# Rare Event Modeling

Tony Tran

*A rare event is an observation that occurs with a low probability under a given distribution*

**Rare events can be expensive**

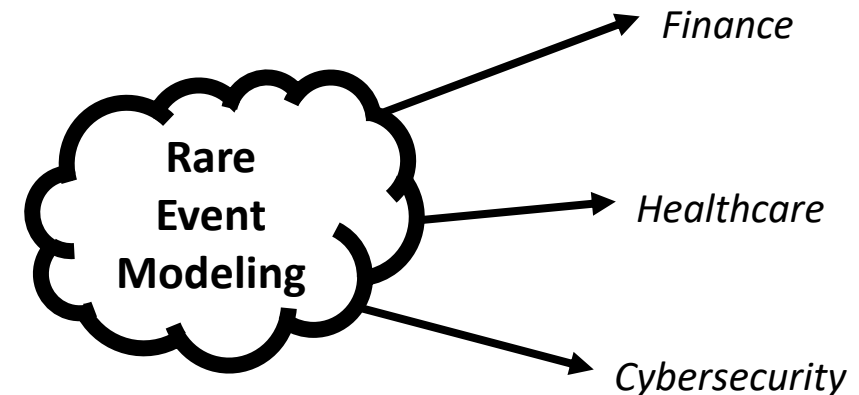**Understand underlying cause**

*Chernobyl 1986*


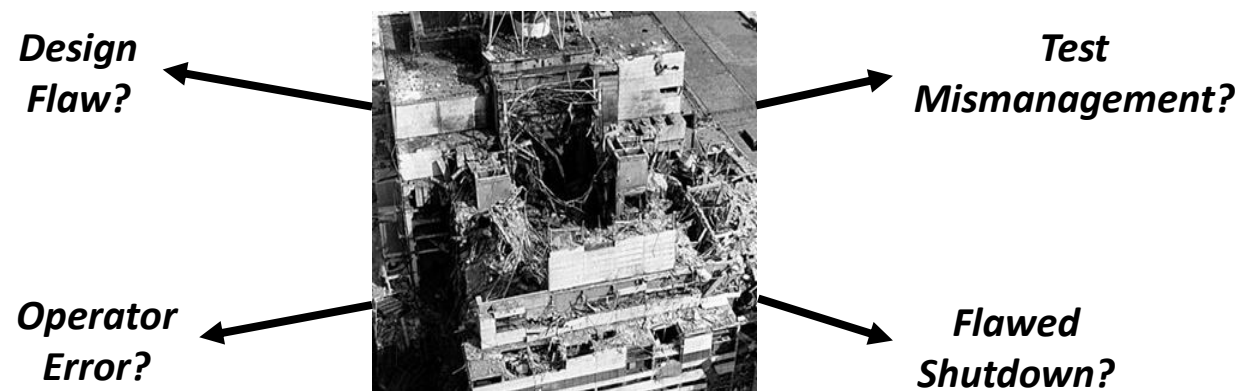
*$700 billion USD*

*30 immediate deaths*

*design flaw?*

*operator error?*

**Key in safety-critical systems**

**Broad relevance across domains**

*Collision Warning*



*Finance*

**Rare Event Modeling**

*Healthcare*

*Cybersecurity*

*Question: What can we learn given the rare event has occurred?*

*Design Flaw?*

*Test Mismanagement?*

*Operator Error?*

*Flawed Shutdown?*

*Task: Model Posterior Over Latent Variables $p(z|x)$*

$$\underbrace{q_\phi(z)}_{\substack{approximate \\ posterior}} \approx \underbrace{p(z|x)}_{\substack{true \\ posterior}} \text{ with } q_\phi(z) \text{ modeled via normalizing flow}$$

*Where $z$ is the latent variable representing hidden causes and $x$ is the observed rare event*

$$z = f_\theta(u)$$

$$z \sim q_\phi(z) \approx p(z|x)$$

$$u \sim \mathcal{N}(0, I_d)$$

*Where $f_\theta$ is a composition of invertible, differentiable functions estimated by a neural network*

**Simple Distribution**
$$\mathcal{N}(0, I_d)$$

**Estimated Distribution**
$$q_\phi(z) \approx p(z|x)$$

UNIVERSITY OF
HOUSTON

**Inverse Bayesian Problem**

Posterior is intractable $\longrightarrow$ $\underbrace{p(z|\mathcal{D})}_{ground\ truth} = \dfrac{p(\mathcal{D}|z)p(z)}{p(\mathcal{D})}$ $\longleftarrow$ Because denominator term is intractable

Where $\mathcal{D}$ is the dataset

***Use Variational Inference!***

**Conditional Evidence Lower BOund (ELBO) Objective**

$$\mathcal{L}(\phi, c, \mathcal{D}) = \mathbb{E}_{(x,y)\in\mathcal{D}, z \sim q_\phi(z|c)}\big[\log p(x,z|y) - \log q_\phi(z|c)\big]$$

$c$ is the conditional (guidance) vector $\longrightarrow$ 

$$\phi^* = \arg\max_\phi \mathcal{L}(\phi, c, \mathcal{D})$$

$\log p(x|y)$ $\longleftarrow$ Maximum Likelihood Estimation maximizes likelihood or "evidence" directly

Posterior Learning aims to minimize KL divergence between learned and true posterior $\longrightarrow$ $D_{KL}\big(q_\phi(z|c) \| p(z|\mathcal{D})\big)$

Variational Inference maximizes "ELBO"

