# Learning with Gaps

A Domain-Adaptive SBI Framework for Mapping Young Stars from Incomplete, Multi-Survey Data
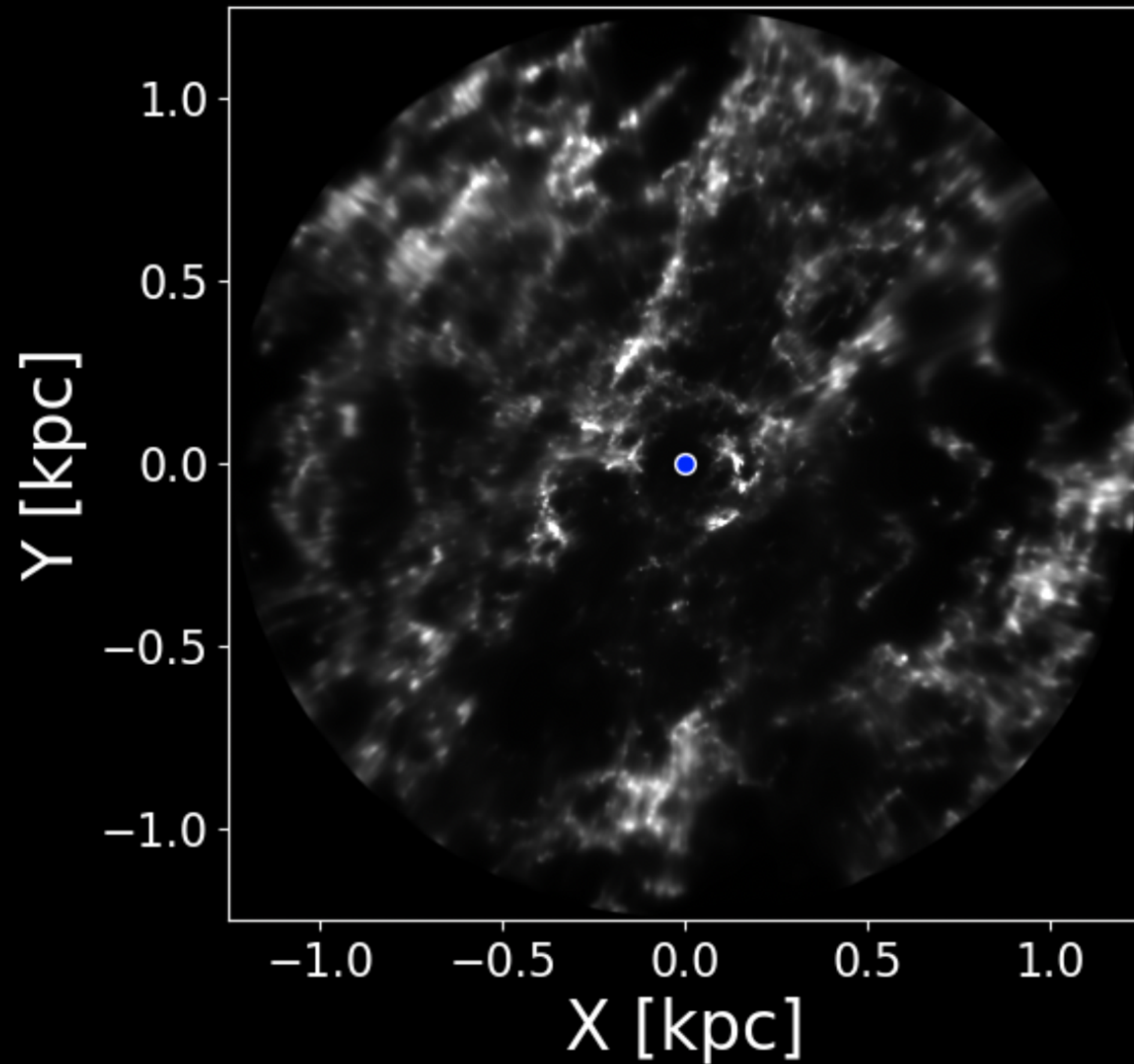
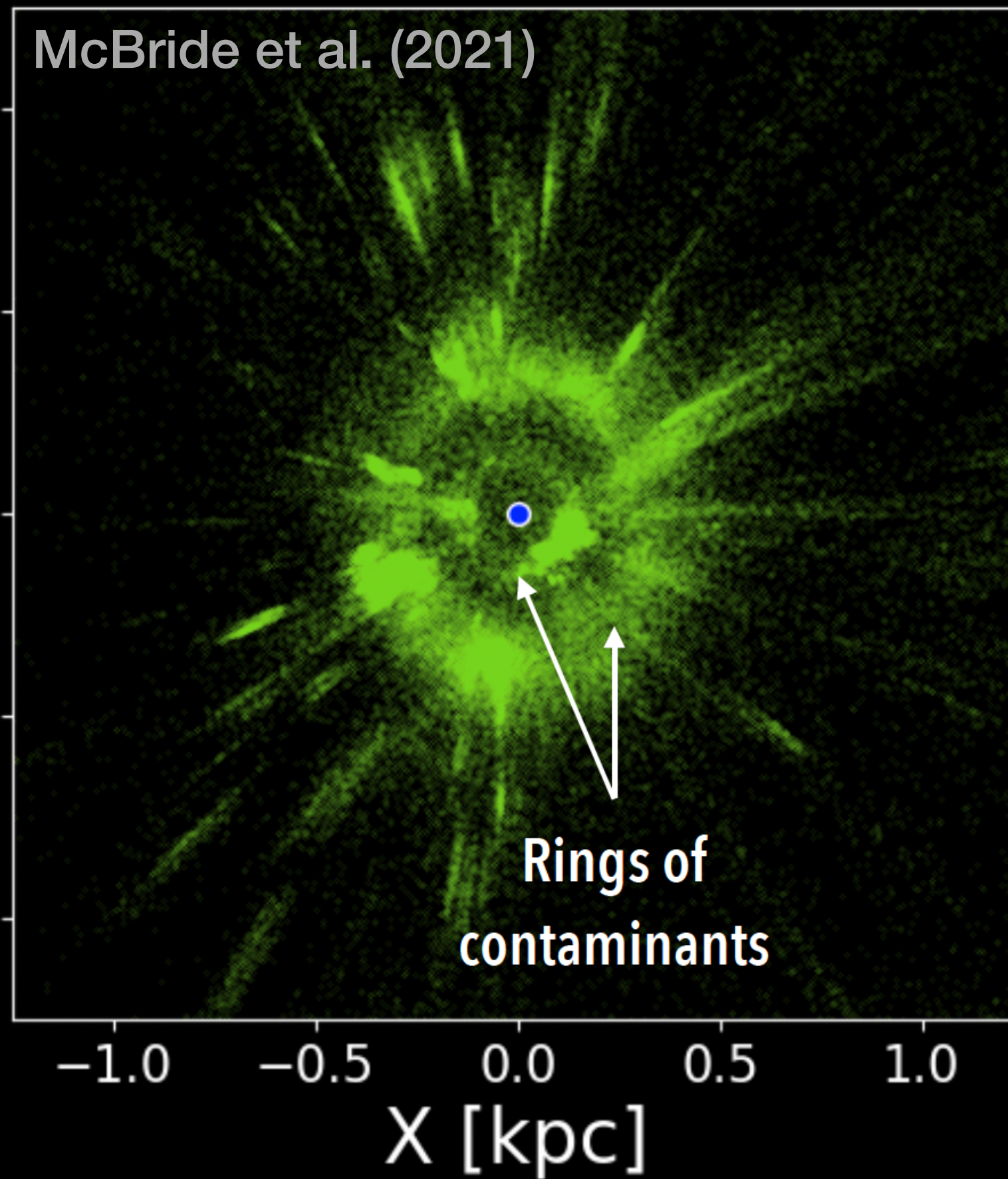**Sebastian Ratzenböck @CfA**

