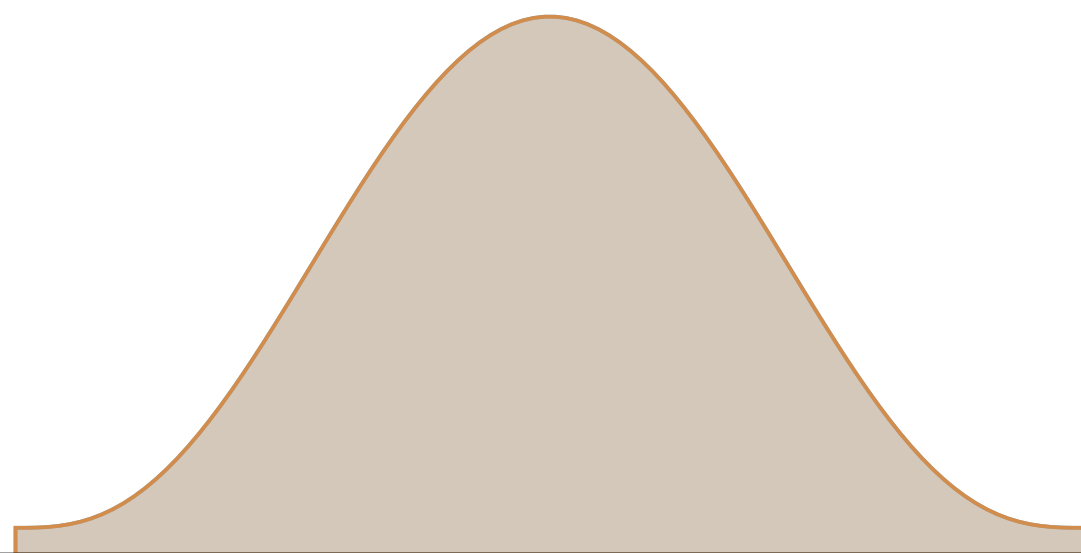


Background: distribution alignment

Given $\mathcal{P} = \{p^\theta \mid \theta \in \Theta\}$ $\mathcal{Q} = \{q^\theta \mid \theta \in \Theta\}$

Find $\theta^* : p^{\theta^*} = q^{\theta^*}$

$$p^{\theta^*} = q^{\theta^*}$$



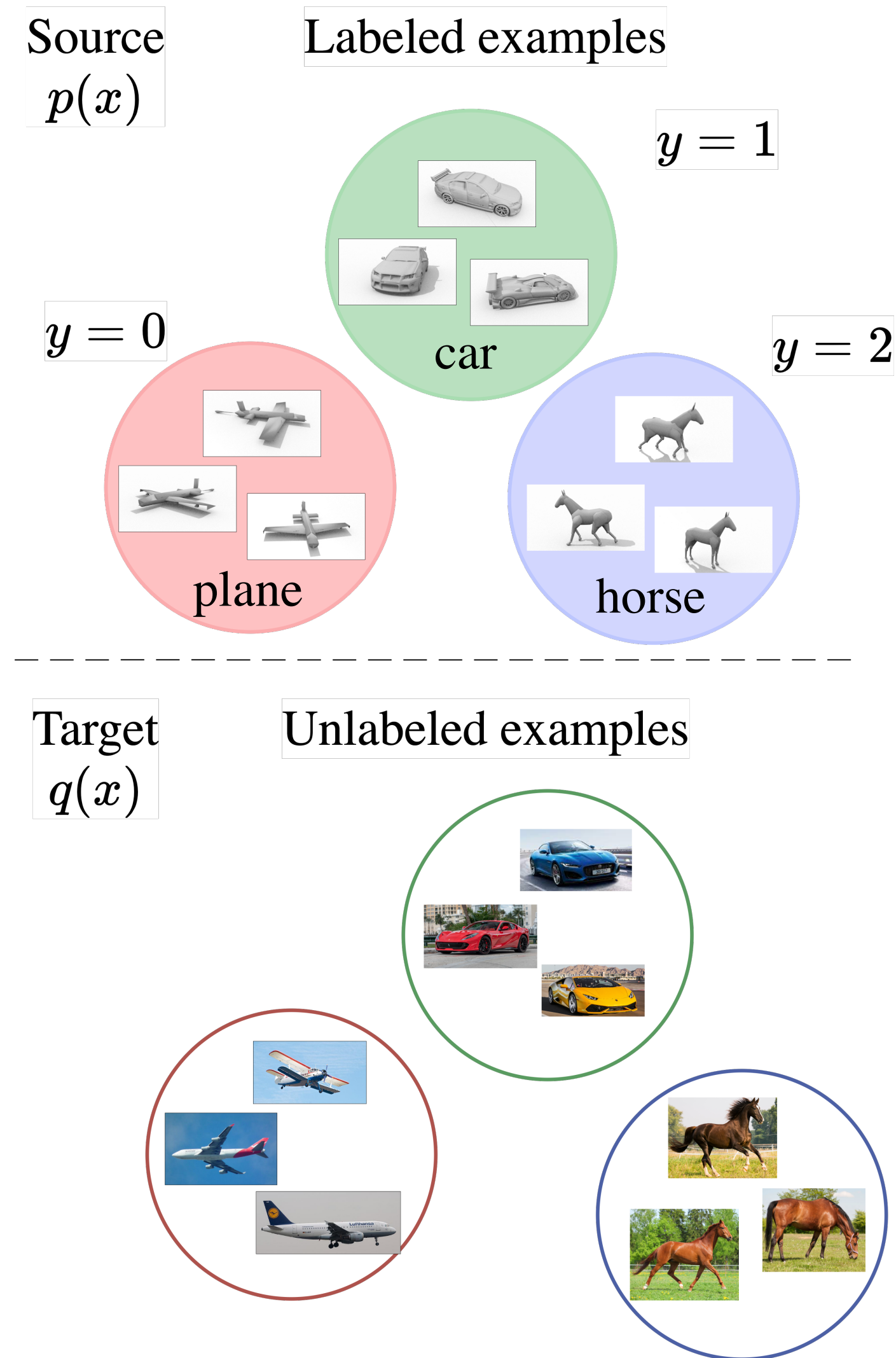
Background: distribution alignment

Given $\mathcal{P} = \{p^\theta \mid \theta \in \Theta\}$ $\mathcal{Q} = \{q^\theta \mid \theta \in \Theta\}$

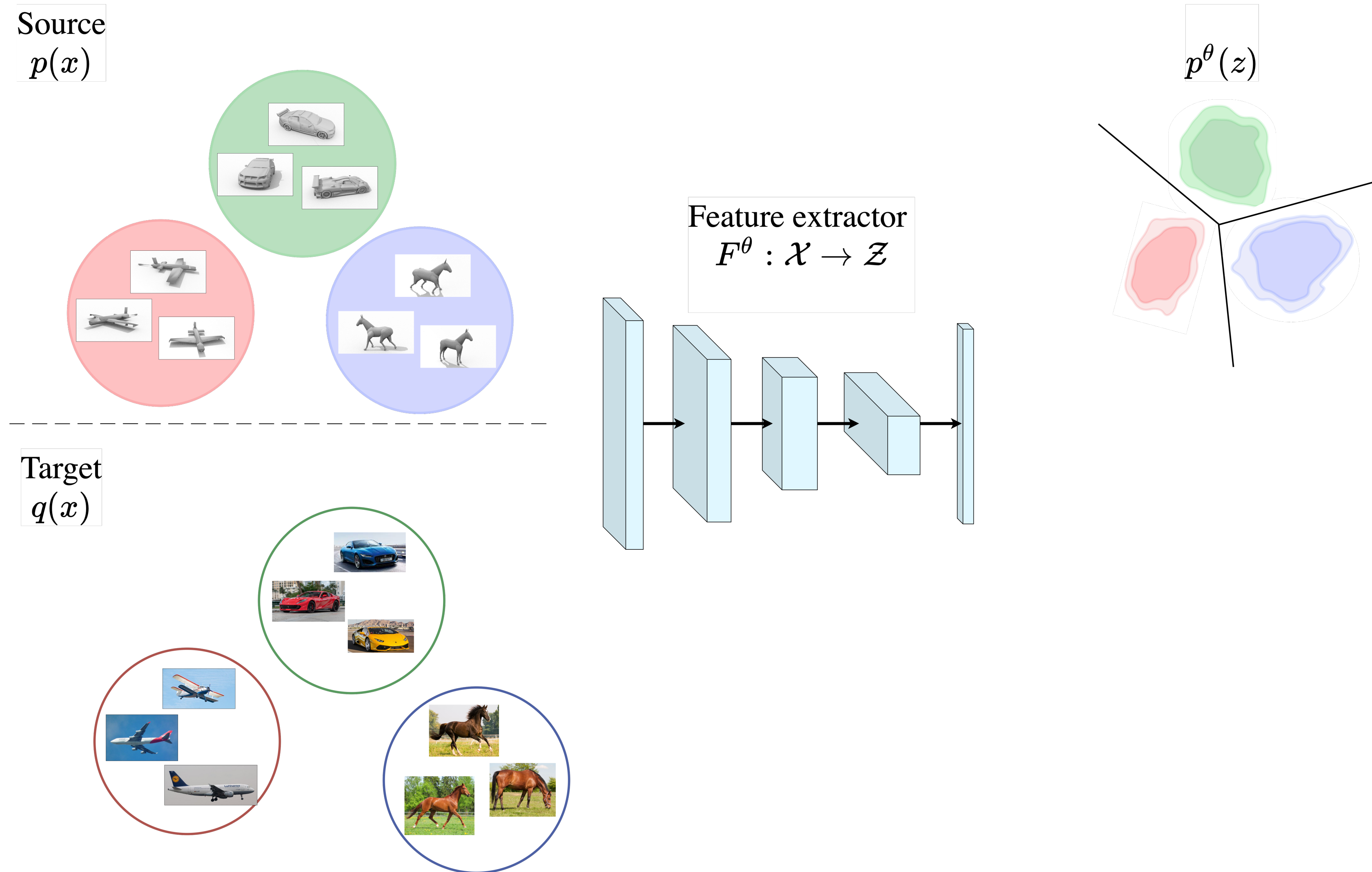
Find $\theta^* : p^{\theta^*} = q^{\theta^*}$

- **Generative Models** (GAN, Goodfellow et al. 2014)
 - alignment of generated and data distributions
- **Domain Adaptation** (DANN, Ganin et al. 2016)
 - alignment of representations across domains

