**Tutorial On**

# PROBABILISTIC U-NETs
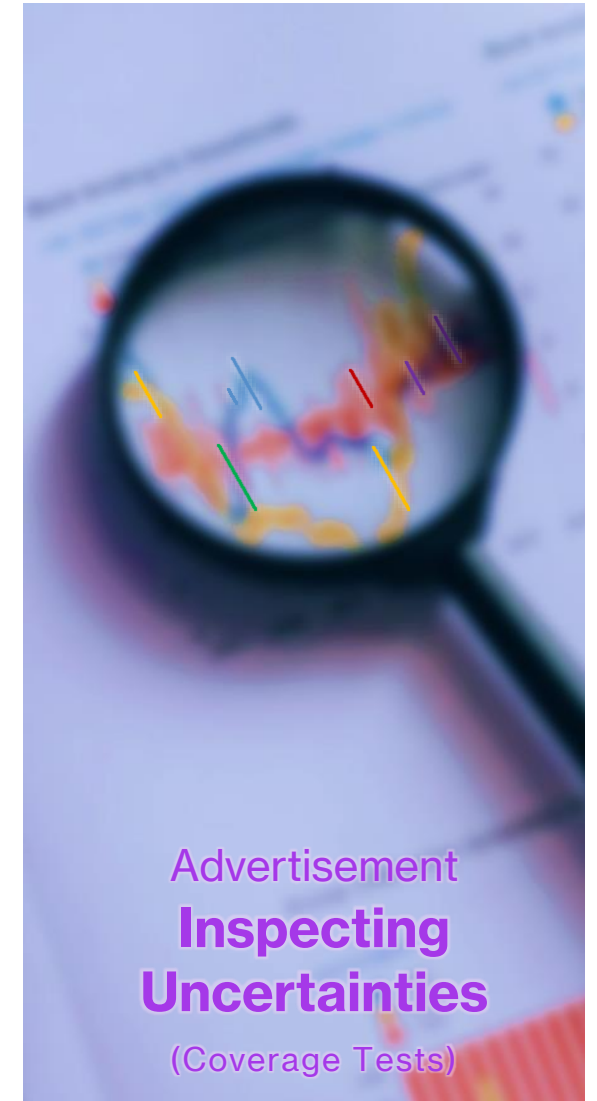
**Hadi Sotoudeh**

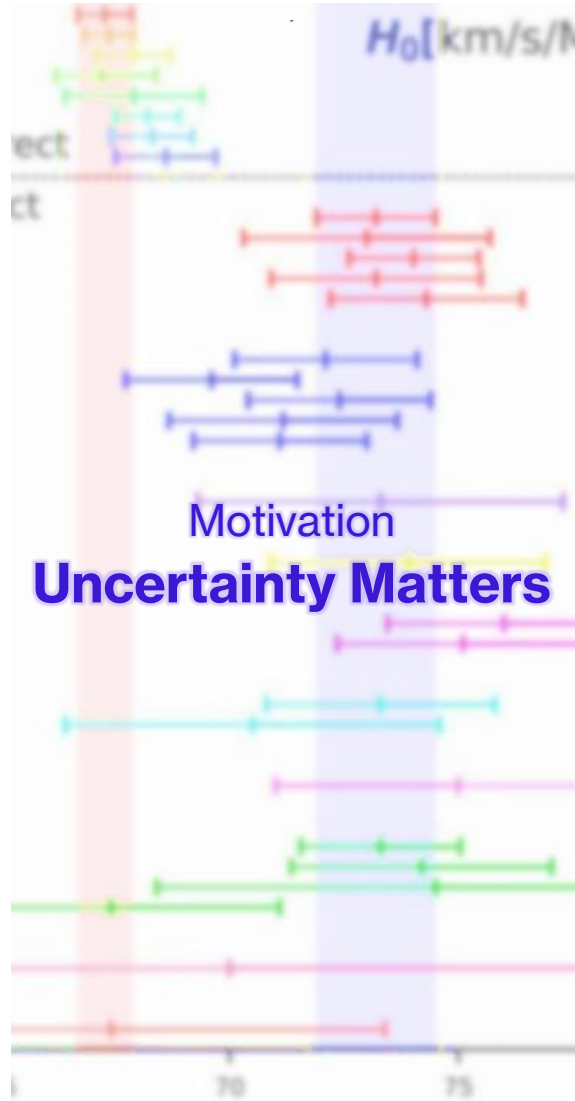KITP Program on Building a Physical Understanding of Galaxy Evolution with Data-driven Astronomy       23 Feb 2023
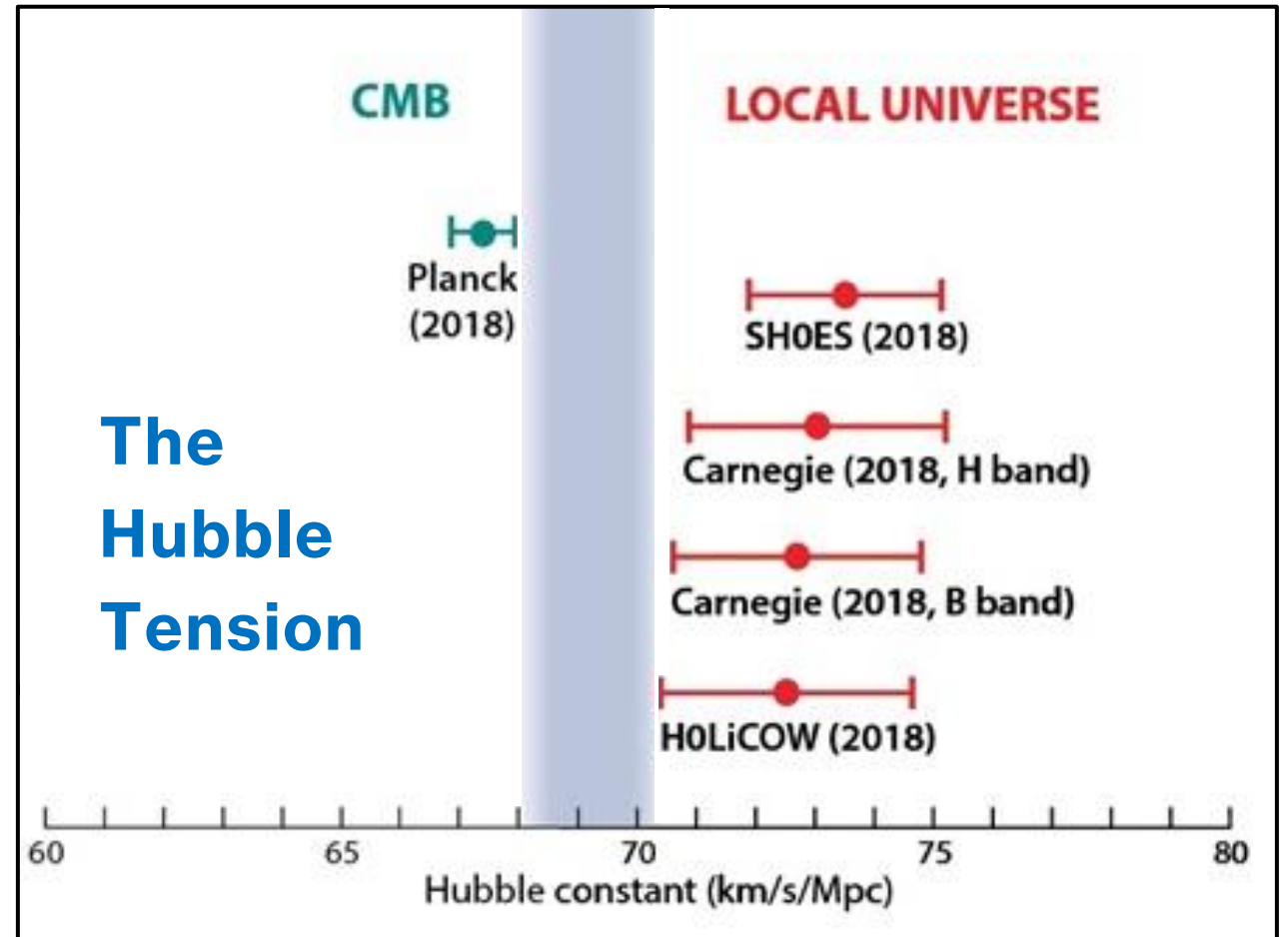
Access Repo