$\mathcal{L}(\phi, c, \mathcal{D})$ $\longleftarrow$

UNIVERSITY OF HOUSTON

*Collection of data is challenging and expensive in some domains*

**Data Scarcity**

**Imbalanced Dataset**

*Environment*

< 0.1% chance of occurrence

> 99.9% chance of occurrence

*few rare event samples*

*abundant nominal samples*

**Overfitting**

**Distribution Shift**

*few rare samples*

*Simulated*

*Real*

*Good Fit*

*Over Fit*

*few rare samples*

Domain Shift

$$N_t \ll N_0$$



Assumption: Nominal and rare observations are from the same domain

Nominal
Target

**True Distribution**

**Few Rare Events**

**Overfit Distribution**

**Few Data Causes Overfitting**

**Regularizes**

**Solution:**

**1. Train a single flow model that represents both nominal and target distributions, guided by a guidance vector**

**2. Regularize the target distribution using the nominal distribution**

$c = [0, 0, 0, 0, 0]$

$c = [1, 1, 1, 1, 1]$

**Guidance Vector $c$**

**Nominal Distribution $q_{\phi_0}$**
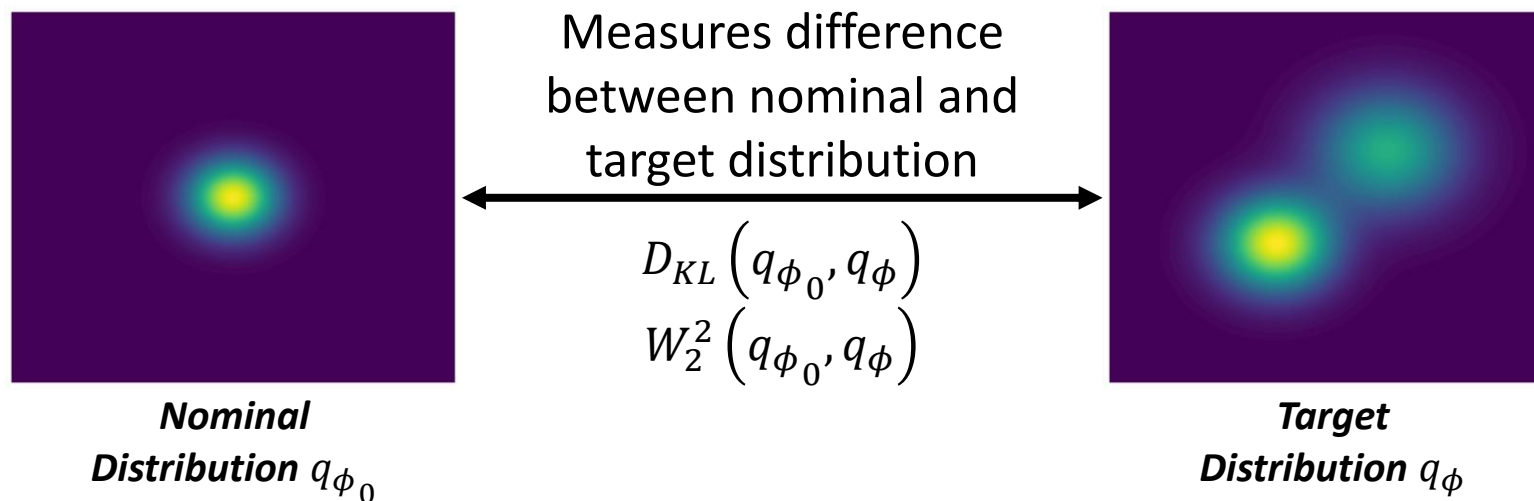
**Target Distribution $q_\phi$**

UNIVERSITY OF HOUSTON

Method 1: KL Divergence

$$J(\phi) = \underbrace{\mathcal{L}(\phi, \mathbf{1}, \mathcal{D}_t)}_{target\ distribution} + \underbrace{\mathcal{L}(\phi, \mathbf{0}, \mathcal{D}_0)}_{nominal\ distribution} - \underbrace{\beta D_{KL}\left(q_{\phi_0}, q_\phi\right)}_{KL\ divergence\ penalty}$$
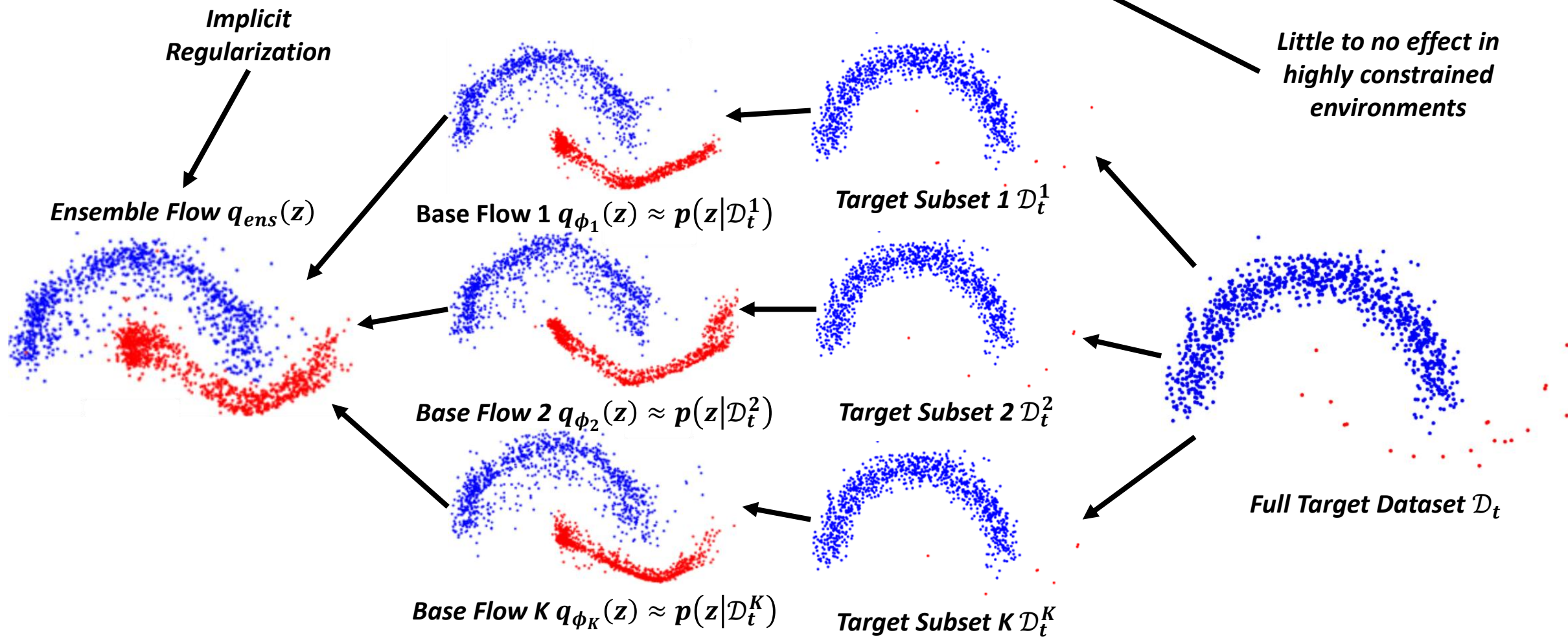
β difficult to tune