In collaboration with
*Catherine Zucker (CfA), Joshua Speagle (UToronto), Phillip Cargile (CfA), Philipp Frank (Stanford), Andrew Saydjari (Princeton)*

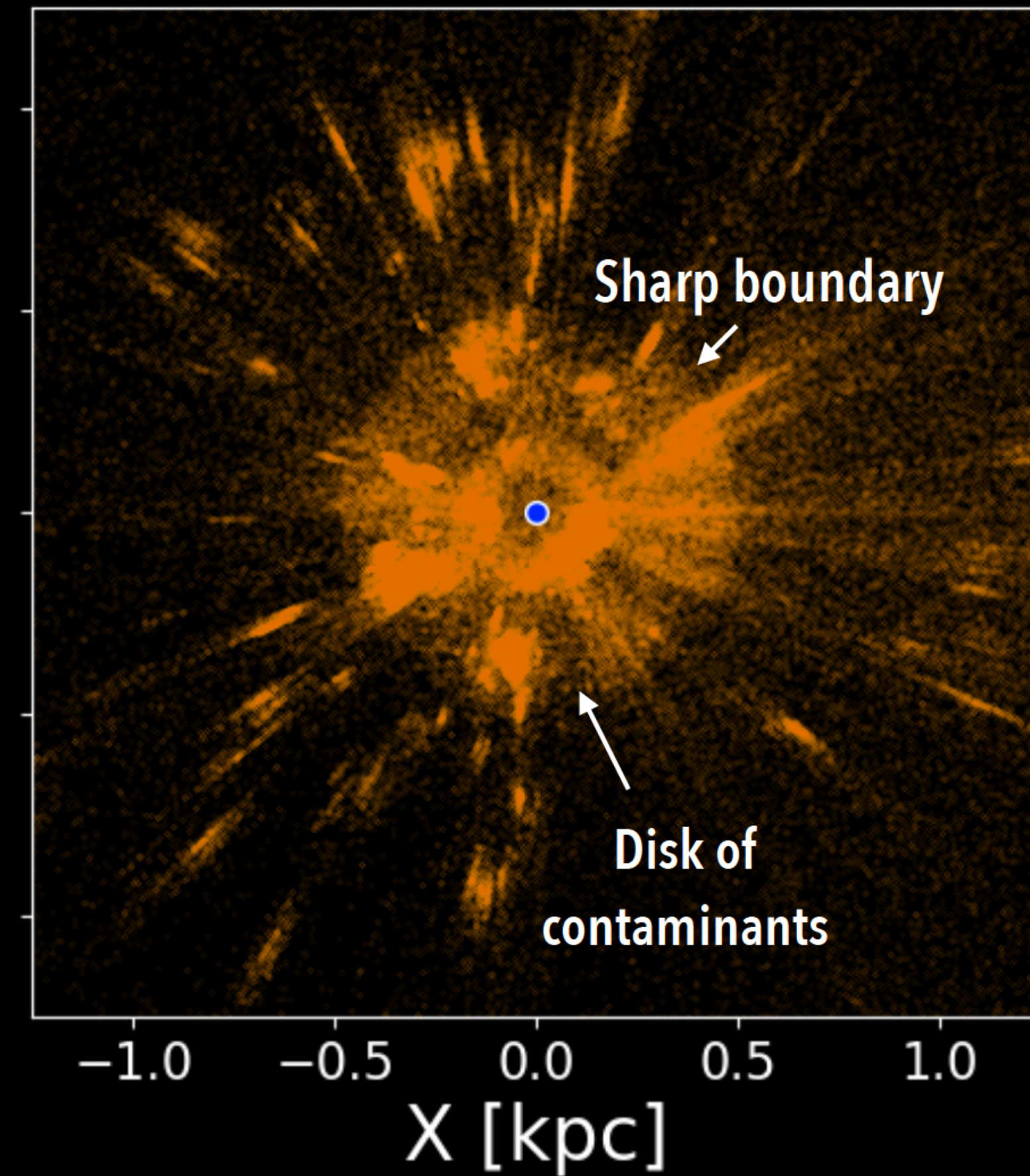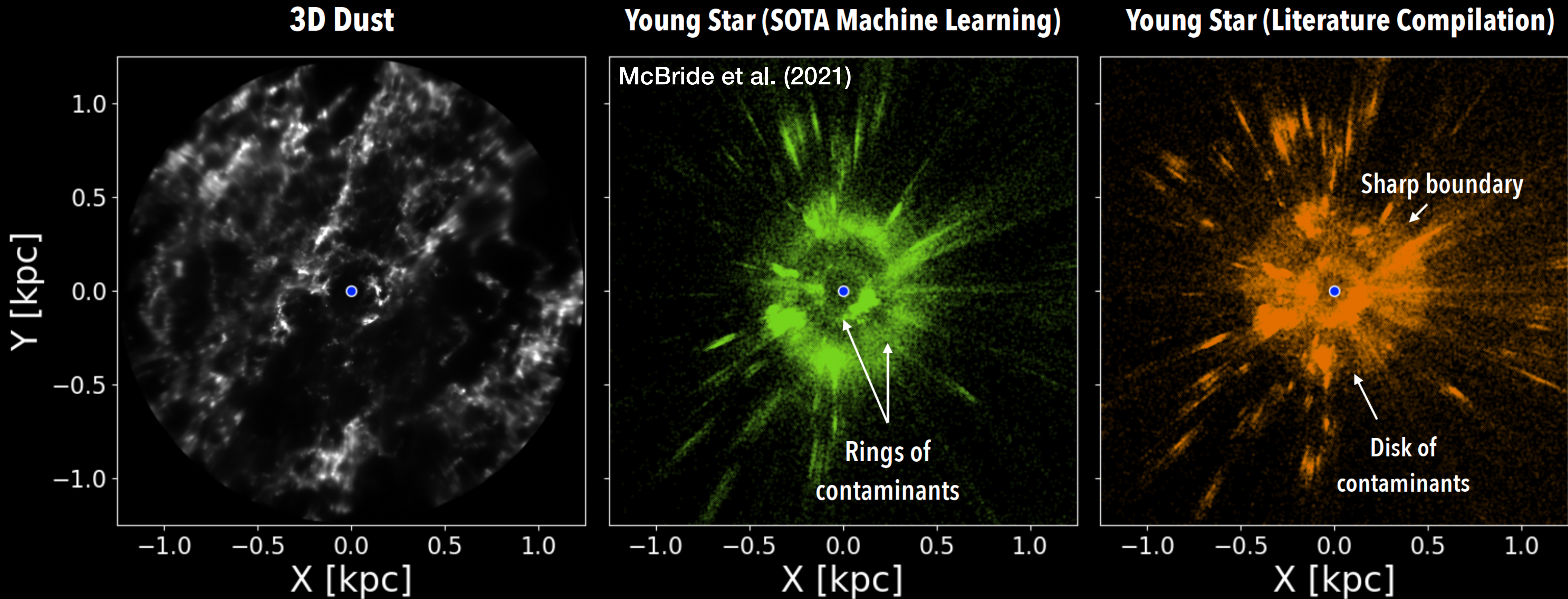# YSOs: Critical link to understanding Galactic baryon cycle