# Outline



Motivation
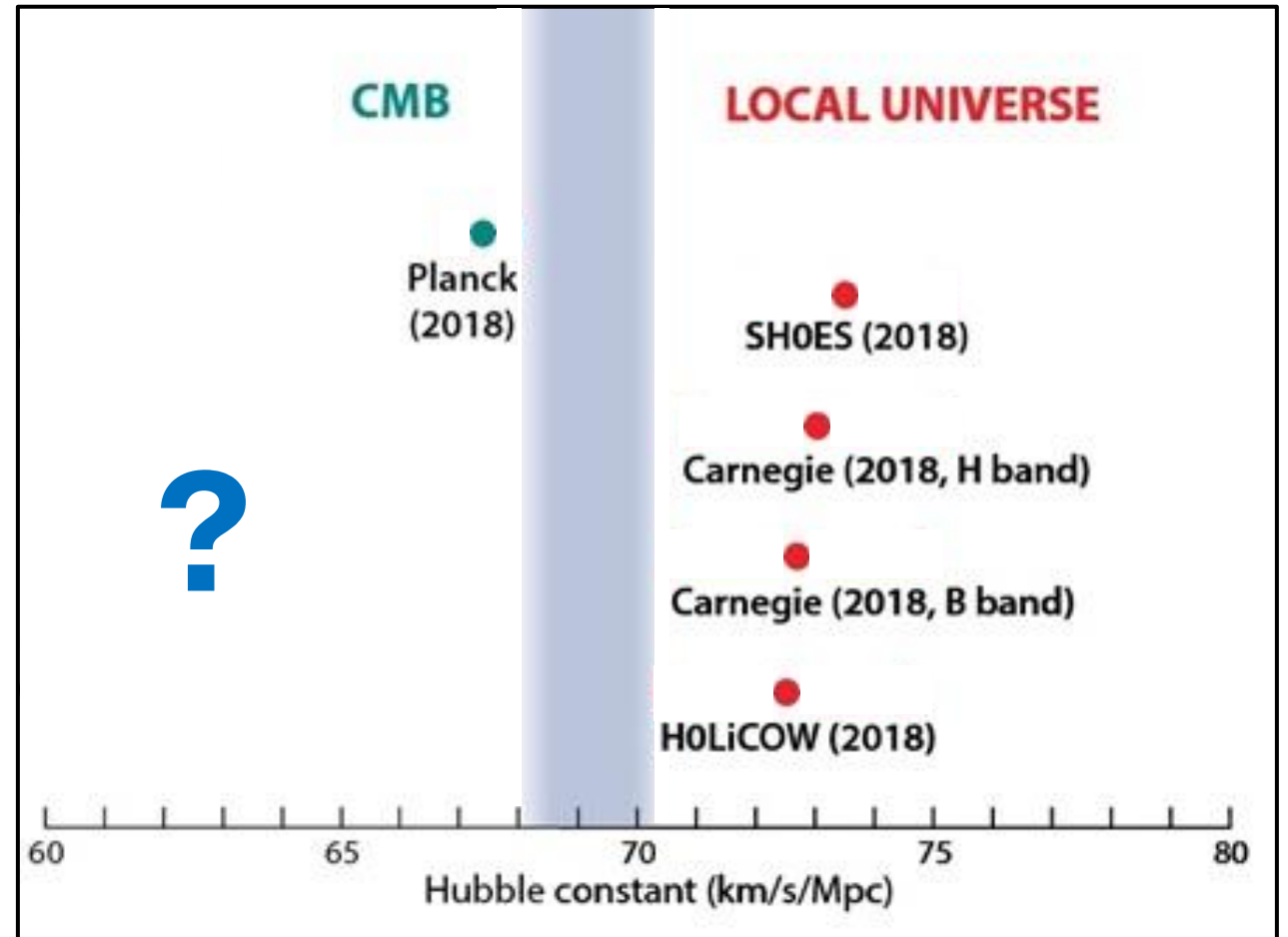**Uncertainty Matters**

Part 1
**The "Probabilistic" U-Net**
(Architecture & Training)

Part 2
**Rescue the Randomness**
(Loss Functions)

Advertisement
**Inspecting Uncertainties**
(Coverage Tests)

# Uncertainty Matters

- Every physical measurement is meaningful with an uncertainty estimate.