Method 2: Wasserstein Divergence

$$J(\phi) = \underbrace{\mathcal{L}(\phi, \mathbf{1}, \mathcal{D}_t)}_{target\ distribution} + \underbrace{\mathcal{L}(\phi, \mathbf{0}, \mathcal{D}_0)}_{nominal\ distribution} - \underbrace{\beta W_2^2\left(q_{\phi_0}, q_\phi\right)}_{Wasserstein\ distance\ penalty}$$



Measures difference between nominal and target distribution

$$D_{KL}\left(q_{\phi_0}, q_\phi\right)$$
$$W_2^2\left(q_{\phi_0}, q_\phi\right)$$

**Nominal Distribution** $q_{\phi_0}$

**Target Distribution** $q_\phi$

Method 3: Ensemble Method

$$\underbrace{q_{ens}(z)}_{\text{ensemble flow}} = \frac{1}{K} \sum_{i=1}^{K} \underbrace{q_{\phi_i}(z)}_{\text{base flows}}$$



**Implicit Regularization**

**Little to no effect in highly constrained environments**

**Ensemble Flow $q_{ens}(z)$**

**Base Flow 1 $q_{\phi_1}(z) \approx p(z|\mathcal{D}_t^1)$**

**Target Subset 1 $\mathcal{D}_t^1$**

**Base Flow 2 $q_{\phi_2}(z) \approx p(z|\mathcal{D}_t^2)$**

**Target Subset 2 $\mathcal{D}_t^2$**

**Base Flow K $q_{\phi_K}(z) \approx p(z|\mathcal{D}_t^K)$**

**Target Subset K $\mathcal{D}_t^K$**

**Full Target Dataset $\mathcal{D}_t$**

*Question:*

**1. How to adaptively choose regularization strength $\beta$?**

$$J(\phi) = \underbrace{\mathcal{L}(\phi, \mathbf{1}, \mathcal{D}_t)}_{target\ distribution} + \underbrace{\mathcal{L}(\phi, \mathbf{0}, \mathcal{D}_0)}_{nominal\ distribution} - \underbrace{\beta D_{KL}\left(q_{\phi_0}, q_\phi\right)}_{KL\ divergence\ penalty}$$

$$J(\phi) = \underbrace{\mathcal{L}(\phi, \mathbf{1}, \mathcal{D}_t)}_{target\ distribution} + \underbrace{\mathcal{L}(\phi, \mathbf{0}, \mathcal{D}_0)}_{nominal\ distribution} - \underbrace{\beta W_2^2\left(q_{\phi_0}, q_\phi\right)}_{Wasserstein\ distance\ penalty}$$

**2. How to share information between flows to learn robustly?**

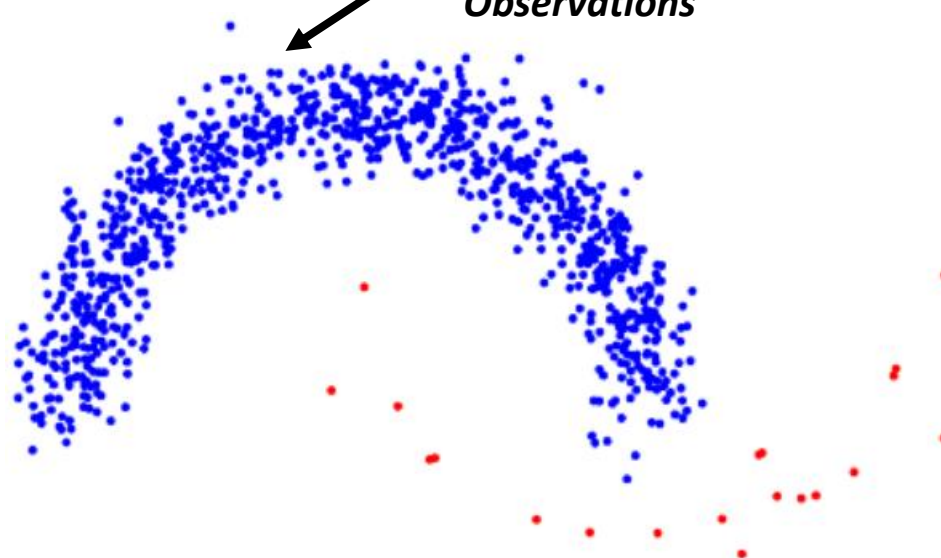$$\underbrace{q_{ens}(z)}_{ensemble\ flow} = \frac{1}{K}\sum_{i=1}^{K}\underbrace{q_{\phi_i}(z)}_{base\ flows}$$