Background: distribution alignment in domain adaptation

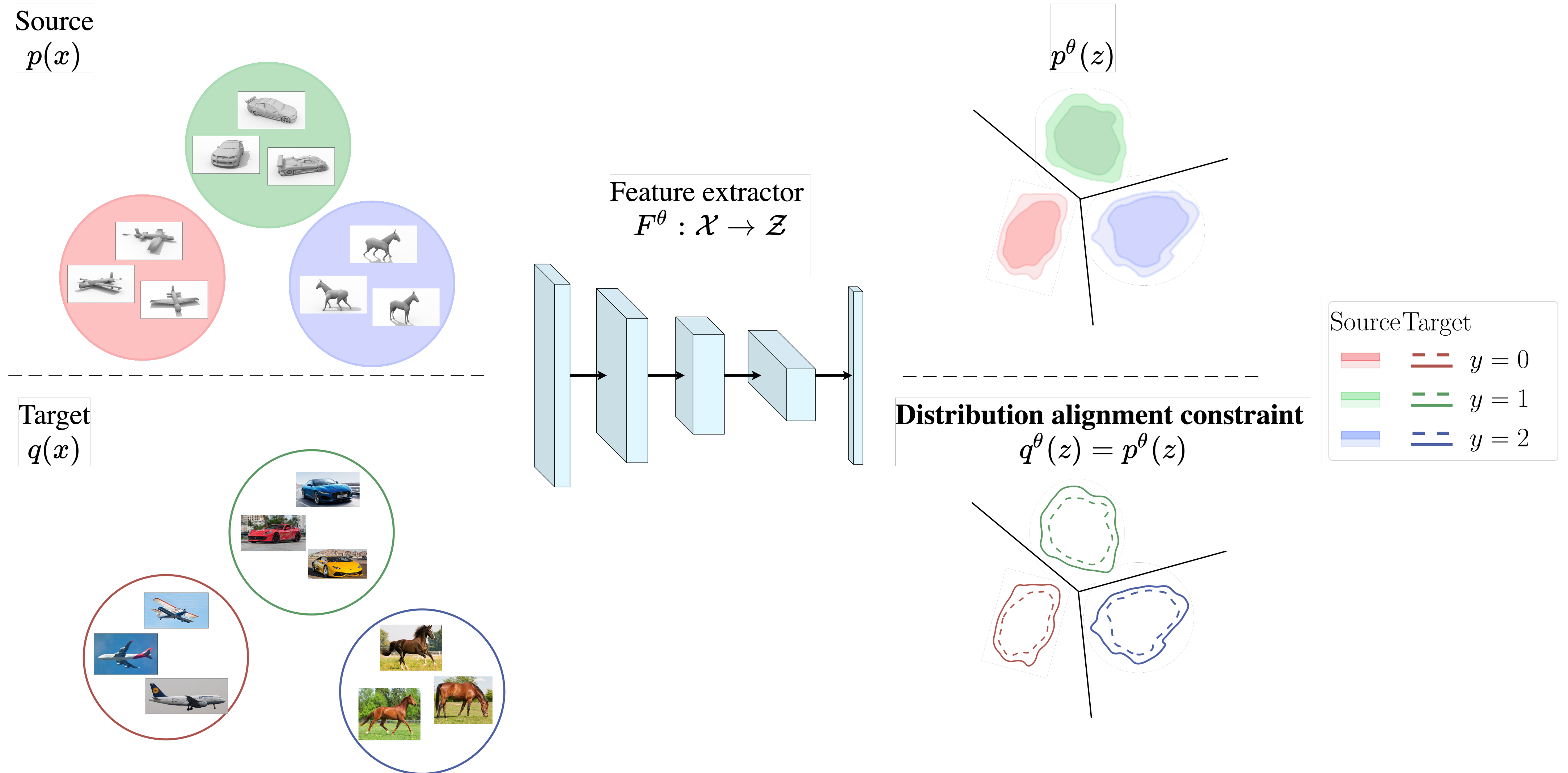


Background: distribution alignment in domain adaptation



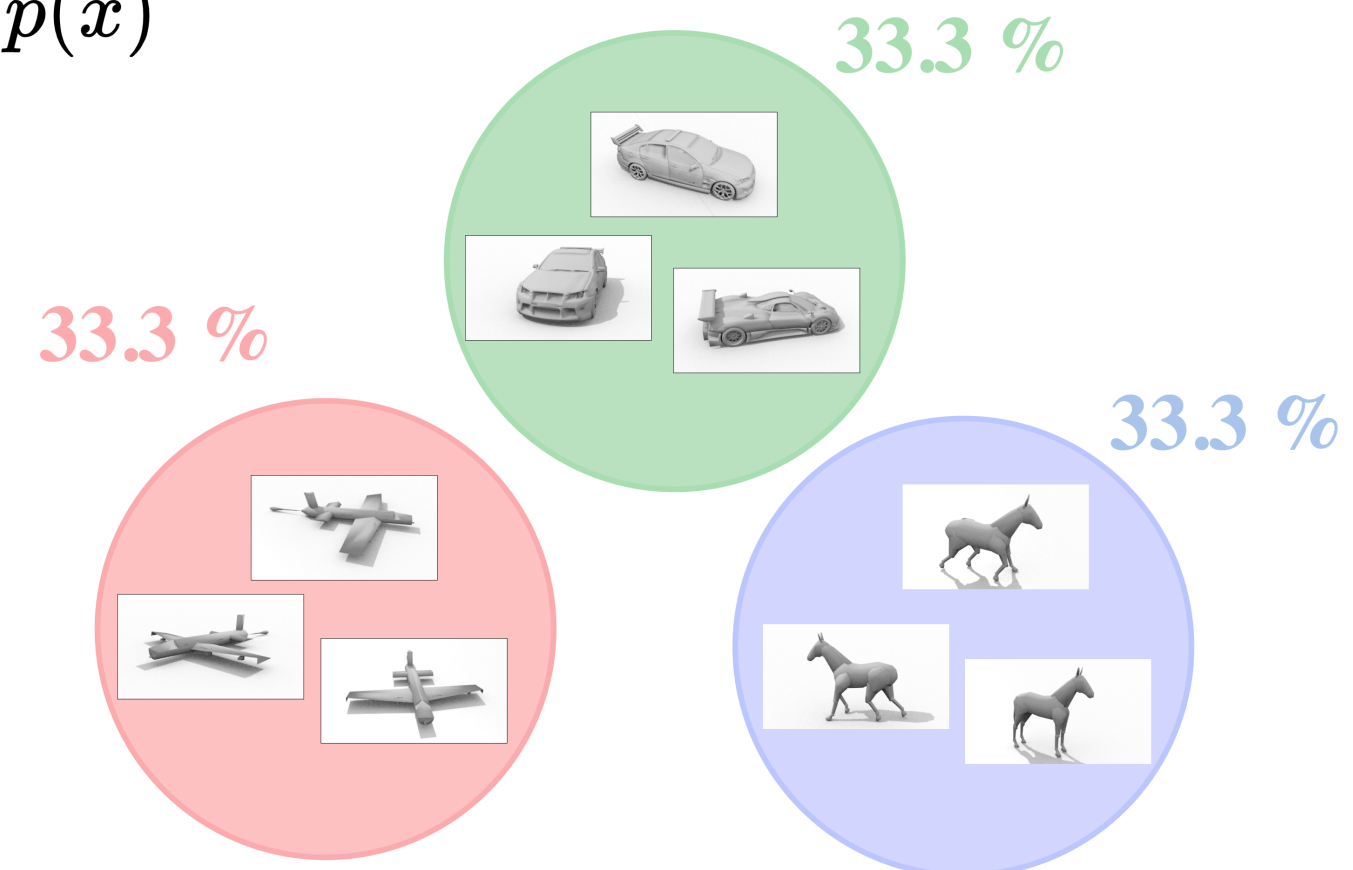
DANN: domain adversarial neural networks [Ganin et al., 2016]