**The Hubble Tension**

**Credit:** Roen Kelly, *Astronomy* Magazine – Figure taken with modification.
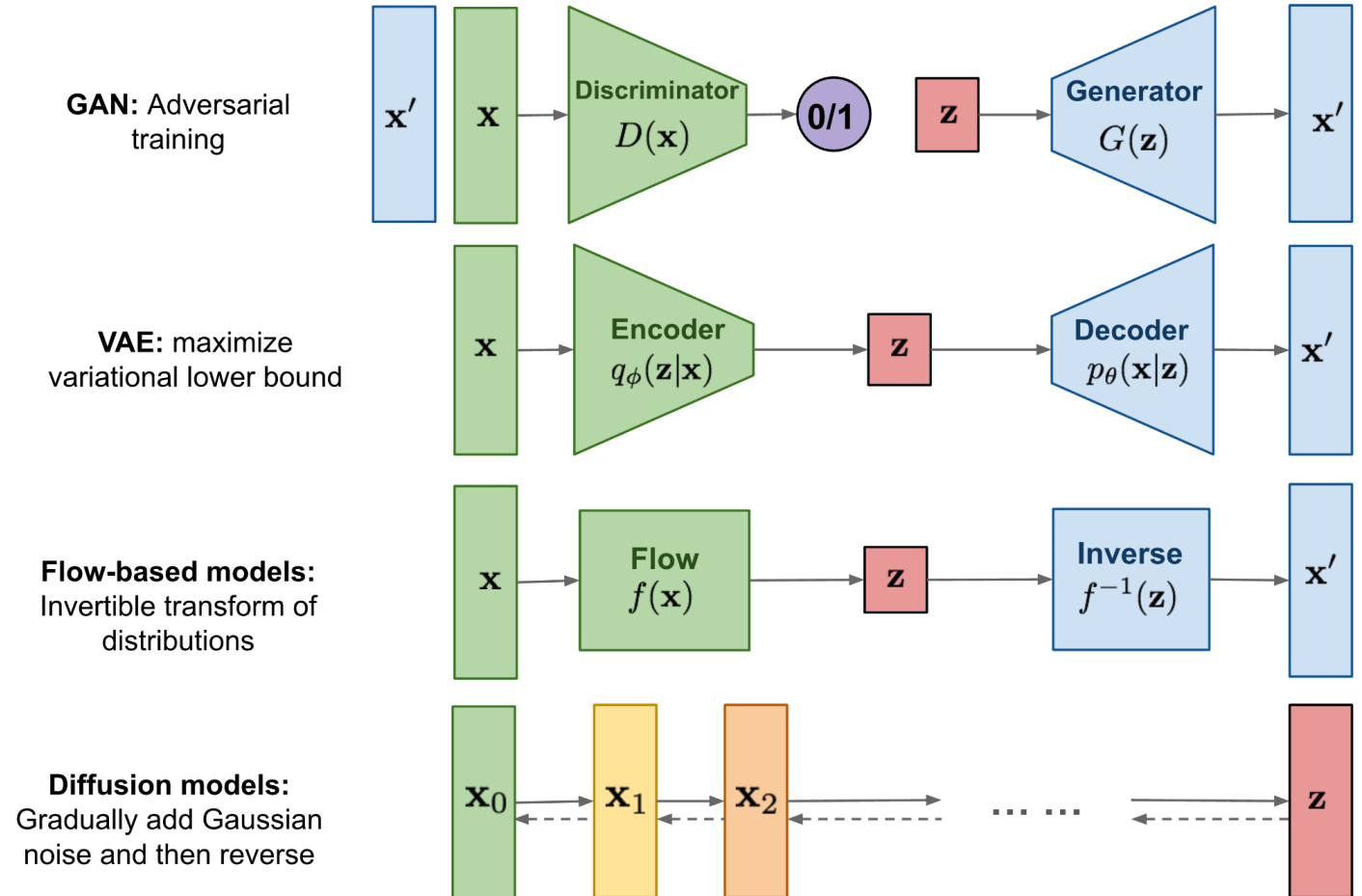
# Deep Learning for Physical Discoveries

- Traditional deep learning methods → No uncertainty estimate

- Physical applications require models capable of quantifying uncertainties.



CMB — Planck (2018)

LOCAL UNIVERSE — SH0ES (2018), Carnegie (2018, H band), Carnegie (2018, B band), H0LiCOW (2018)

Hubble constant (km/s/Mpc)