Method 4: Self-Regularization

$$J(\phi, c) = \frac{1}{K}\sum_{i=1}^{K}\underbrace{\mathcal{L}(\phi, \mathbf{1}_i, \mathcal{D}_t^i)}_{base\ flows} + \underbrace{\mathcal{L}(\phi, \mathbf{0}_K, \mathcal{D}_0)}_{nominal\ flow} + \underbrace{\mathcal{L}(\phi, c, \mathcal{D}_t)}_{ensemble\ flow} - \beta\sum_{i\neq j}^{K}\underbrace{D_{KL}\left(q_\phi(\cdot; \mathbf{1}_i), q_\phi(\cdot; \mathbf{1}_j)\right)}_{penalty\ between\ base\ flows}$$
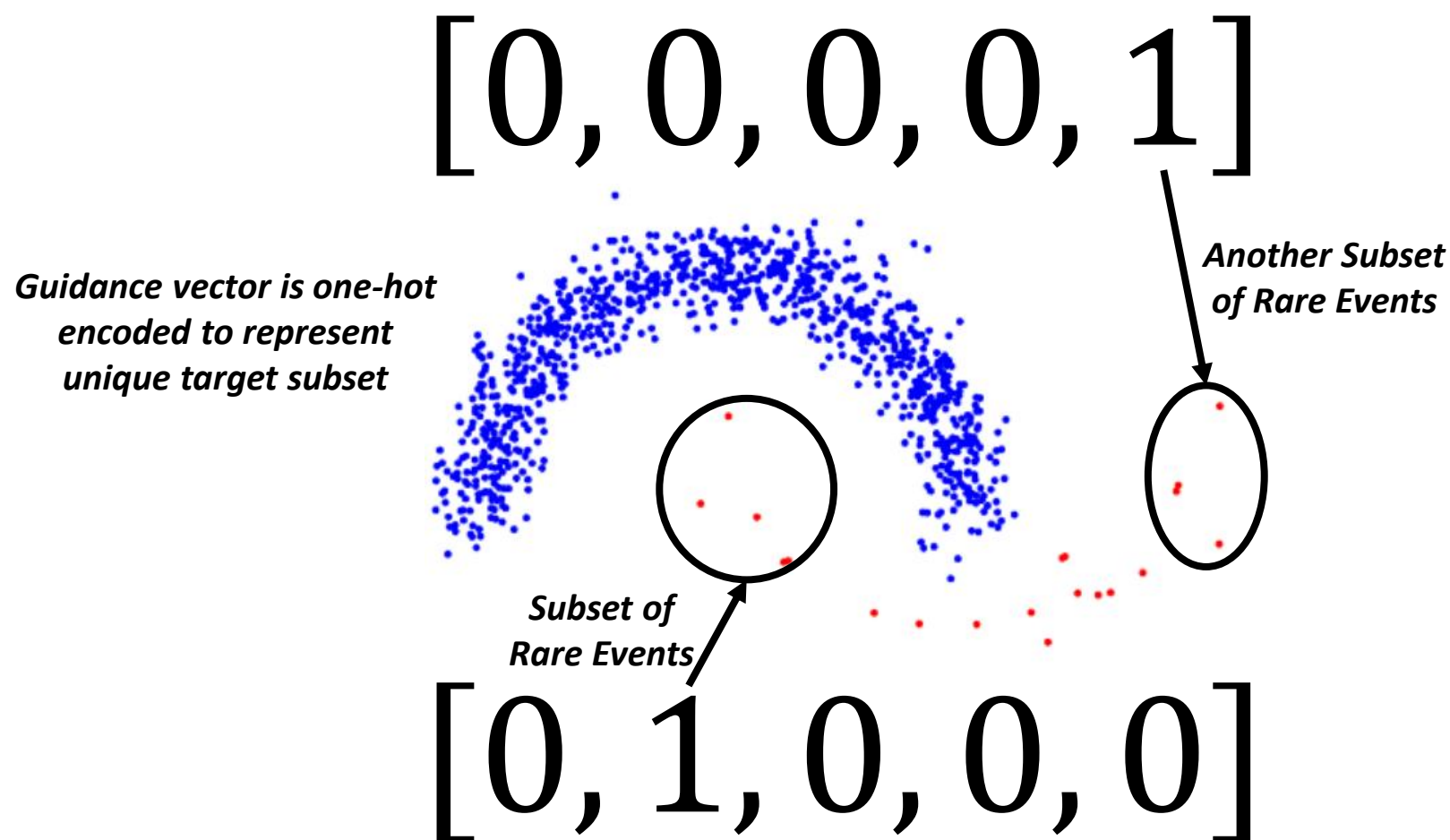
$$J(\phi, c) = \frac{1}{K} \sum_{i=1}^{K} \underbrace{\mathcal{L}(\phi, \mathbf{1}_i, \mathcal{D}_t^i)}_{base\ flows} + \underbrace{\mathcal{L}(\phi, \mathbf{0}_K, \mathcal{D}_0)}_{nominal\ flow} + \underbrace{\mathcal{L}(\phi, c, \mathcal{D}_t)}_{ensemble\ flow} - \beta \sum_{i \neq j}^{K} \underbrace{D_{KL}\left(q_\phi(\cdot; \mathbf{1}_i), q_\phi(\cdot; \mathbf{1}_j)\right)}_{penalty\ between\ base\ flows}$$