Background: distribution alignment in domain adaptation



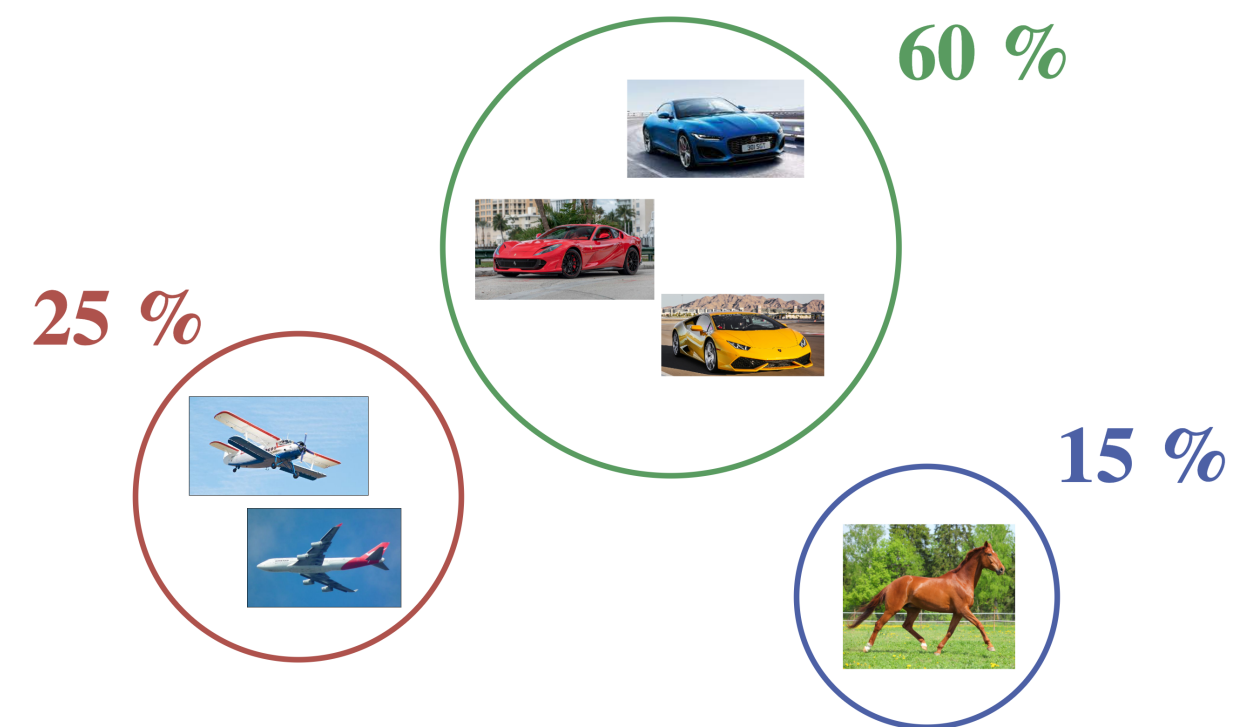
Problem: label distribution shift

Source
 $p(x)$

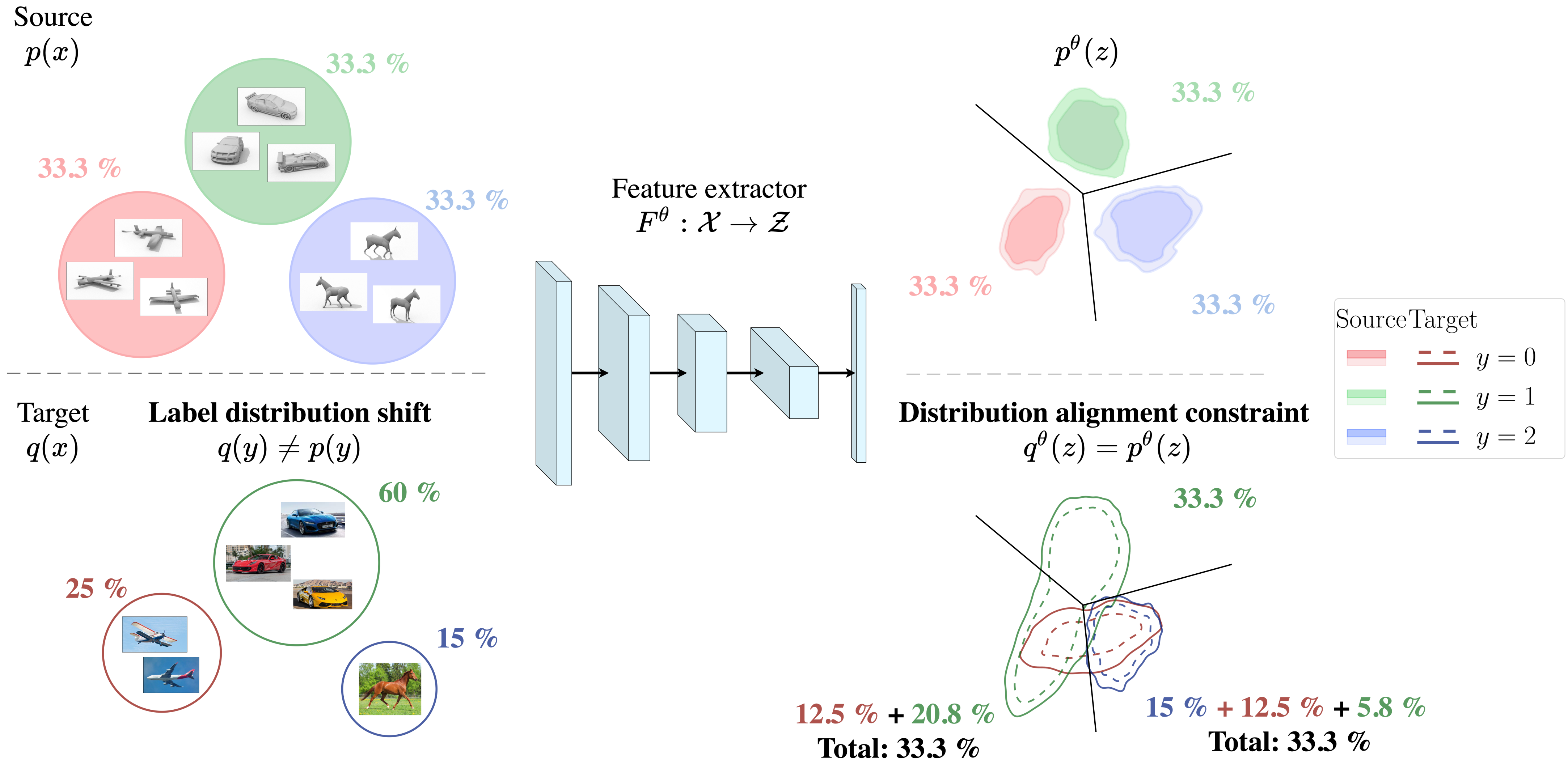


Target
 $q(x)$

Label distribution shift
 $q(y) \neq p(y)$



Problem: label distribution shift

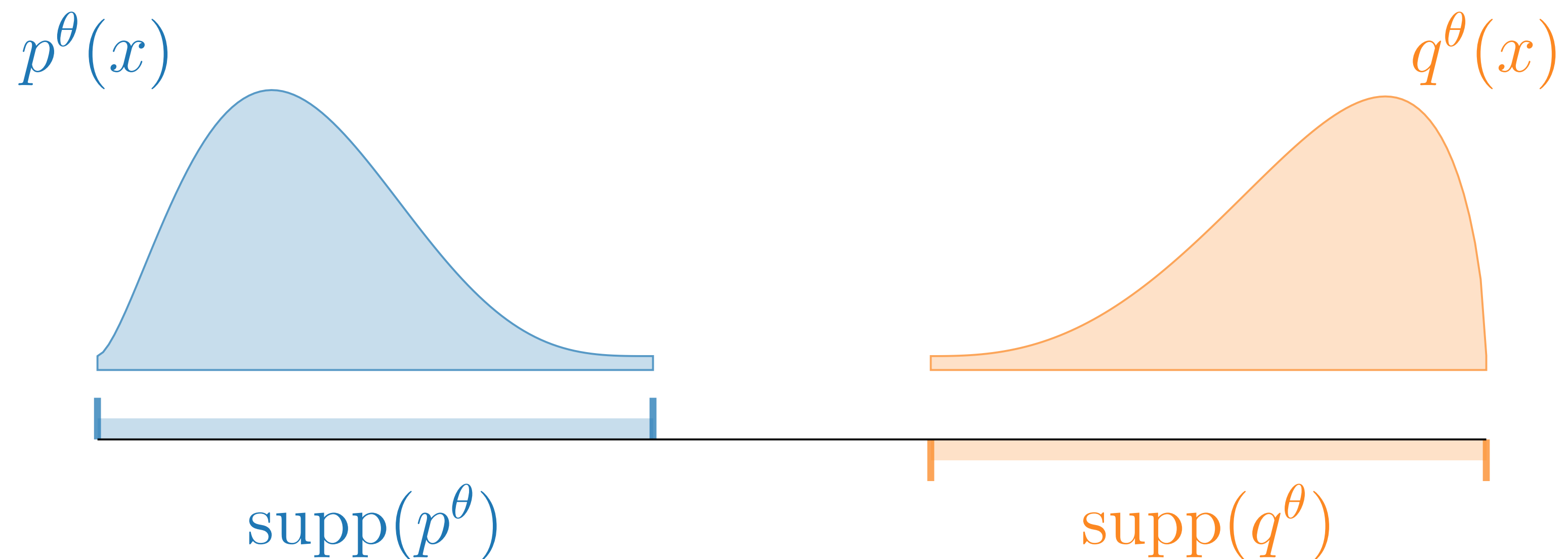


[Zhao et al., 2019; Wu et al., 2019; Tachet des Combes et al, 2020, ...]

Support alignment

We study the problem of aligning the supports of distributions

Given $\mathcal{P} = \{p^\theta \mid \theta \in \Theta\}$ $\mathcal{Q} = \{q^\theta \mid \theta \in \Theta\}$

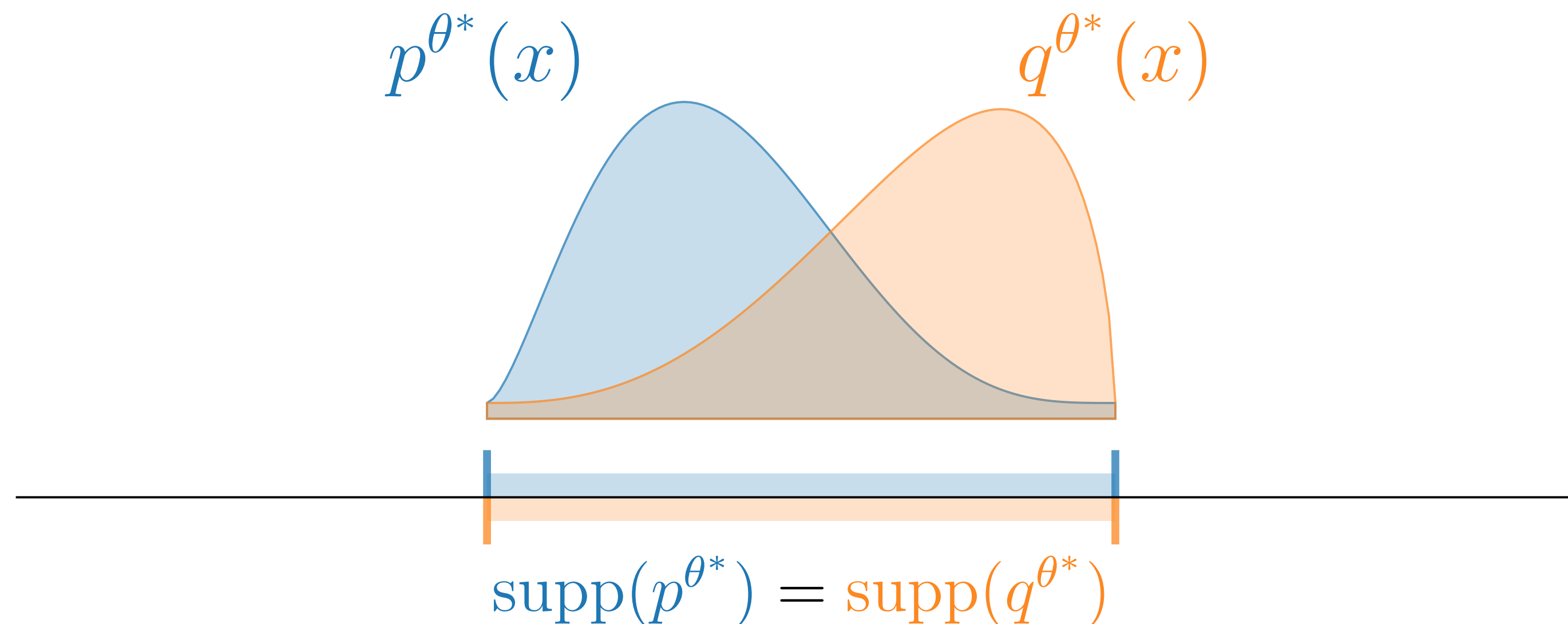


Support alignment

We study the problem of aligning the supports of distributions

Given $\mathcal{P} = \{p^\theta \mid \theta \in \Theta\}$ $\mathcal{Q} = \{q^\theta \mid \theta \in \Theta\}$