**?**

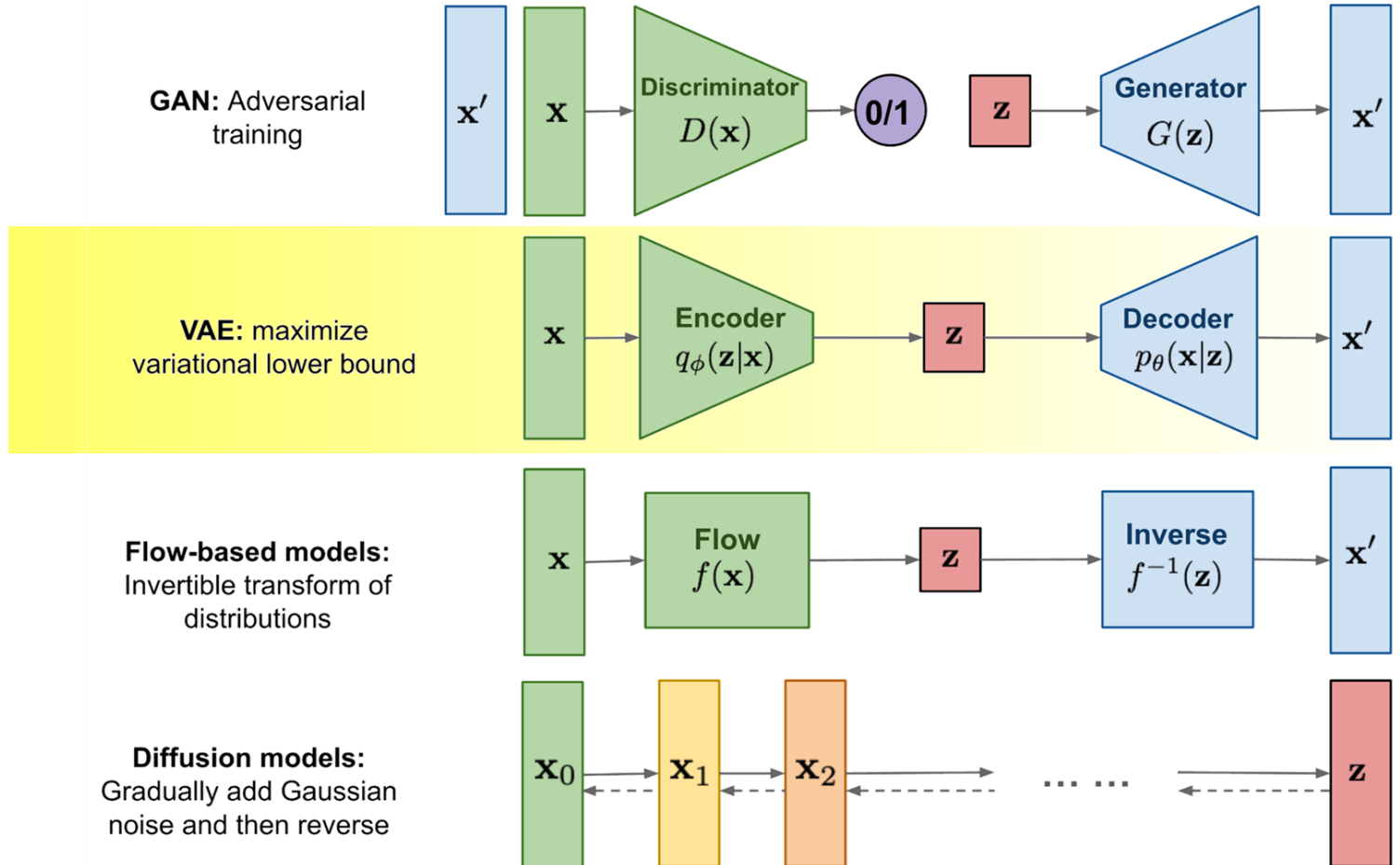#uncertainty_quantification

# Deep Generative Models

- Can learn to generate new data that resembles a distribution → Can encode uncertainties

5

# Prob. U-Net

- Advantages:

  - Relatively lower computational cost during both training & inference

  - Well-understood theoretical framework

Part 1
# The "Probabilistic" U-Net

Autoencoder, VAE & cVAE

U-Net

Probabilistic U-Net

Training

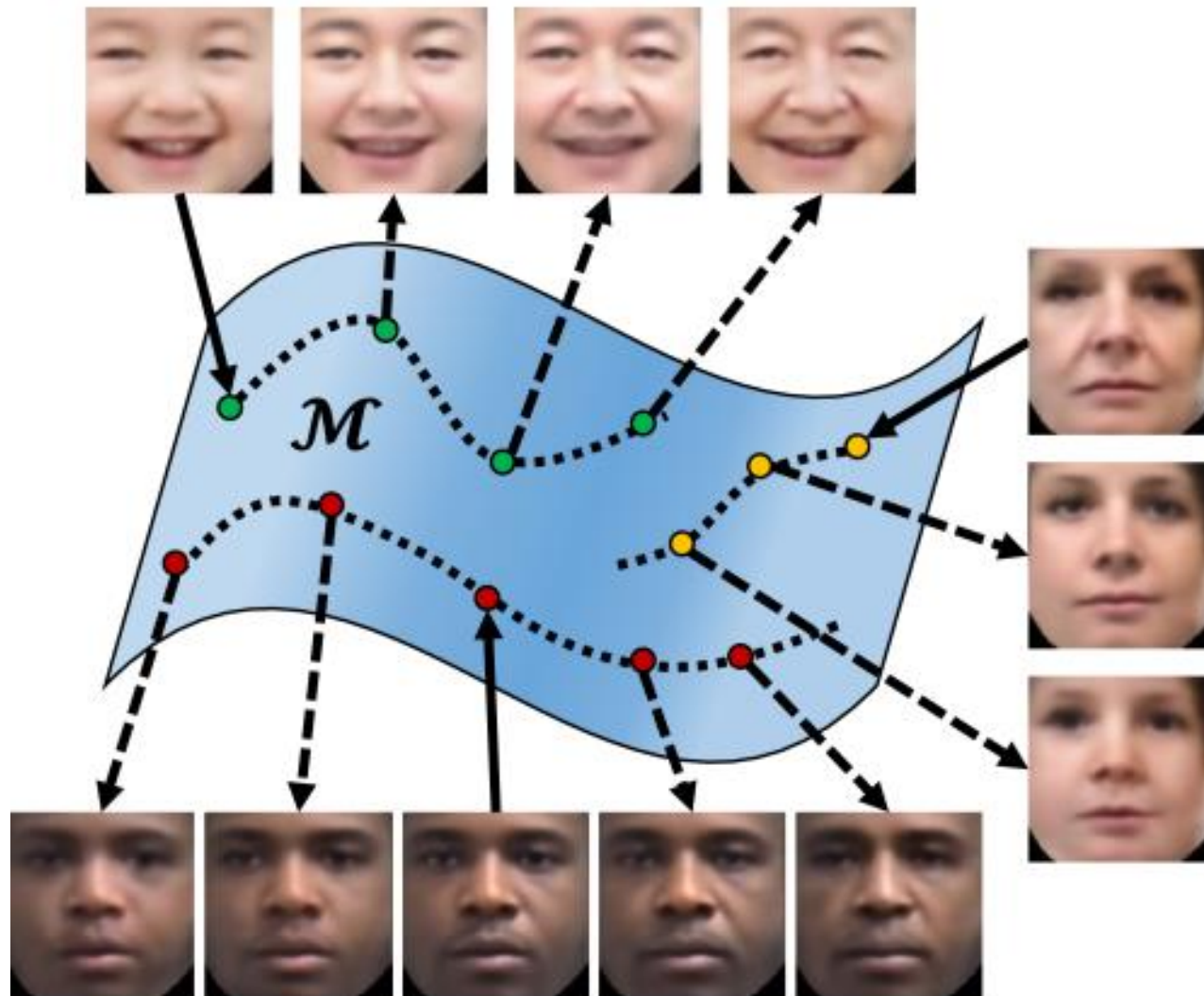Toy Problem 1: Source Reconstruction

# Curse of Dimensionality
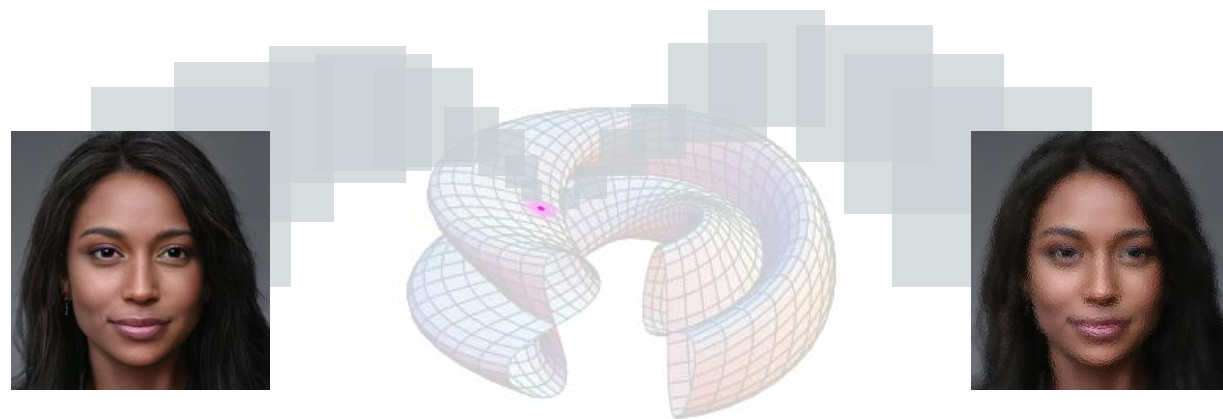


pixel space is overparameterized!

# Manifold Hypothesis

- Many real-world high-dim datasets lie along low-dim **latent manifolds** inside that space

- Manifold of valid human faces
  - If accessible, can easily draw samples from the distribution of valid faces



**Credit:** Zhifei Zhang et al. (2017)

# Latent Space

- Dimensionality lower than data space ($\ell < m$)

- Defined by the **encoder** & **decoder**



**data space**  **latent space**

$$\boldsymbol{x} \in \mathbb{R}^m \qquad \boldsymbol{z} \in \mathbb{R}^\ell$$
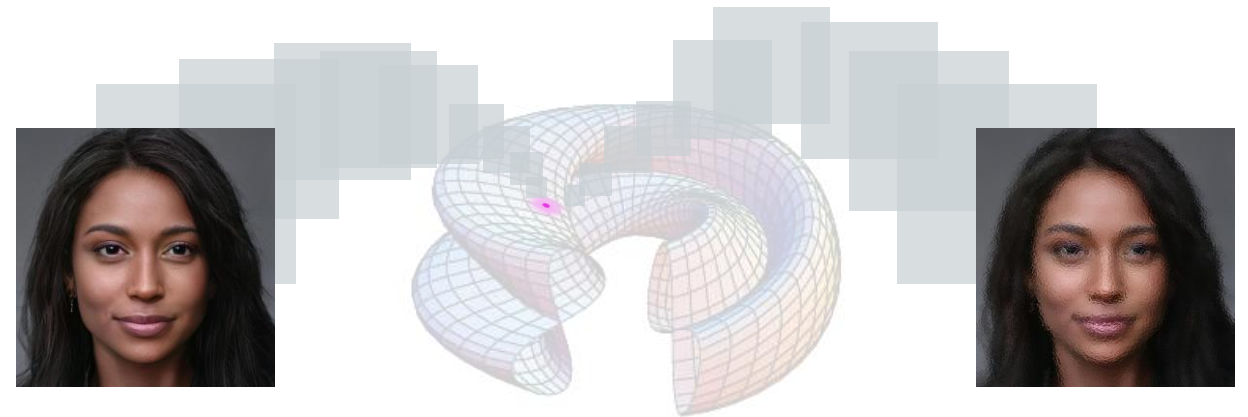


encoder **e**  decoder **d**

# Latent Space

- **Goal:** Best encoder-decoder pair

  - keep maximum information

  - minimize reconstruction error

- Reconstructed Image: $\widehat{x} := d\big(e(x)\big)$

- Examples of reconstruction error:

  - $\mathrm{MSE}(x, \widehat{x}) = \|x - \widehat{x}\|^2$

  - $\mathrm{BCE}(x, \widehat{x}) = -\sum x_i \log \widehat{x}_i + (1 - x_i)\log(1 - \widehat{x}_i)$
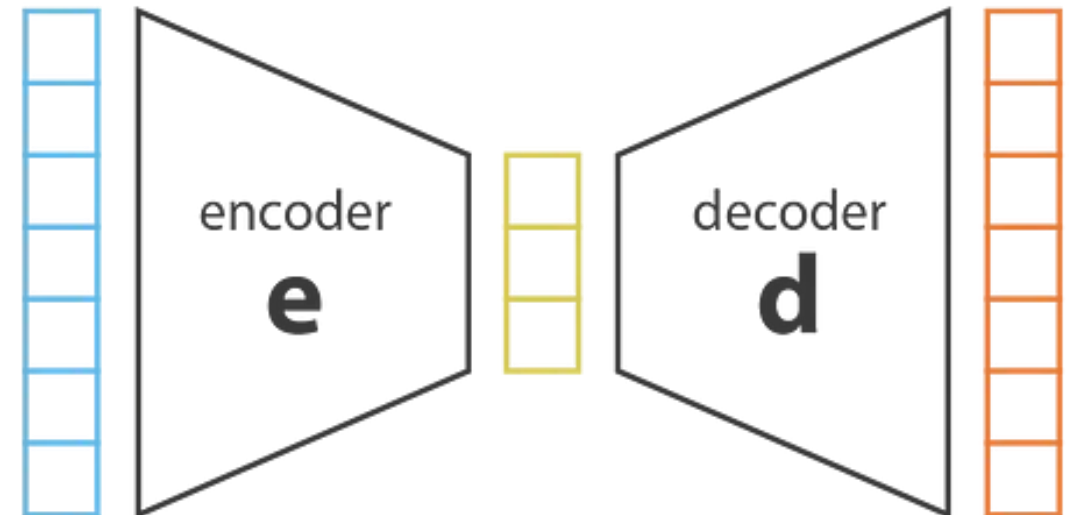
#dimensionality_reduction

#representation_learning



**data space**        **latent space**

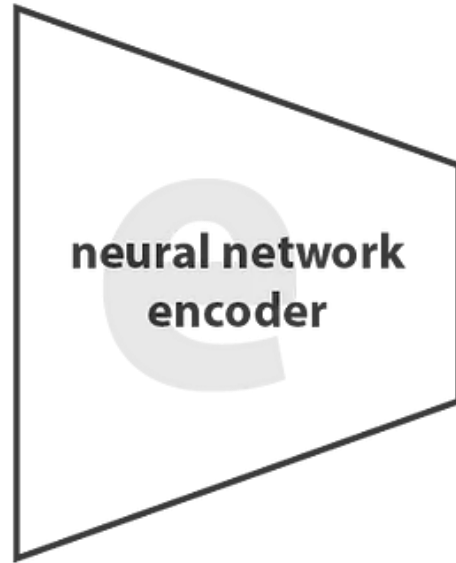$$x \in \mathbb{R}^m \qquad\qquad z \in \mathbb{R}^\ell$$

encoder **e**      decoder **d**

# Autoencoder



$$\mathcal{L}_{\mathrm{rec}} = \mathrm{MSE}(\boldsymbol{x}, \widehat{\boldsymbol{x}})$$

# Problem: Irregular Latent Space

- Autoencoders only focus on reconstruction → Don't care about the structure of latent space

  - Tend to learn **punctual distributions**

- Latent space should be **continuous** and **complete**

# Ideal Latent Space



**Frey Face Dataset**

# Probabilistic Setup

- Learn a mapping from some latent distribution on $z$ to a complicated distribution on $x$

latent space       data space



- Sample from the prior distribution in latent space → Map the sample to data space

$$p(z) = \text{something simple} \qquad p(x|z) \text{ modeled by generator}$$

- Learn representation such that the marginal **data likelihood** (evidence) is maximized:

$$p(x) = \int p(x, z)\, dz \qquad \text{where} \qquad p(x, z) = p(x|z)\, p(z)$$

likelihood    prior

# **Variational Inference**

Problem:  $p(x) = \int p(x, z)\, dz$   is intractable

- Variational inference approach: Find a lower bound for the integral using an auxiliary distribution ($q$)

variational posterior    true posterior

$$\ln p(\boldsymbol{x}) = \underbrace{\int q(\boldsymbol{z}|\boldsymbol{x}) \ln\left(\frac{p(\boldsymbol{x}, \boldsymbol{z})}{q(\boldsymbol{z}|\boldsymbol{x})}\right) d\boldsymbol{z}}_{\textbf{E}\text{vidence } \textbf{L}\text{ower } \textbf{BO}\text{und (ELBO)}} - \underbrace{\int q(\boldsymbol{z}|\boldsymbol{x}) \ln\left(\frac{p(\boldsymbol{z}|\boldsymbol{x})}{q(\boldsymbol{z}|\boldsymbol{x})}\right) d\boldsymbol{z}}_{\text{Variational Gap}}$$

$$\ln p(\boldsymbol{x}) \quad \geq \quad \text{ELBO}$$

This is what we will try to maximize!

# Variational Autoencoder

- Based on variational inference

"implicitly"
models variational posterior

$$q_\phi(z|x)$$

"implicitly"
models likelihood

$$p_\theta(x|z)$$

neural network
encoder

neural network
decoder

x

$\hat{x} = d(z)$

# Variational Autoencoder



$$\mathrm{ELBO}(\theta, \phi, \boldsymbol{x}) = \int q_\phi(\boldsymbol{z}|\boldsymbol{x}) \ln\left(\frac{p_\theta(\boldsymbol{x}, \boldsymbol{z})}{q_\phi(\boldsymbol{z}|\boldsymbol{x})}\right) d\boldsymbol{z}$$

$$\mathcal{L}_{\mathrm{ELBO}} = -\mathrm{ELBO}(\theta, \phi, \boldsymbol{x}) = \mathbb{E}_{q_\phi(\boldsymbol{z}|\boldsymbol{x})}[-\ln p_\theta(\boldsymbol{x}|\boldsymbol{z})] + D_{\mathrm{KL}}\left(q_\phi(\boldsymbol{z}|\boldsymbol{x}) \parallel p_\theta(\boldsymbol{z})\right)$$

$\mathcal{L}_{\mathbf{rec}}$ (*reconstruction* term)     $\mathcal{L}_{\mathrm{KL}}$ (*regularization* term)

# Conditional VAE

- How to generate samples of a particular class?

**VAE**

**cVAE**



Decoder

Decoder

latent space

latent space

3

# Conditional VAE

- Used to generate samples of a particular class on demand

# U-Net

- A type of convolutional neural network architecture → learn image to image mapping



Input

Prediction

Contracting Path

Expanding Path

◻ n x Res-Block   ↘/↗ Down- / Up-Sampling   ☐ Concatenation   → Skip Connection

# What We Have So Far

## cVAE



## U-Net



a deep generative model to **generate new data** based on a noise vector and a set of conditional inputs

a convolutional neural network to learn **image to image mapping**

# What We Need

- Problem Space...