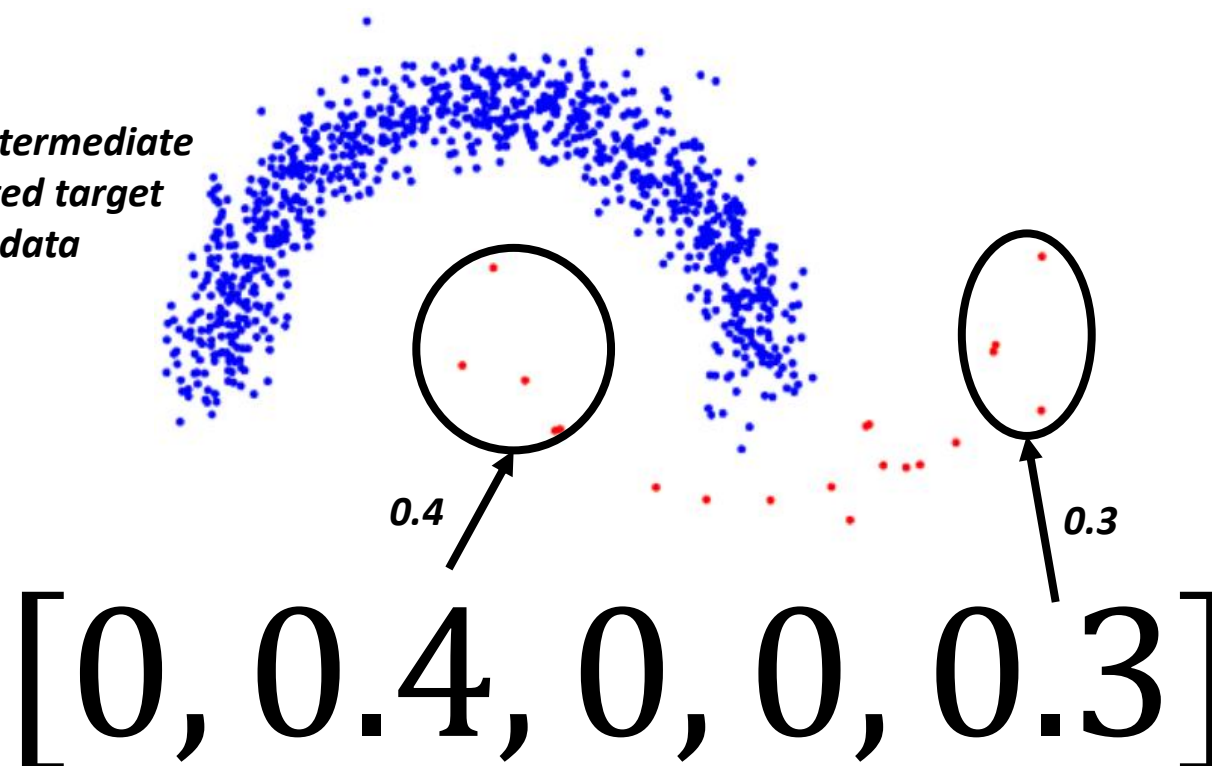$$[0, 0, 0, 0, 0]$$ *Guidance vector*

*Nominal Observations*

$$J(\phi, c) = \frac{1}{K} \sum_{i=1}^{K} \underbrace{\mathcal{L}(\phi, \mathbf{1}_i, \mathcal{D}_t^i)}_{base\ flows} + \underbrace{\mathcal{L}(\phi, \mathbf{0}_K, \mathcal{D}_0)}_{nominal\ flow} + \underbrace{\mathcal{L}(\phi, c, \mathcal{D}_t)}_{ensemble\ flow} - \beta \sum_{i \neq j}^{K} \underbrace{D_{KL}\left(q_\phi(\cdot; \mathbf{1}_i), q_\phi(\cdot; \mathbf{1}_j)\right)}_{penalty\ between\ base\ flows}$$

$$[0, 0, 0, 0, 1]$$

*Guidance vector is one-hot encoded to represent unique target subset*

**Another Subset of Rare Events**

**Subset of Rare Events**

$$[0, 1, 0, 0, 0]$$

$$J(\phi, c) = \frac{1}{K} \sum_{i=1}^{K} \underbrace{\mathcal{L}(\phi, \mathbf{1}_i, \mathcal{D}_t^i)}_{base\ flows} + \underbrace{\mathcal{L}(\phi, \mathbf{0}_K, \mathcal{D}_0)}_{nominal\ flow} + \underbrace{\textcolor{red}{\mathcal{L}(\phi, c, \mathcal{D}_t)}}_{\textcolor{red}{ensemble\ flow}} - \beta \sum_{i \neq j}^{K} \underbrace{D_{KL}\Big(q_\phi(\cdot; \mathbf{1}_i), q_\phi(\cdot; \mathbf{1}_j)\Big)}_{penalty\ between\ base\ flows}$$
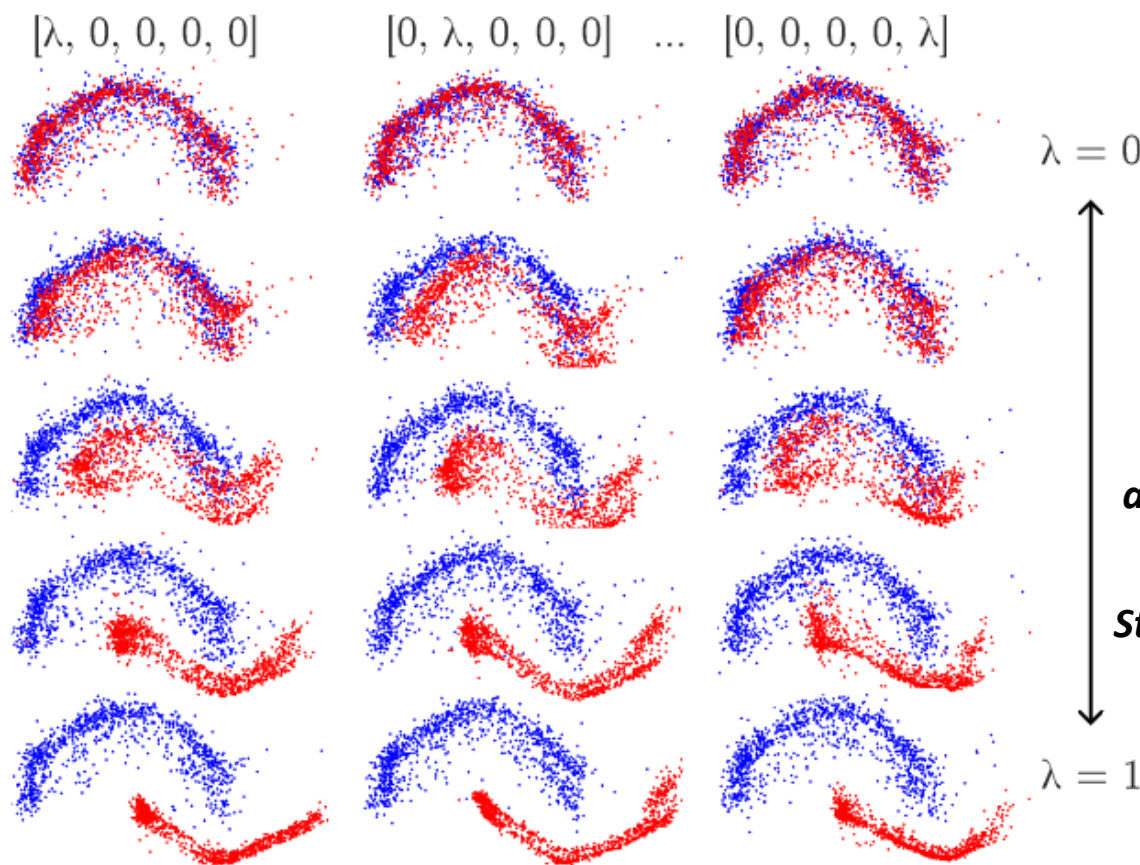
*Ensemble flow learns intermediate flows between weighted target subsets of target data*