Find $\theta^* : \text{supp}(p^{\theta^*}) = \text{supp}(q^{\theta^*})$



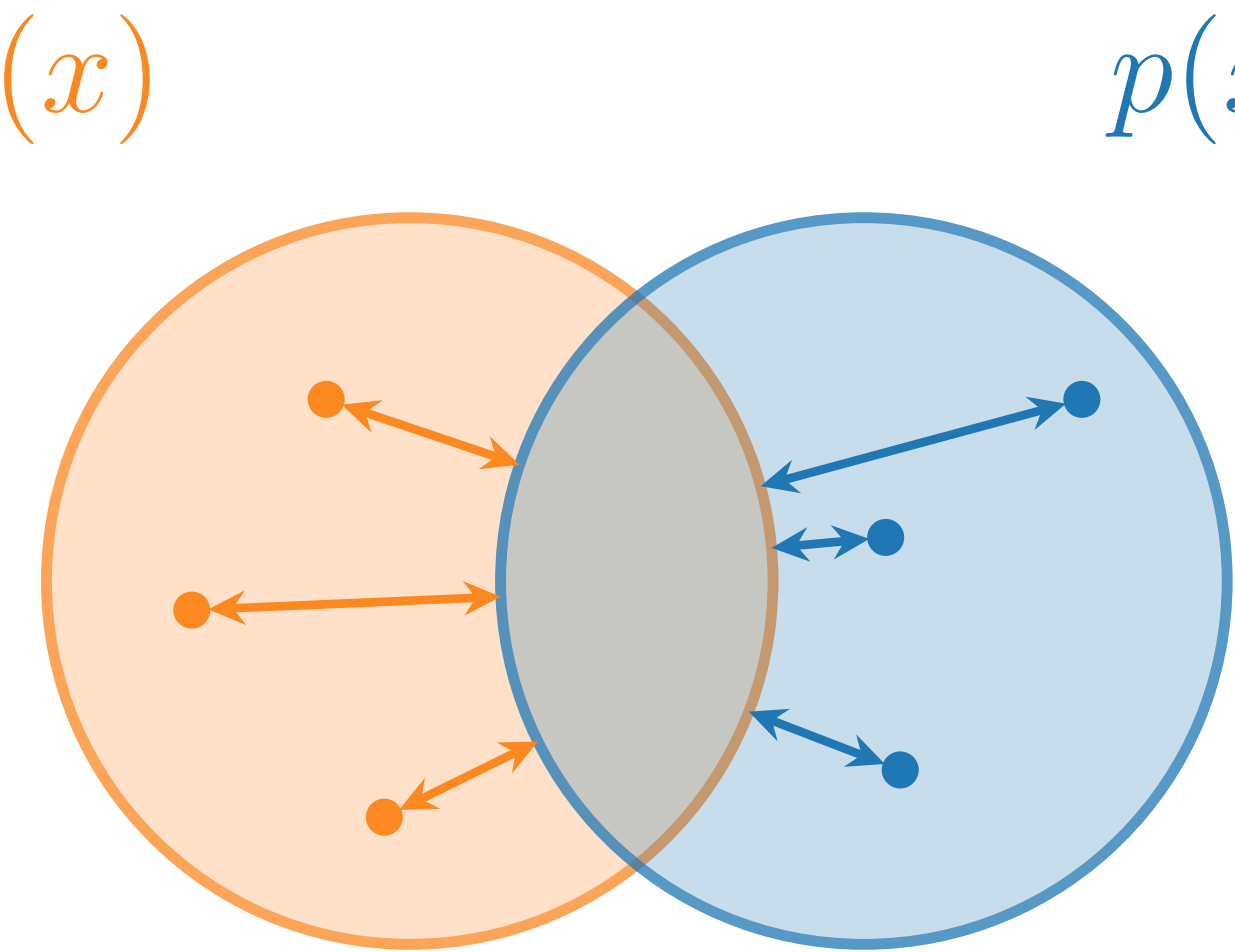
Support divergence

Symmetric Support Difference (SSD) divergence

$$\mathcal{D}_{\Delta}(p, q) = \mathbb{E}_{x^q \sim q} [d(x^q, \text{supp}(p))] + \mathbb{E}_{x^p \sim p} [d(x^p, \text{supp}(q))] \quad q(x) \quad p(x)$$

$$d(x^q, \text{supp}(p)) = \inf_{x^p \in \text{supp}(p)} d(x^q, x^p)$$

$$d(x^p, \text{supp}(q)) = \inf_{x^q \in \text{supp}(q)} d(x^p, x^q)$$



- 1) $\mathcal{D}_{\Delta}(p, q) \geq 0 \quad \forall p, q;$
- 2) $\mathcal{D}_{\Delta}(p, q) = 0 \iff \text{supp}(p) = \text{supp}(q)$

Support alignment via log-loss discriminator

$$\sup_{f: \mathcal{X} \rightarrow [0,1]} \mathbb{E}_{x \sim p} [\log f(x)] + \mathbb{E}_{y \sim q} [\log(1 - f(y))]$$

$$f^*(x) = \frac{p(x)}{p(x) + q(x)}$$

Support alignment via log-loss discriminator

$$\sup_{f: \mathcal{X} \rightarrow [0,1]} \mathbb{E}_{x \sim p} [\log f(x)] + \mathbb{E}_{y \sim q} [\log(1 - f(y))]$$

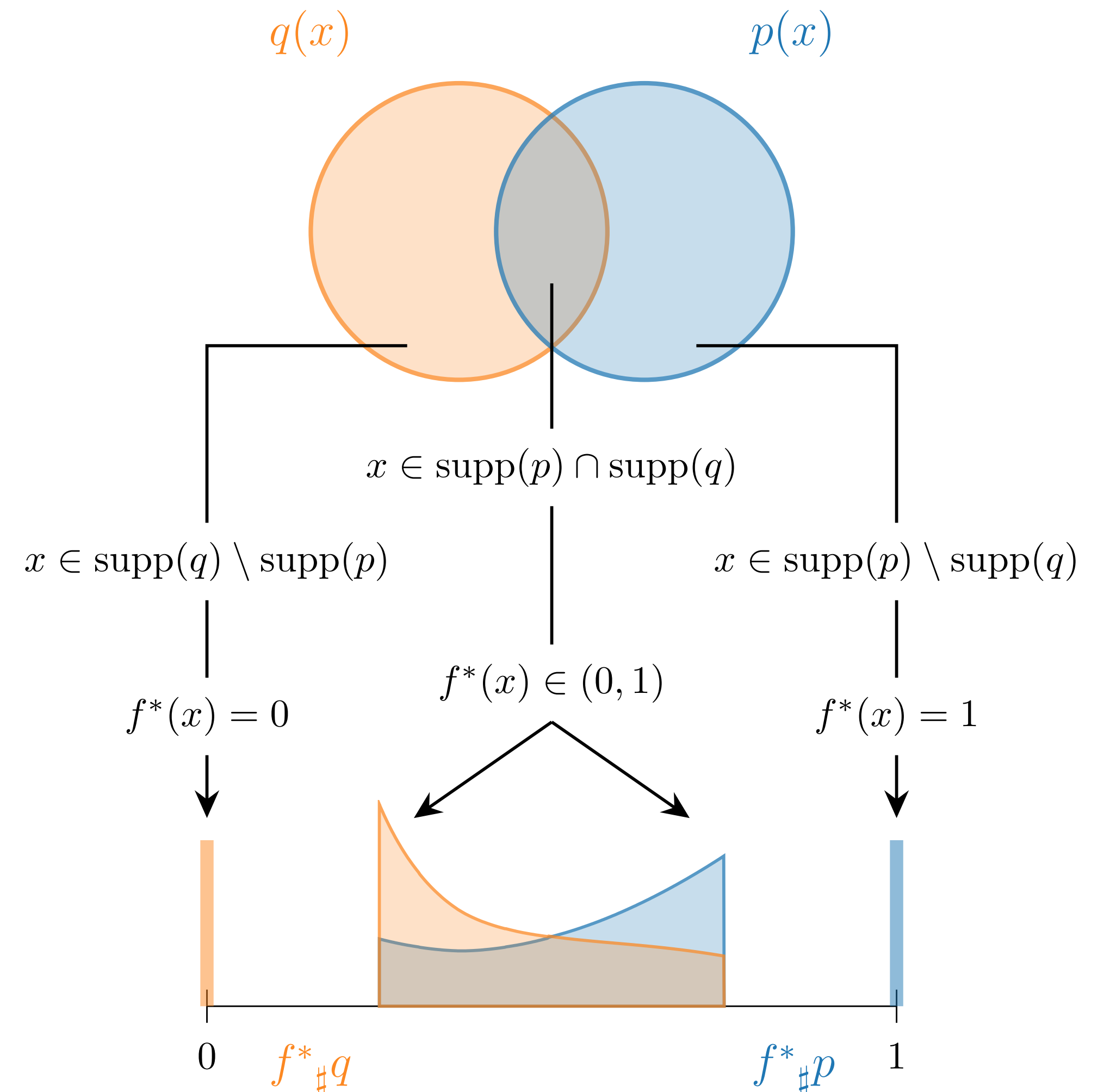
$$f^*(x) = \frac{p(x)}{p(x) + q(x)}$$

Theorem

The mapping $f^* : \mathcal{X} \rightarrow [0, 1]$
 realized by the optimal discriminator
 preserves support discrepancy

$$\mathcal{D}_\Delta(p, q) = 0 \iff \mathcal{D}_\Delta(f^*_{\#}p, f^*_{\#}q) = 0$$

$f^*_{\#}p, f^*_{\#}q$ — pushforward distributions



Support alignment via log-loss discriminator

$$\sup_{f: \mathcal{X} \rightarrow [0,1]} \mathbb{E}_{x \sim p} [\log f(x)] + \mathbb{E}_{y \sim q} [\log(1 - f(y))]$$