  - high-dimensional observations and parameters

  - noisy observations

  - a **manifold of parameters** consistent with a given observation instead of a deterministic prediction (underconstrained problem)

- We Need...

  - high-dimensional inference

  - quantify uncertainties

  - model variability

# Applications

## Model Output Variability

**Training Set:** 1 observation ↔ multiple predictions

**Example:** Different doctors assign different lesion areas on lung CT scans



CT Scan

$x$

Segmentation Samples
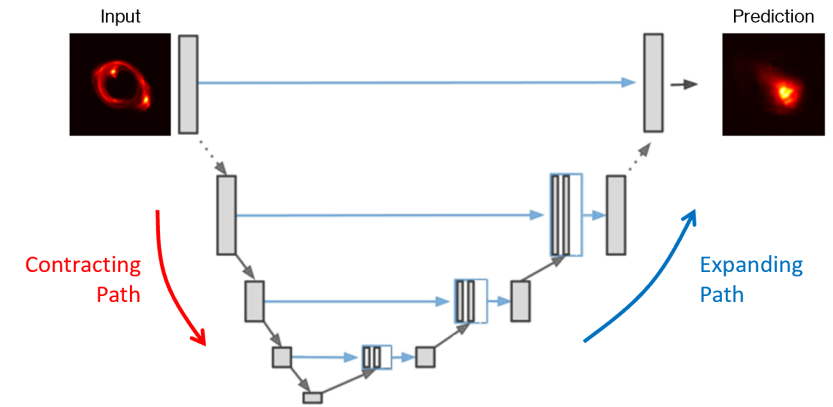
$y$

## Inverse Problems

**Training Set:** multiple observations ↔ 1 prediction

**Example:** Reconstruct the initial conditions of the Universe

forward model

$$x = f(y) + \text{noise}$$



Initial Conditions

$y$

Observed Galaxy Distribution

$x$

# Applications

in both cases, we are interested in modelling

$$p(y|x)$$

CT Scan

Segmentation Samples

Initial Conditions

Observed Galaxy Distribution

forward model

$x = f(y) + \text{noise}$

$x$

$y$

$y$

$x$ $x$

25

# Prob. U-Net

- Combination of cVAE & U-Net

- Latent spaces at several "scales" of the expanding path



**Credit:** Simon Kohl et al. (2019) – Figure taken with modification.

# Prob. U-Net

- Prior "conditioned" on the observation and latents of previous scales

$$z_i \sim p(z_i | z_{<i}, x)$$

- Joint prior decomposes into priors of each scale

$$p(z_0, \ldots, z_L \mid x) = p(z_L \mid z_{<L}, x) \cdot \ldots \cdot p(z_0 \mid x)$$



Observation

Predictions

**Prior Net**

**Question:** Which part(s) resemble the VAE component that models

**prior** / **likelihood** / **variational posterior**?

27

# Prob. U-Net



**Used in Training & Inference**

$$\boldsymbol{z}_i \sim p(\boldsymbol{z}_i | \boldsymbol{z}_{<i}, \boldsymbol{x})$$

$$p(\boldsymbol{z}_0, ..., \boldsymbol{z}_L \mid \boldsymbol{x}) = p(\boldsymbol{z}_L | \boldsymbol{z}_{<L}, \boldsymbol{x}) \cdot ... \cdot p(\boldsymbol{z}_0 | \boldsymbol{x})$$

**Used in Training**

$$\boldsymbol{z}_i \sim q(\boldsymbol{z}_i | \boldsymbol{z}_{<i}, \boldsymbol{x}, \boldsymbol{y})$$

$$q(\boldsymbol{z}_0, ..., \boldsymbol{z}_L \mid \boldsymbol{x}, \boldsymbol{y}) = q(\boldsymbol{z}_L | \boldsymbol{z}_{<L}, \boldsymbol{x}, \boldsymbol{y}) \cdot ... \cdot q(\boldsymbol{z}_0 | \boldsymbol{x}, \boldsymbol{y})$$

# Prob. U-Net



Prior Net



Latents are pixelwise Gaussians

Posterior Net

Posterior Net has a "truncated" decoder

# Training

- Means and STDs predicted using both networks

  - Used to calculate KL

- Samples drawn from Posterior Net latents and inserted into the Prior Net

- **Objective:** Maximize evidence

$$\mathcal{L}_{\text{ELBO}} = \mathbb{E}_{z \sim Q}[-\ln p(y \mid x, z)]$$

$$+ \sum_{i=0}^{L} D_{\text{KL}}\big(q_i(z_i \mid z_{<i}, x, y) \parallel p_i(z_i \mid z_{<i}, x)\big)$$

$$\mathcal{L}_{\text{ELBO}} = \mathcal{L}_{\text{rec}} + \mathcal{L}_{\text{KL}}$$

# Toy Problem 1

- **Input:** Observations of Lensed Galaxies

- **Goal:** Reconstruct the Source Galaxy



lensed image seen of
background galaxy

background galaxy

foreground galaxy

Looking further into the past

**Credit:** ALMA (ESO/NRAO/NAOJ), Luis Calçada (ESO), Yashar Hezaveh et al.

# Toy Problem 1

**Input:** Observation of a Lens-Source System



**Goal:** Find the Undistorted Image of the Source Galaxy



- Different ways to lens the source galaxy → Problem is underconstrained

- *More Precise* **Goal:** Draw samples from the posterior distribution of reconstructed source images

#source_reconstruction

#posterior_sampling

# A Subtle Difference!

**Variational Posterior**

Defined in Latent Space

$$p(z|x, y)$$

<span style="color:blue">we mean this when we say **posterior network**!</span>

**Parameters Posterior**

Defined in Parameter Space

$$p(y|x)$$

<span style="color:purple">we mean this when we say **posterior sampling**!</span>

$x \in$ data space    $y \in$ parameter space    $z \in$ latent space

# Part 2
# **Rescue the Randomness**

**KL Vanishing Problem**

**ELBO Loss with $\beta$**

**GECO Loss**

**Toy Problem 2: One-hot Flipping**

# KL Vanishing Problem

Training

Validation



- KL Term vanishes early in the training due to non-informative latents

- Model ignores the cause of probabilistic behavior by setting respective weights to 0 → Deterministic

# KL Vanishing Problem

- KL Term vanishes early in the training due to non-informative latents

- Model ignores the cause of probabilistic behavior by setting respective weights to 0 → Deterministic

- Happens when two types of paths exist:

  A. **Latent Path:** Conditioned on the latent space (same as VAEs)

  B. **Leaky Path:** Does not pass through latent spaces;
  Leaks the ground truth information



Observation +
Ground Truth

# ELBO with $\beta$

- **Idea:** Prevent the optimization scheme from caring too much about the KL term before having meaningful latents.