0.4

0.3

$$[0, 0.4, 0, 0, 0.3]$$

$$J(\phi, c) = \frac{1}{K}\sum_{i=1}^{K}\underbrace{\mathcal{L}(\phi, \mathbf{1}_i, \mathcal{D}_t^i)}_{base\ flows} + \underbrace{\mathcal{L}(\phi, \mathbf{0}_K, \mathcal{D}_0)}_{nominal\ flow} + \underbrace{\mathcal{L}(\phi, c, \mathcal{D}_t)}_{ensemble\ flow} - \beta\sum_{i\neq j}^{K}\underbrace{D_{KL}\Big(q_\phi(\cdot; \mathbf{1}_i), q_\phi(\cdot; \mathbf{1}_j)\Big)}_{penalty\ between\ base\ flows}$$

$[\lambda, 0, 0, 0, 0]$  $[0, \lambda, 0, 0, 0]$  ...  $[0, 0, 0, 0, \lambda]$
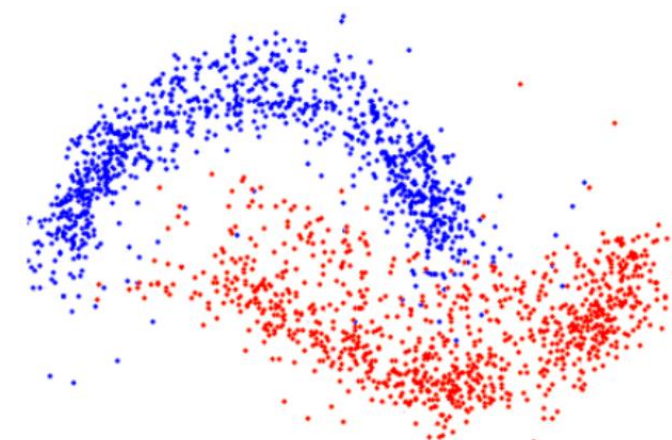


$\lambda = 0$

**Step 1: Yield a family of flows**

*$\lambda$ represents regularization coefficient which weights the nominal and target distribution*
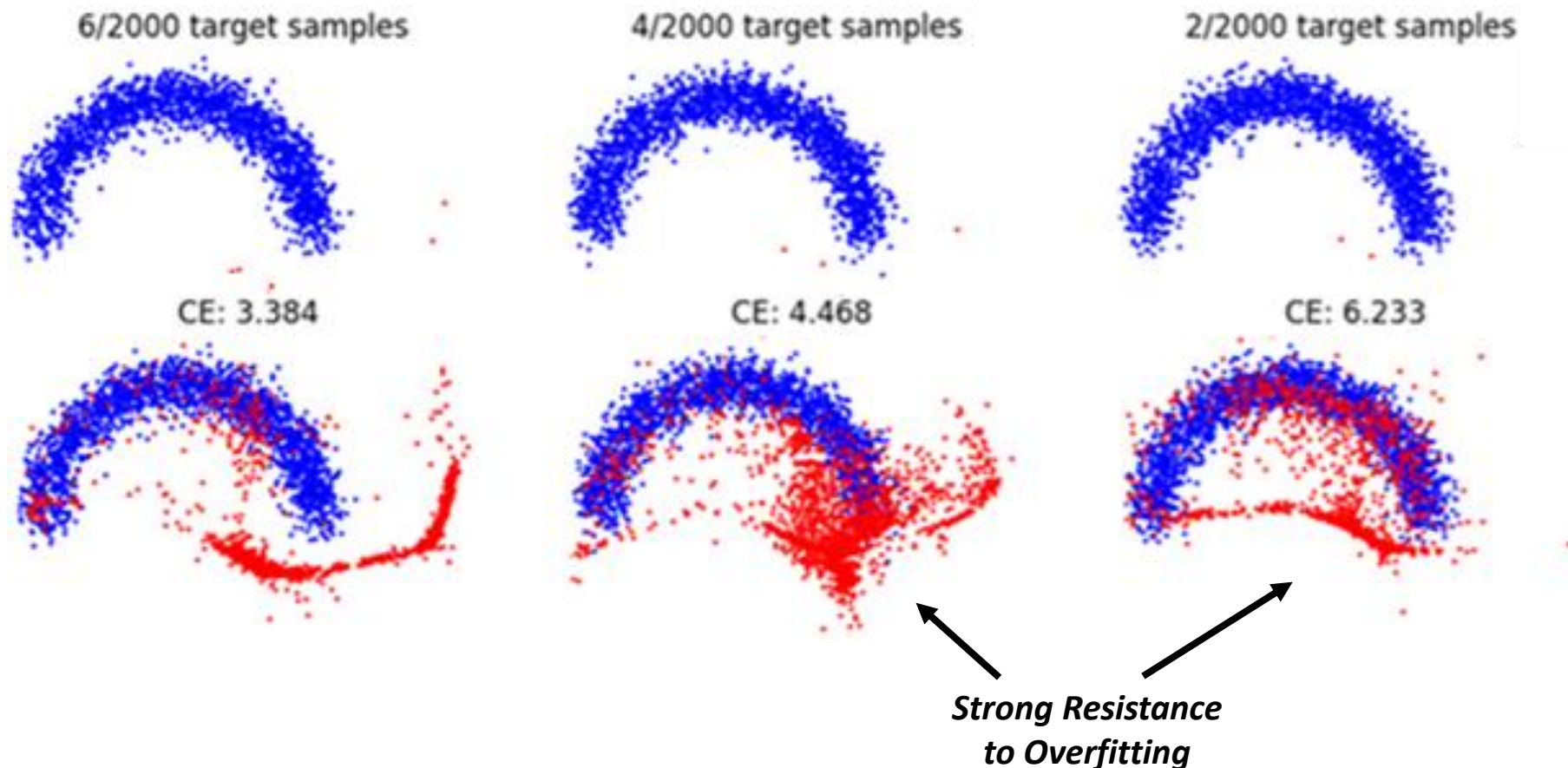
*Bigger $\lambda$ Stronger Regularization*