$$f^*(x) = \frac{p(x)}{p(x) + q(x)}$$

Theorem

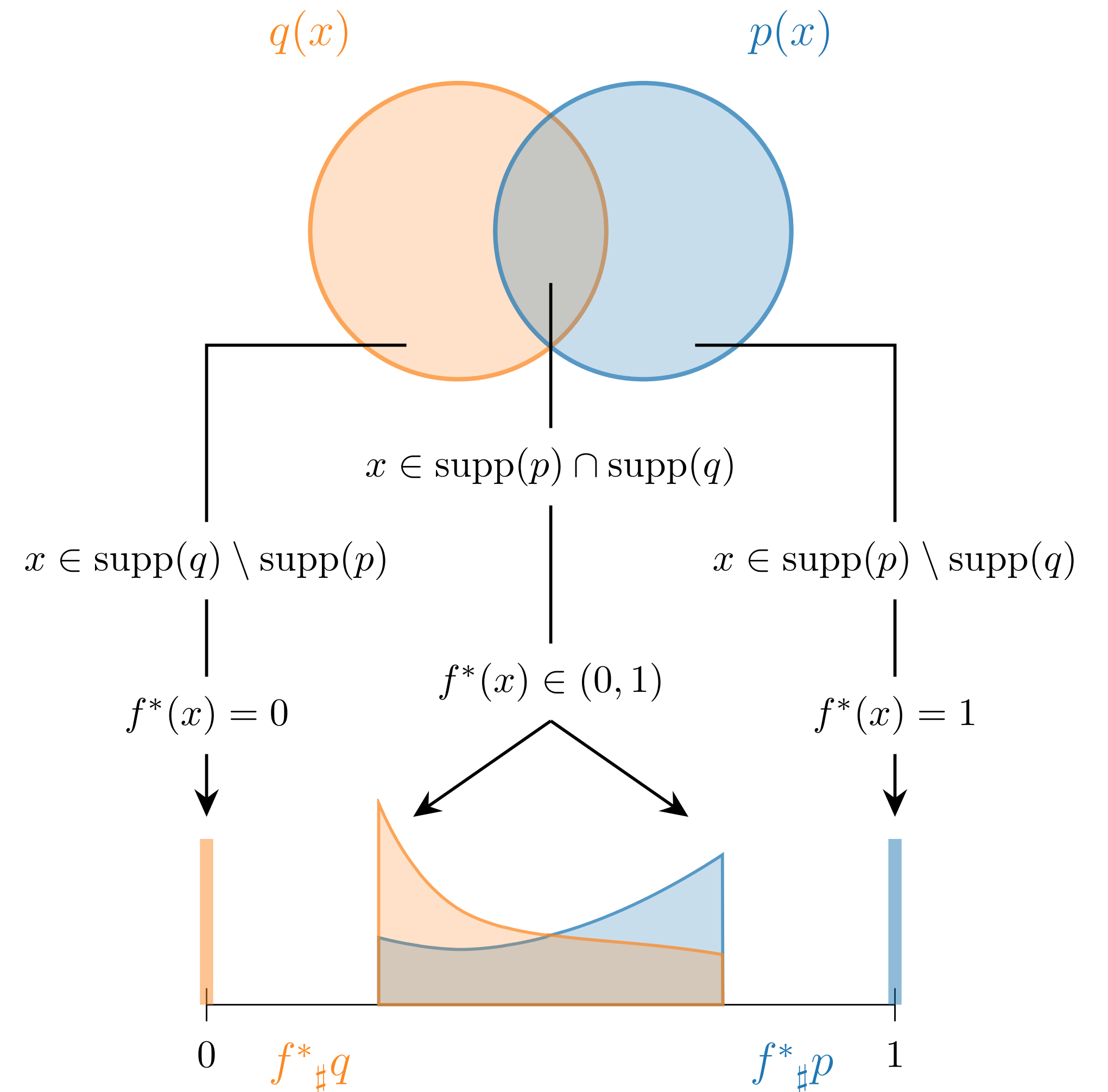
The mapping $f^* : \mathcal{X} \rightarrow [0, 1]$
realized by the optimal discriminator
preserves support discrepancy

$$\mathcal{D}_{\Delta}(p, q) = 0 \iff \mathcal{D}_{\Delta}(f^*_{\#}p, f^*_{\#}q) = 0$$

$f^*_{\#}p, f^*_{\#}q$ — pushforward distributions

Remark: the result holds for $g : \mathcal{X} \rightarrow \mathbb{R}$

$$g(x) : f(x) = \text{sigmoid}(g(x))$$



Adversarial Distribution Alignment (GAN/DANN)

Goal: find $\theta : p^\theta = q^\theta$

Adversarial Support Alignment (ASA)

Goal: find $\theta : \text{supp}(p^\theta) = \text{supp}(q^\theta)$

Discriminator objective: $\min_{g: \mathcal{X} \rightarrow \mathbb{R}} \mathcal{L}_D(\theta, g)$

$$\mathcal{L}_D(\theta, g) = \mathbb{E}_{x \sim p^\theta} \left[\log \left(1 + e^{-g(x)} \right) \right] + \mathbb{E}_{x \sim q^\theta} \left[\log \left(1 + e^{g(x)} \right) \right]$$

Alignment objective: $\min_{\theta} \mathcal{L}_A(\theta, g)$

$$\mathcal{L}_A(\theta, g) = \mathbb{E}_{x \sim p^\theta} \left[\log \left(1 + e^{g(x)} \right) \right] + \mathbb{E}_{x \sim q^\theta} \left[\log \left(1 + e^{-g(x)} \right) \right]$$

$$\mathcal{L}_A(\theta, g) = \mathbb{E}_{x \sim p^\theta} [d(g(x), \text{supp}(g_\# q))] + \mathbb{E}_{x \sim q^\theta} [d(g(x), \text{supp}(g_\# p))]$$

Adversarial Distribution Alignment (GAN/DANN)

Goal: find $\theta : p^\theta = q^\theta$

Adversarial Support Alignment (ASA)

Goal: find $\theta : \text{supp}(p^\theta) = \text{supp}(q^\theta)$

1. Sample examples

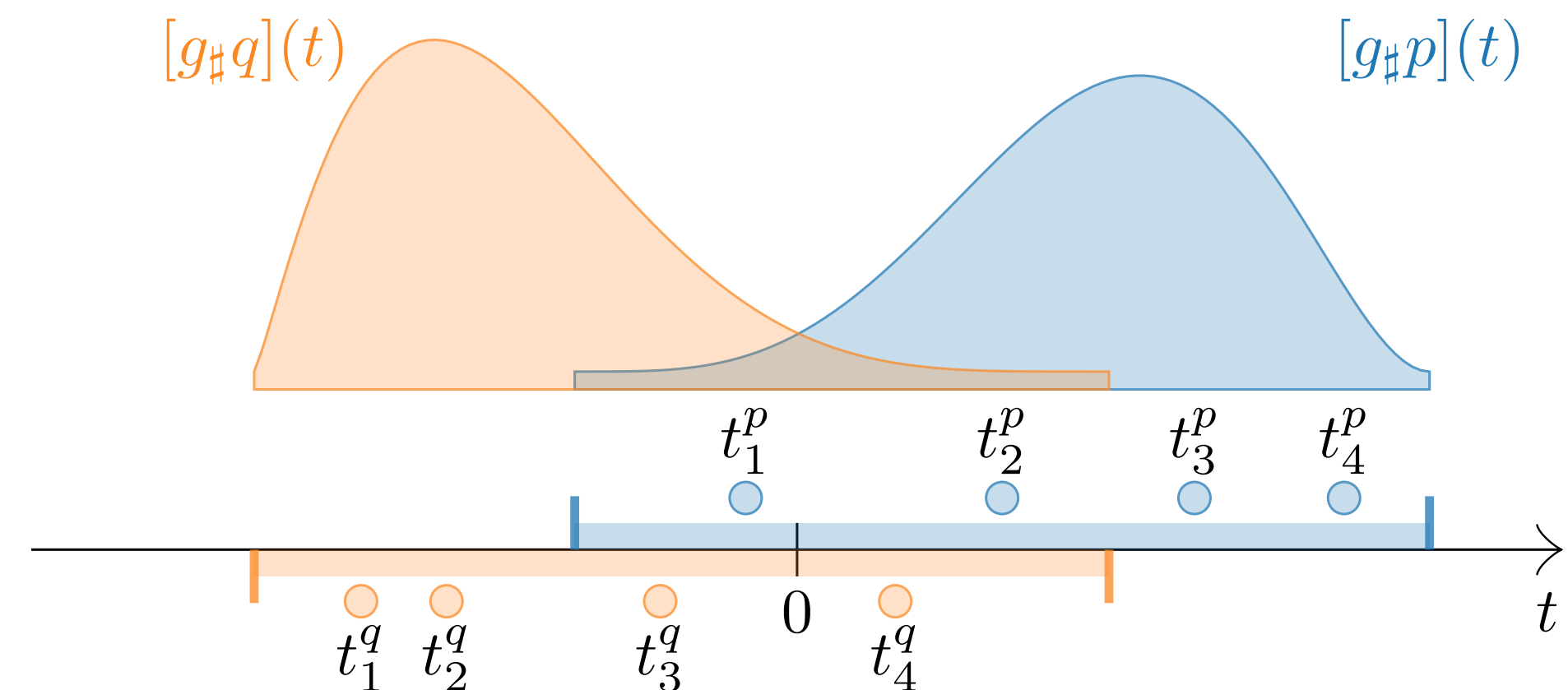
$$x_1^p, x_2^p, \dots, x_N^p \sim p^\theta(x)$$

$$x_1^q, x_2^q, \dots, x_N^q \sim q^\theta(x)$$

2. Apply discriminator mapping

$$p^\theta(x) \xrightarrow{g} [g_\# p^\theta](t) : t_i^p = g(x_i^p)$$

$$q^\theta(x) \xrightarrow{g} [g_\# q^\theta](t) : t_i^q = g(x_i^q)$$



Adversarial Distribution Alignment (GAN/DANN)

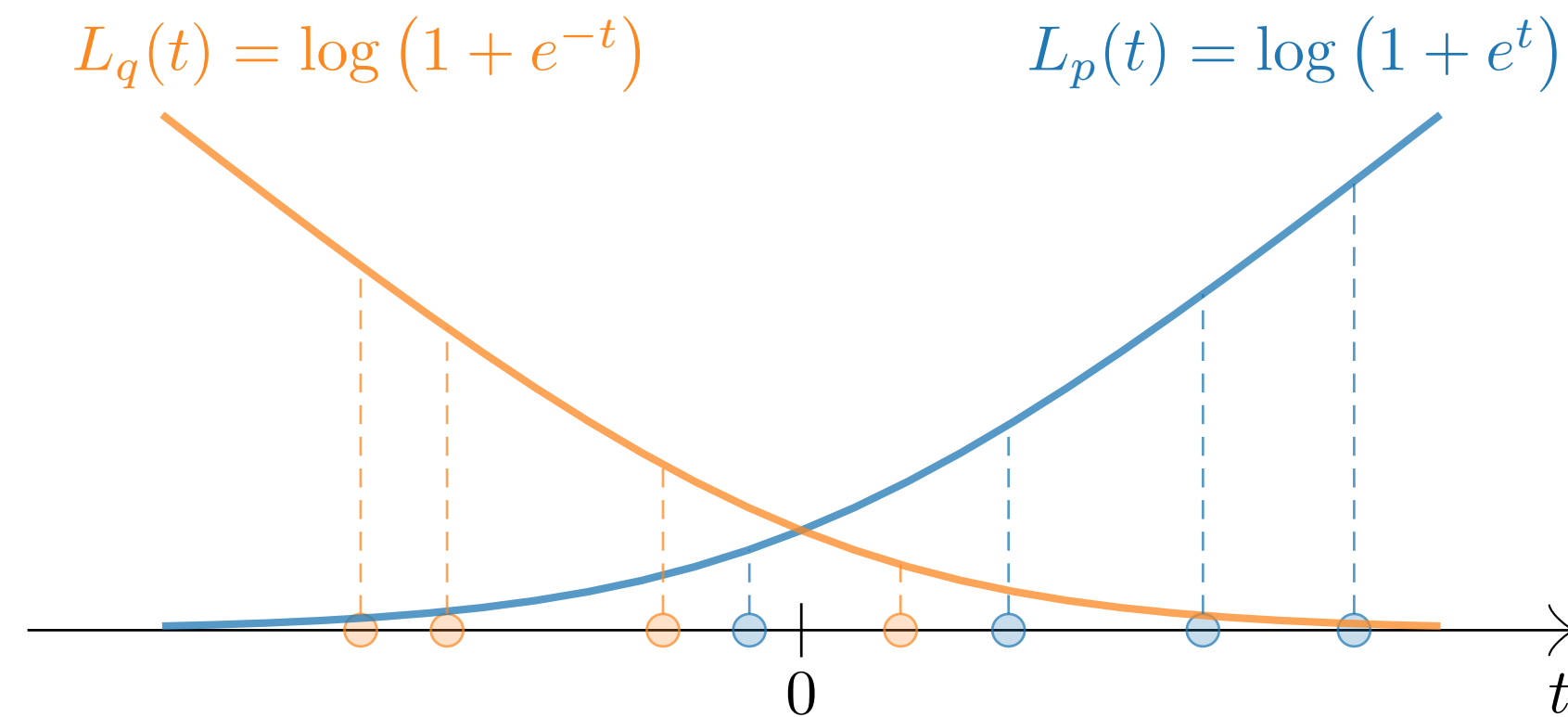
Goal: find $\theta : p^\theta = q^\theta$