- Possible Approaches:

  - Set $0 < \beta < 1$

  - Start with $\beta = 0$ and Gradually Increase it (**Beta Annealing**)

  - Other ways of scheduling $\beta$ (e.g., **Cyclical Schedule**)

- What is the best way to schedule $\beta$?

  - Variety of choices

  - Depends on the specific problem

$$\mathcal{L}_{\mathrm{ELBO}} = \mathcal{L}_{\mathrm{rec}} + \beta \, \mathcal{L}_{\mathrm{KL}}$$

governs the amount of regularization

37

# GECO

- **G**eneralized **E**LBO with **C**onstrained **O**ptimization

- Constrained Optimization Framework

  - Minimize the KL Term under a set of reconstruction constraints

- $\lambda$ plays the role of $\beta \rightarrow$ automatically updated during training

  - (Usually) tend to focus on the reconstruction loss early in the training until it reaches $\kappa$;

  - Then moves the pressure over on the KL Term.

- Advantages:

  - Hyperparameter ($\kappa$) defined in data space $\rightarrow$ More intuitive

  - $\beta$ is updated automatically

$$\mathcal{L}_{\mathrm{ELBO}} = \mathcal{L}_{\mathrm{rec}} + \beta \, \mathcal{L}_{\mathrm{KL}}$$

$$\mathcal{L}_{\mathrm{GECO}} = \lambda \, (\mathcal{L}_{\mathrm{rec}} - \kappa) + \mathcal{L}_{\mathrm{KL}}$$

Lagrange multiplier

reconstruction threshold

$$\lambda \equiv \frac{1}{\beta}$$

38

# Toy Problem 2

- Training Set: all 32-bit **one-hot** vectors

**Input**



**Ground Truth**

Flip All Sequences

Global Uncertainty
Add Noise $\sim \mathcal{N}(0, \sigma)$

Local Uncertainty
Roll Each Sequence
by $r$ Pixels

$p(r = 0) = 0.4$
$p(+1) = p(-1) = 0.2$
$p(+2) = p(-2) = 0.1$

- New realizations generated at each training step

# Visualizing Latents

- Assign a unique color to each input

- Sample a bunch of latent representations for each input
  - For an arbitrary scale of the Prob. U-Net
  - Can have many dimensions

- Plot the first two **principal components** (orthogonal directions with most variability)



$PC_2$

$PC_1$

"Advertisement"
**Inspecting Uncertainties**

Coverage Probability Test

# Don't Get Too Excited!

- Having a model with probabilistic behavior is not enough!

- Require comprehensive statistical analysis that goes beyond the model's assumptions

- To make sure uncertainties are appropriately quantified



DALL·E's impression of

**A Robot Thinking About Statistics**

# Coverage Probability Test

**High-level explanation:**

1. Repeatedly sample from the model

2. Calculate a confidence interval using samples (expected coverage)

3. Check if the true value falls within the interval

4. Repeat steps 1-3 for multiple "samples - true value" combinations

5. Calculate the fraction of times that the true value falls within each confidence interval (true coverage)

6. Plot true coverage vs. expected coverage



**Credit:** Pablo Lemos et al. (2022)

# 📢 Advertisement

## Sampling-Based Accuracy Testing of Posterior Estimators for General Inference

Pablo Lemos [1,2,3,4,*]   Adam Coogan [1,2,3,*]   Yashar Hezaveh [1,2,3]   Laurence Perreault-Levasseur [1,2,3]

**arXiv:** 2302.03026

- Method to estimate coverage probabilities of generative posterior estimators without posterior evaluations (by just using samples).

- Necessary and sufficient to show that a posterior estimator is optimal.

- `pip`-installable package on the way!

# Thank You For Your Attention!

**Contact:** hadi.sotoudeh@umontreal.ca

**Collaborators:** Laurence Perreault-Levasseur, Pablo Lemos, Ève Campeau-Poirier, Charles Wilson, Alexandre Adam

This work was made possible through the generous support of:

Université de Montréal · Mila · Ciela Institute · SIMONS FOUNDATION · Digital Research Alliance of Canada · IVADO

# To Read More…

- **Prob. U-Net:** Kohl, S. A. A., "A Probabilistic U-Net for Segmentation of Ambiguous Images", arXiv: 1806.05034 🔗

- **Hierarchical Prob. U-Net:** Kohl, S. A. A., "A Hierarchical Probabilistic U-Net for Modeling Multi-Scale Ambiguities", arXiv: 1905.13077 🔗

- **KL Vanishing and Cyclical $\beta$:** Fu, H., Li, C., Liu, X., Gao, J., Celikyilmaz, A., and Carin, L., "Cyclical Annealing Schedule: A Simple Approach to Mitigating KL Vanishing", arXiv: 1903.10145 🔗

- **GECO:** Jimenez Rezende, D. and Viola, F., "Taming VAEs", arXiv: 1810.00597 🔗

- **Coverage Test:** Lemos, P., Coogan, A., Hezaveh, Y., and Perreault-Levasseur, L., "Sampling-Based Accuracy Testing of Posterior Estimators for General Inference", arXiv: 2302.03026 🔗

- **VAEs:** Rocca, J., Blog Post on "Understanding Variational Autoencoders (VAEs)", 🔗

- **Conditional VAEs:** Dykeman, I., Blog Post on "Conditional Variational Autoencoders", 🔗