$\lambda = 1$

Step 2:

Solve optimal Guidance Vector
$$c^* = \arg\max_c \mathcal{L}(\phi^*, c, \mathcal{D}_t)$$
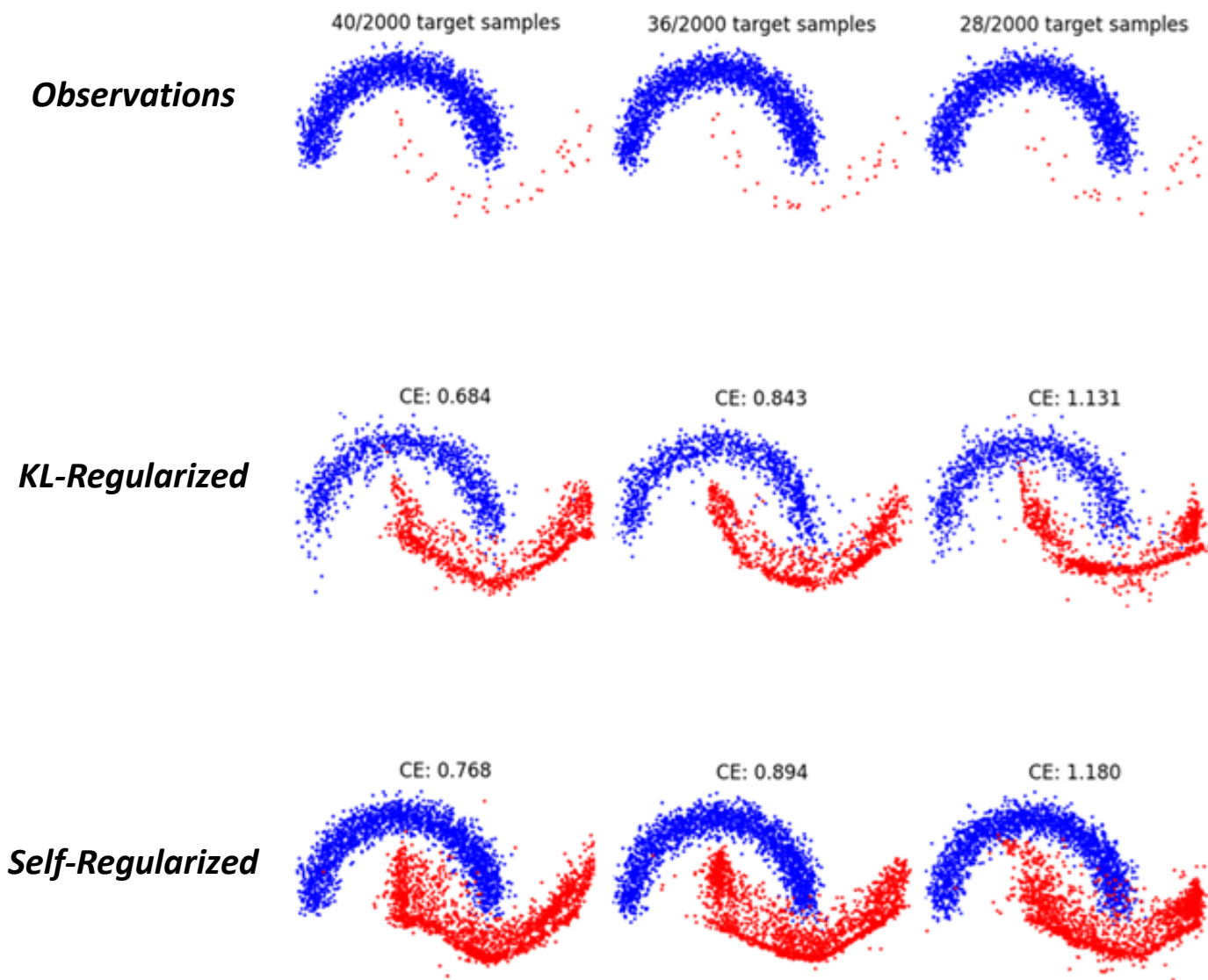
*Limitations?*

**UNIVERSITY OF HOUSTON**

*Hypothesis: Self-regularization may still overfit when rare event observations are extremely limited.*



6/2000 target samples

CE: 3.384

4/2000 target samples

CE: 4.468

2/2000 target samples

CE: 6.233

*Strong Resistance to Overfitting*

**Hypothesis: Simple prior regularization methods can outperform self-regularization with more rare event samples.**



**Observations**

40/2000 target samples     36/2000 target samples     28/2000 target samples

Cross Entropy Test:

$$CE\left(p(z|x), q_\phi(z)\right) = -\sum p(z|x) \log q_\phi(z)$$

Lower cross entropy is better.