Adversarial Support Alignment (ASA)

Goal: find $\theta : \text{supp}(p^\theta) = \text{supp}(q^\theta)$

3. Compute loss function

$$\mathcal{L}_A = \frac{1}{N} \sum_i \underbrace{\log(1 + e^{-t_i^q})}_{L_q(t_i^q)} + \frac{1}{N} \sum_i \underbrace{\log(1 + e^{t_i^p})}_{L_p(t_i^p)}$$



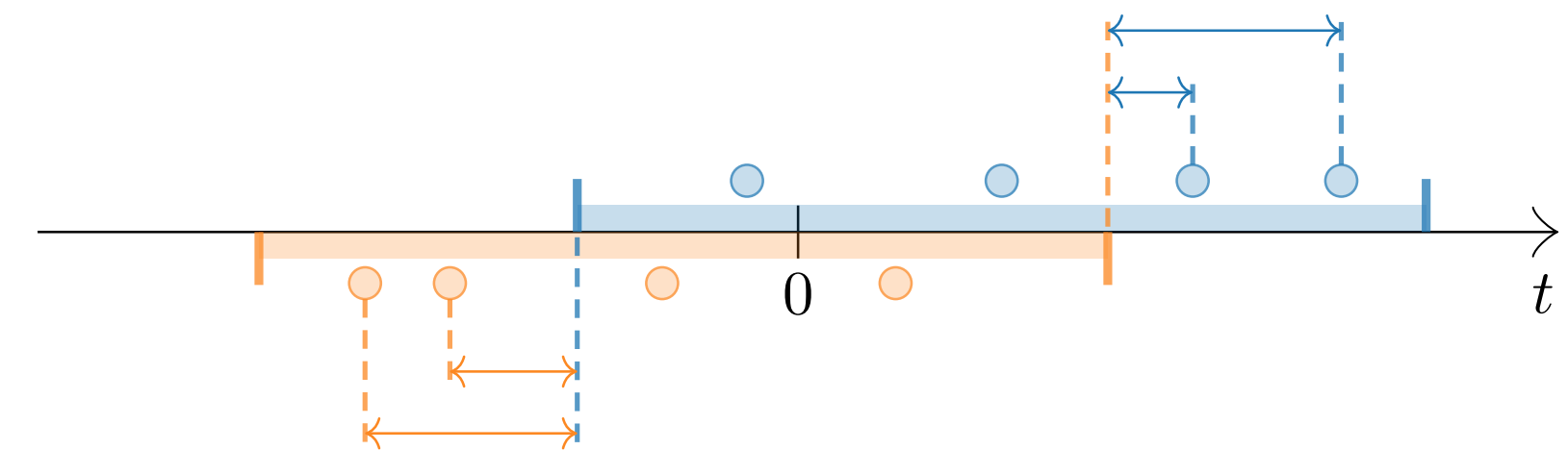
$$\mathcal{L}_A = \frac{1}{N} \sum_i \underbrace{d(t_i^q, \text{supp}(g_\#p))}_{L_q(t_i^q)} + \frac{1}{N} \sum_i \underbrace{d(t_i^p, \text{supp}(g_\#q))}_{L_p(t_i^p)}$$

$$L_q(t) = d(t, \text{supp}(g_\#p)) = |t - \pi_p^*(t)|$$

$$\pi_p^*(t) = \underset{t^p \in \text{supp}(g_\#p)}{\text{argmin}} |t - t^p|$$

$$L_p(t) = d(t, \text{supp}(g_\#q)) = |t - \pi_q^*(t)|$$

$$\pi_q^*(t) = \underset{t^q \in \text{supp}(g_\#q)}{\text{argmin}} |t - t^q|$$



Adversarial Distribution Alignment (GAN/DANN)

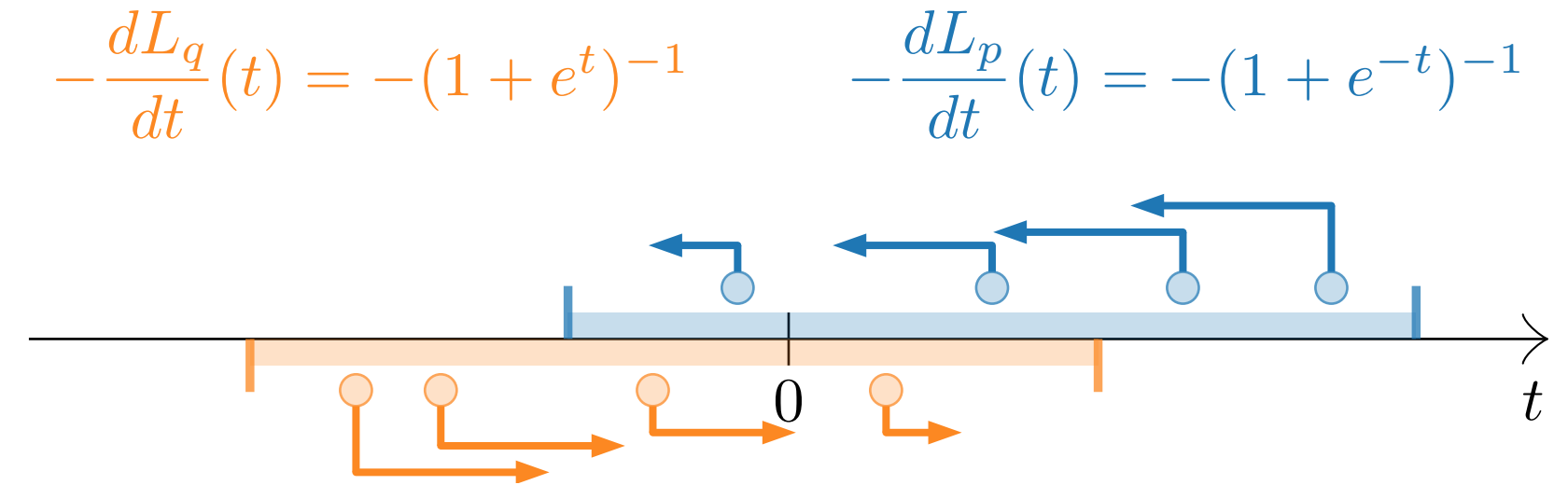
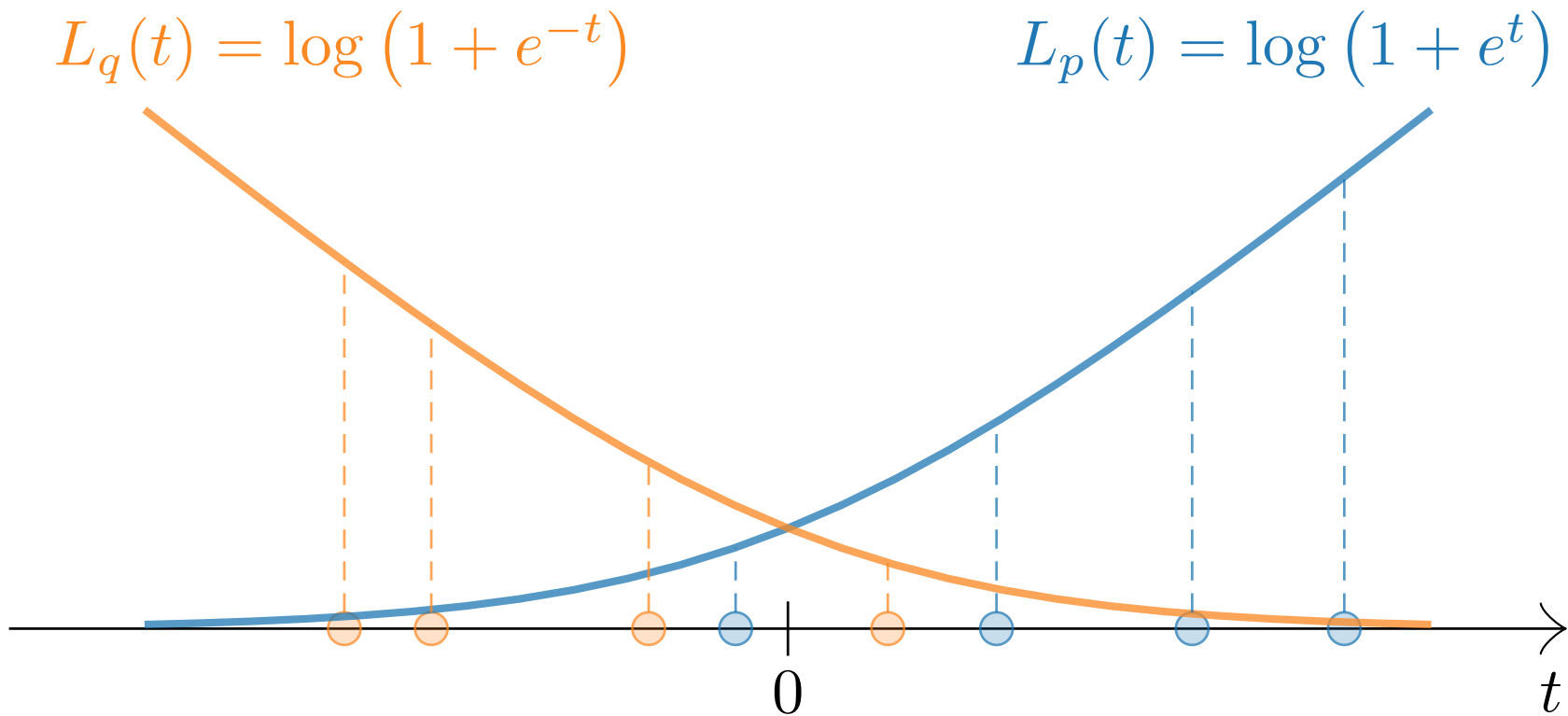
Goal: find $\theta : p^\theta = q^\theta$

Adversarial Support Alignment (ASA)

Goal: find $\theta : \text{supp}(p^\theta) = \text{supp}(q^\theta)$

4. Compute and propagate gradients

$$\mathcal{L}_A = \frac{1}{N} \sum_i \underbrace{\log(1 + e^{-t_i^q})}_{L_q(t_i^q)} + \frac{1}{N} \sum_i \underbrace{\log(1 + e^{t_i^p})}_{L_p(t_i^p)}$$



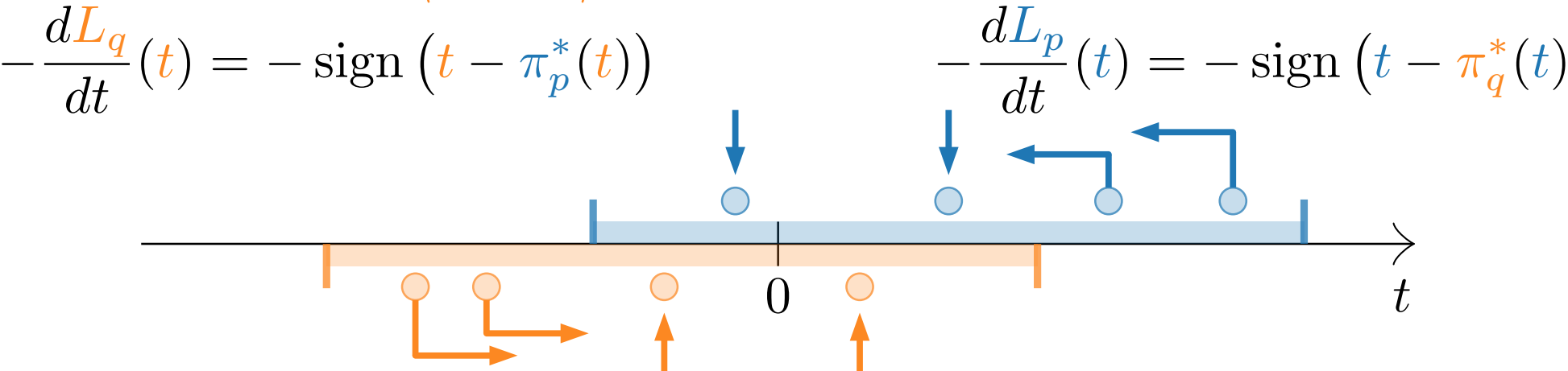
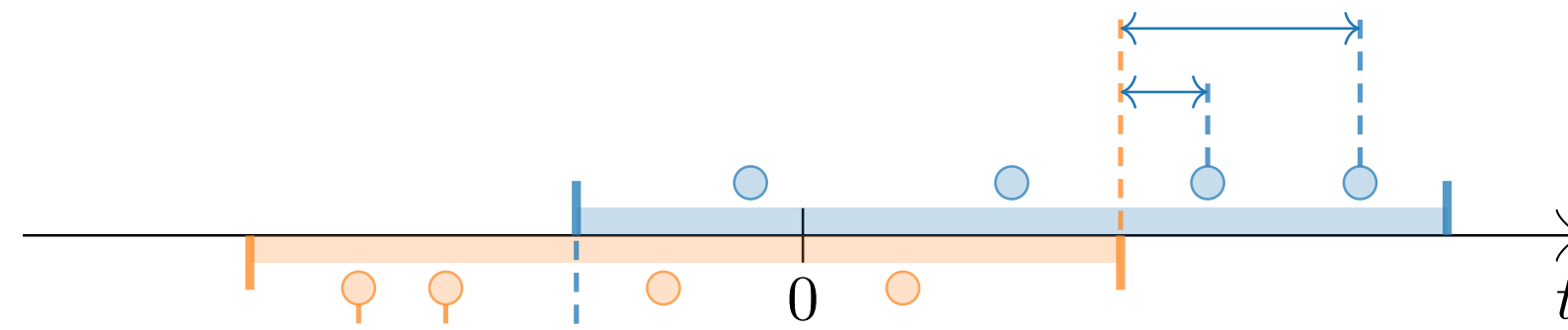
$$\mathcal{L}_A = \frac{1}{N} \sum_i \underbrace{d(t_i^q, \text{supp}(g_{\#}p))}_{L_q(t_i^q)} + \frac{1}{N} \sum_i \underbrace{d(t_i^p, \text{supp}(g_{\#}q))}_{L_p(t_i^p)}$$

$$L_q(t) = d(t, \text{supp}(g_{\#}p)) = |t - \pi_p^*(t)|$$

$$L_p(t) = d(t, \text{supp}(g_{\#}q)) = |t - \pi_q^*(t)|$$

$$\pi_p^*(t) = \underset{t^p \in \text{supp}(g_{\#}p)}{\text{argmin}} |t - t^p|$$

$$\pi_q^*(t) = \underset{t^q \in \text{supp}(g_{\#}q)}{\text{argmin}} |t - t^q|$$



Domain adaptation results: 2D embeddings visualization

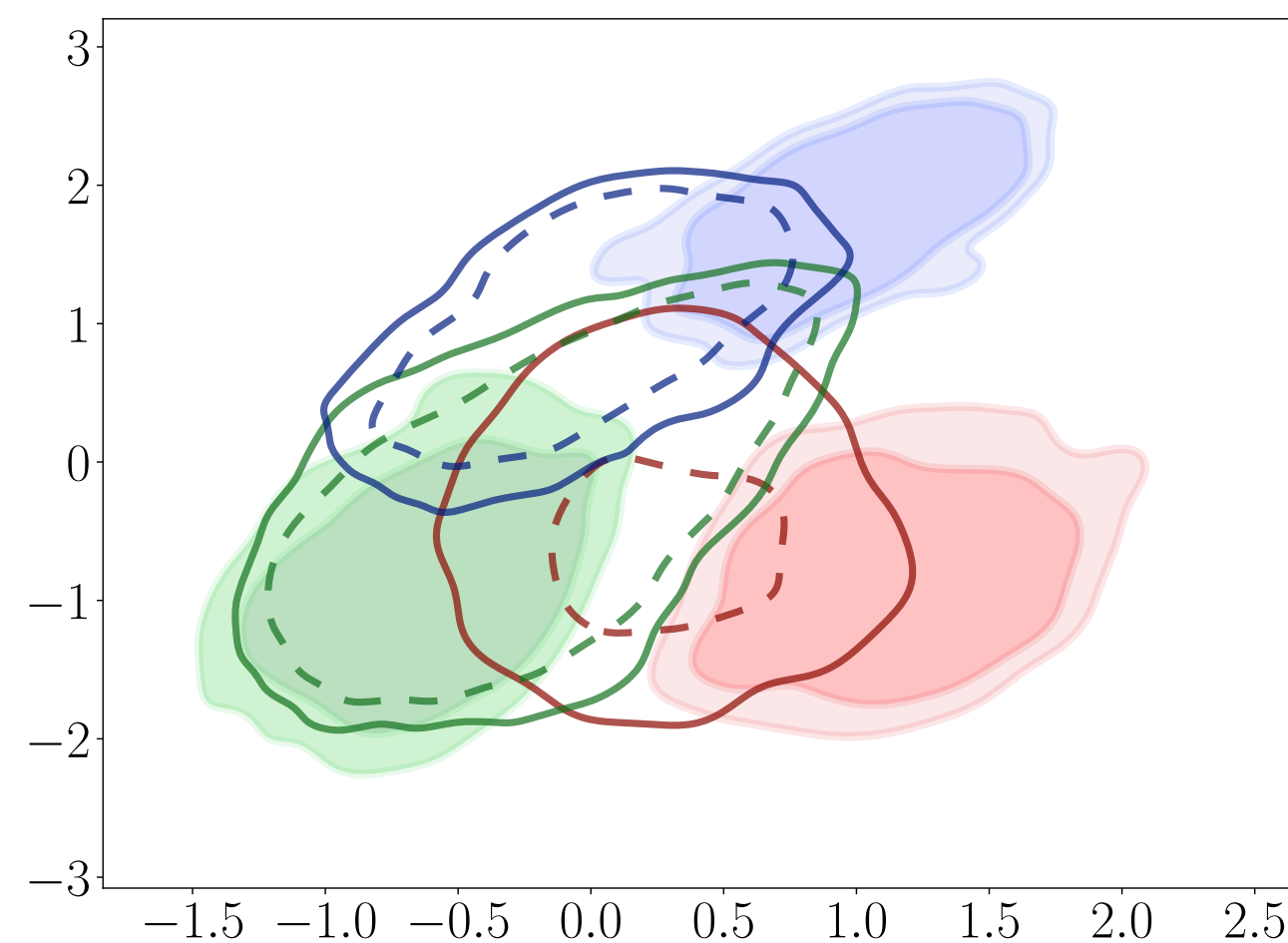
USPS → MNIST, 3 classes

Source class distribution

[33%, 33%, 33%]