**There exist instances in which simple KL-regularized paradigms outperform CalNF.**

**KL-Regularized**

CE: 0.684     CE: 0.843     CE: 1.131

**Self-Regularized**

CE: 0.768     CE: 0.894     CE: 1.180

$$J(\phi, c) = \frac{1}{K} \sum_{i=1}^{K} \underbrace{\mathcal{L}(\phi, \mathbf{1}_i, \mathcal{D}_t^i)}_{specialization\ term} + \mathcal{L}(\phi, \mathbf{0}_K, \mathcal{D}_0) + \underbrace{\mathcal{L}(\phi, c, \mathcal{D}_t)}_{generalization\ term} - \beta \sum_{i \neq j}^{K} D_{KL}\left(q_\phi(\cdot; \mathbf{1}_i), q_\phi(\cdot; \mathbf{1}_j)\right)$$

**Pushes φ to learn flows to specialize per task**

**Hypothesis: There exists inherent conflict between the specialization and generalization terms in the objective within parameters.**

**Pushes φ to learn flows to generalize across task**

Cosine Similarity Test between Gradients: $S_C(\boldsymbol{a}, \boldsymbol{b}) = \frac{\boldsymbol{a} \cdot \boldsymbol{b}}{\|\boldsymbol{a}\| \|\boldsymbol{b}\|}$
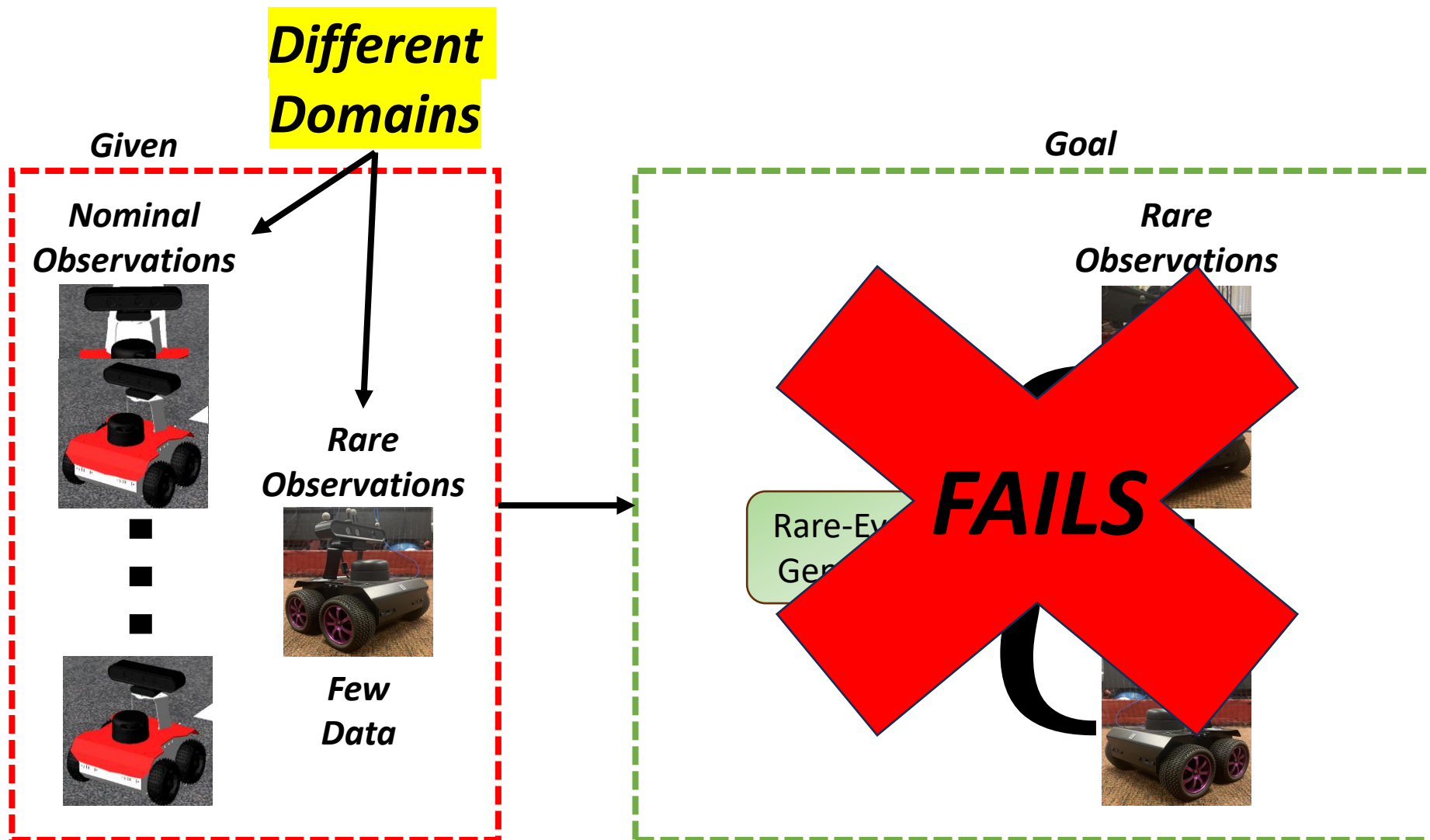
$$K = 2$$

$$S_C\left(\nabla_\theta \mathcal{L}(\phi, \mathbf{1}_1, \mathcal{D}_t^1), \nabla_\theta \mathcal{L}(\phi, c^*, \mathcal{D}_t)\right) = +0.8799$$

$$S_C\left(\nabla_\theta \mathcal{L}(\phi, \mathbf{1}_2, \mathcal{D}_t^2), \nabla_\theta \mathcal{L}(\phi, c^*, \mathcal{D}_t)\right) = +0.6994$$

**Since $S_C$ is positive, gradients roughly align**

Limitation: Nominal and rare observations MUST originate from same domain

# References

[1] Dawson, C., Tran, V., Li, M. Z., & Fan, C. (2025). *Rare event modeling with self-regularized normalizing flows: What can we learn from a single failure?* (arXiv:2502.21110). arXiv. https://doi.org/10.48550/arXiv.2502.21110

[2] Abdollahzadeh, Milad, Touba Malekzadeh, Christopher T. H. Teo, Keshigeyan Chandrasegaran, Guimeng Liu, and Ngai-Man Cheung. "A Survey on Generative Modeling with Limited Data, Few Shots, and Zero Shot." arXiv, July 26, 2023. https://doi.org/10.48550/arXiv.2307.14397.