Target class distribution

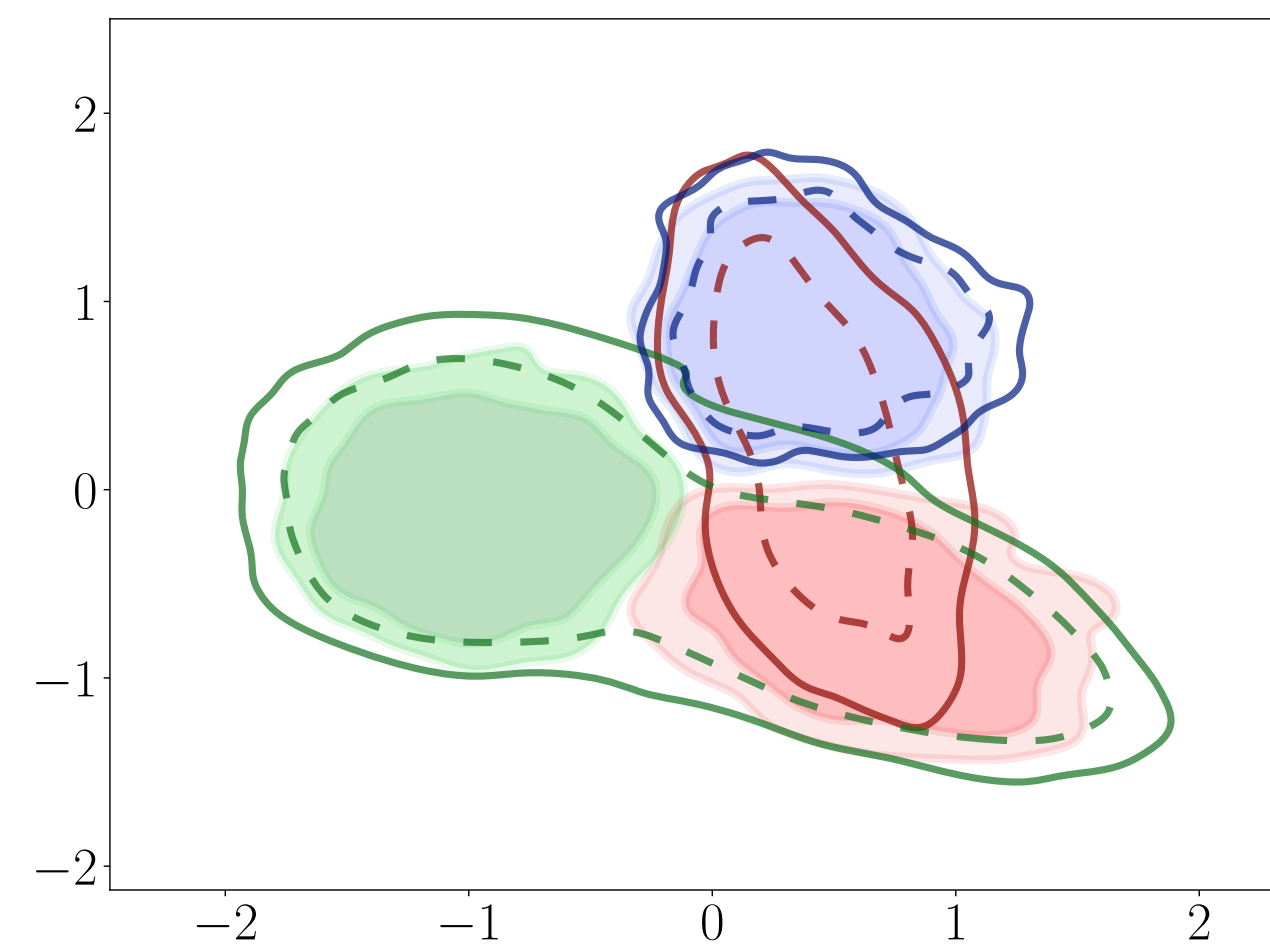
[23%, 65%, 12%]



(a) No DA (avg acc: 63%)

$$\mathcal{D}_W(p_Z^\theta, q_Z^\theta) = 0.78$$

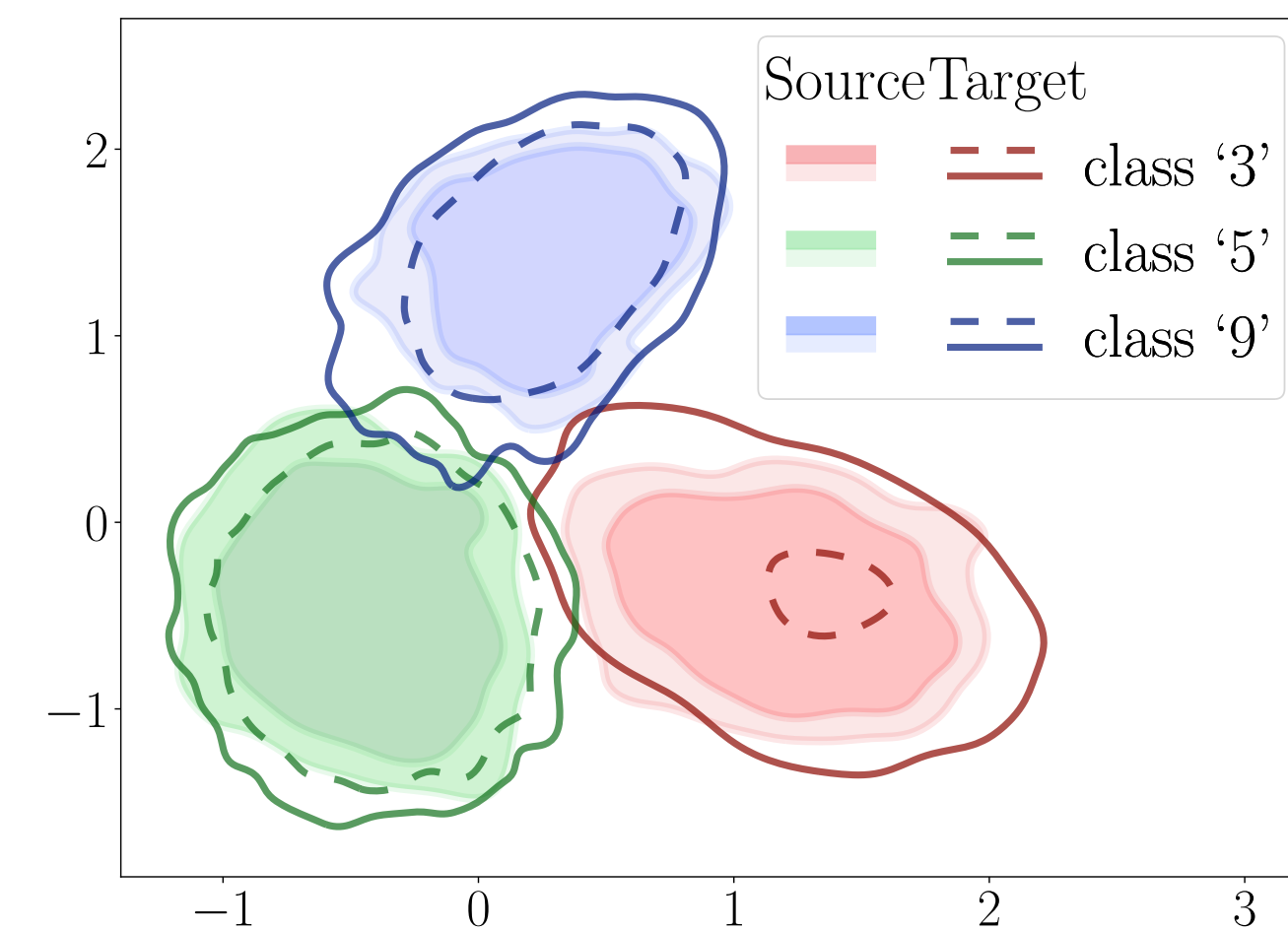
$$\mathcal{D}_\Delta(p_Z^\theta, q_Z^\theta) = 0.10$$



(b) DANN (avg acc: 75%)

$$\mathcal{D}_W(p_Z^\theta, q_Z^\theta) = 0.07$$

$$\mathcal{D}_\Delta(p_Z^\theta, q_Z^\theta) = 0.02$$



(c) ASA-abs (avg acc: 94%)

$$\mathcal{D}_W(p_Z^\theta, q_Z^\theta) = 0.59$$

$$\mathcal{D}_\Delta(p_Z^\theta, q_Z^\theta) = 0.03$$

Domain adaptation results: USPS→MNIST, LeNet

Average and minimum class accuracy (%) on USPS→MNIST across different levels of shifts in label distributions (α).

Algorithm	$\alpha = 0.0$ no shift		$\alpha = 1.0$		$\alpha = 1.5$		$\alpha = 2.0$ severe shift	
	average	min	average	min	average	min	average	min
No DA	71.9	20.3	72.9	25.8	71.3	27.5	71.3	16.6
DANN	97.8	96.0	83.5	25.1	70.0	01.1	57.8	00.9
VADA	98.0	96.2	88.2	48.9	78.2	06.6	61.9	01.4

Distribution Alignment

Domain adaptation results: USPS→MNIST, LeNet

Average and minimum class accuracy (%) on USPS→MNIST across different levels of shifts in label distributions (α).

		$\alpha = 0.0$ no shift		$\alpha = 1.0$		$\alpha = 1.5$		$\alpha = 2.0$ severe shift	
Algorithm		average	min	average	min	average	min	average	min
No DA		71.9	20.3	72.9	25.8	71.3	27.5	71.3	16.6
Distribution Alignment	DANN	97.8	96.0	83.5	25.1	70.0	01.1	57.8	00.9
	VADA	98.0	96.2	88.2	48.9	78.2	06.6	61.9	01.4
Relaxed Distribution Alignment	IWDAN	97.5	95.7	95.7	81.3	86.5	15.2	74.4	07.3
	IWCDAN	98.0	96.6	96.7	85.1	91.3	66.5	77.5	22.2
	sDANN-4	87.4	05.6	94.9	85.7	86.8	21.6	81.5	39.3

Domain adaptation results: USPS→MNIST, LeNet

Average and minimum class accuracy (%) on USPS→MNIST across different levels of shifts in label distributions (α).

		$\alpha = 0.0$ no shift		$\alpha = 1.0$		$\alpha = 1.5$		$\alpha = 2.0$ severe shift	
Algorithm		average	min	average	min	average	min	average	min
No DA		71.9	20.3	72.9	25.8	71.3	27.5	71.3	16.6
Distribution Alignment	DANN	97.8	96.0	83.5	25.1	70.0	01.1	57.8	00.9
	VADA	98.0	96.2	88.2	48.9	78.2	06.6	61.9	01.4
Relaxed Distribution Alignment	IWDAN	97.5	95.7	95.7	81.3	86.5	15.2	74.4	07.3
	IWCDAN	98.0	96.6	96.7	85.1	91.3	66.5	77.5	22.2
	sDANN-4	87.4	05.6	94.9	85.7	86.8	21.6	81.5	39.3
Support Alignment (ours)	ASA-sq	93.7	89.2	92.3	83.5	90.9	69.9	87.2	62.5
	ASA-abs	94.1	88.9	92.8	78.9	92.5	82.4	90.4	68.4

Domain adaptation results

STL \rightarrow CIFAR10

Deep CNN

Algorithm	$\alpha = 0.0$ no shift		$\alpha = 2.0$ severe shift	
	average	min	average	min
No DA	69.9	49.8	65.8	43.7
DANN	75.3	54.6	63.3	27.0
VADA	76.7	56.9	63.2	25.5
IWDAN	69.9	50.5	64.4	36.8
IWCDAN	70.1	47.8	64.5	37.0
sDANN-4	71.8	52.1	66.4	39.0
ASA-sq	71.7	52.9	68.1	44.7
ASA-abs	71.6	49.0	67.8	40.9

VisDA-17

ResNet50

Algorithm	$\alpha = 0.0$ no shift		$\alpha = 2.0$ severe shift	
	average	min	average	min
No DA	49.5	22.2	45.3	19.5
DANN	75.4	36.7	43.1	03.6
VADA	75.3	40.5	43.9	08.5
IWDAN	73.2	31.7	45.1	04.6
IWCDAN	71.6	27.6	38.3	00.6
sDANN-4	72.4	37.8	50.7	18.6
ASA-sq	64.9	35.7	51.9	18.3
ASA-abs	64.8	40.6	52.5	19.7

Summary

- **Support alignment** as extreme relaxation of **distribution alignment**
- **Adversarial Support Alignment (ASA)**: method to align distribution supports
- **Optimal Transport perspective**: spectrum of **relaxed alignment** approaches
- Evaluation on domain adaptation datasets under **label distribution shift**