**3D Dust**

**Young Star (SOTA Machine Learning)**

McBride et al. (2021)

Rings of contaminants

**Young Star (Literature Compilation)**

Sharp boundary

Disk of contaminants

# YSOs: Critical link to understanding Galactic baryon cycle



**3D Dust**

**Young Star (SOTA Machine Learning)**

McBride et al. (2021)

Rings of contaminants

**Young Star (Literature Compilation)**
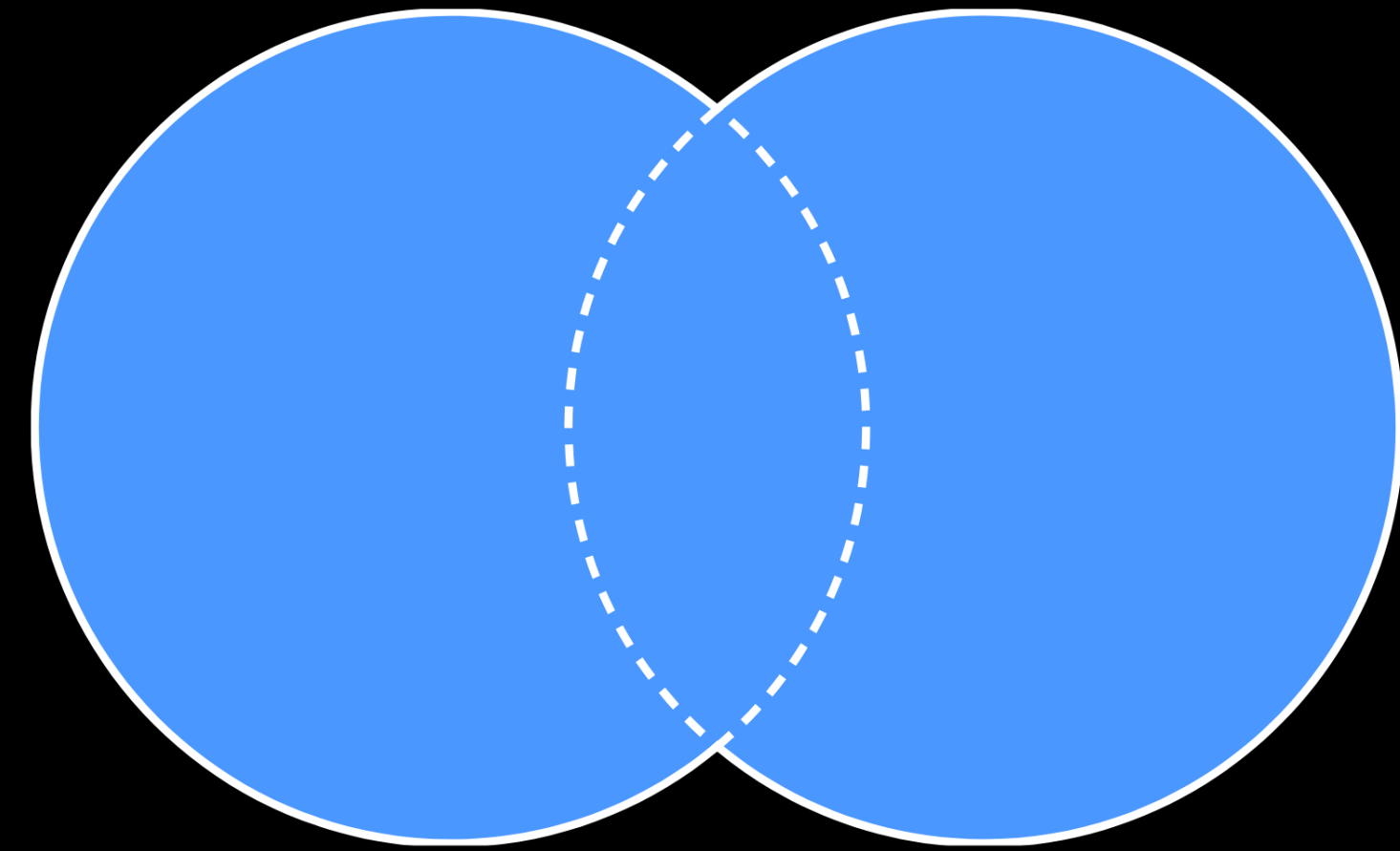
Sharp boundary

Disk of contaminants

— aim to improve this

# Aim to improve YSO catalog

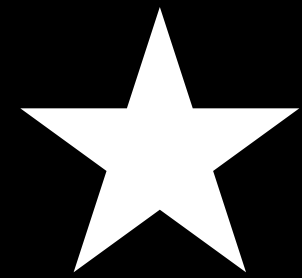- Data fusion: use as many informative data sets as possible



Often works focus
on intersection

Aim to be truely multi-survey

# Aim to improve YSO catalog

- Data fusion: use as many informative data sets as possible

★

Gaia
2MASS
WISE
LAMOST

★

WISE
Spitzer
APOGEE

# Aim to improve YSO catalog

- Data fusion: use as many informative data sets as possible

- Produce well-calibrated posteriors over stellar parameters given spectra & photometric observations

# Aim to improve YSO catalog

- Data fusion: use as many informative data sets as possible

- Produce well-calibrated posteriors over stellar parameters given spectra & photometric observations

- Scale inference to > 1M - 1B stars

# Challenges with "1 model does it all" approach

- Fusing surveys is hard due to different
  - resolutions & depths
  - coverage
  - instrument response
  - noise model
  - …

# Challenges with "1 model does it all" approach

- Fusing surveys is hard due to different
  - resolutions & depths
  - coverage
  - instrument response
  - noise model

- **Model misspecification** leads to domain shift between simulated and real data

# Challenges with "1 model does it all" approach

- Fusing surveys is hard due to different
  - resolutions & depths
  - coverage
  - instrument response
  - noise model

- **Model misspecification** leads to domain shift between simulated and real data

→ *Domain-Adaptive SBI w/ incomplete, multi-survey data*

# Model implementation
## I. SBI model

# Typical ML regression

$$\vec{x} \longrightarrow \text{MLP} \longrightarrow \vec{\theta}$$

MLP… Series of **_learnable_** affine transformations of $\vec{x}$

followed by pointwise non-linear map: $f_\phi(\vec{x}) = \hat{\theta}$

$\phi$…learable paramerters

# Typical ML regression

$$\vec{x} \quad \longrightarrow \quad \text{MLP} \quad \longrightarrow \quad \vec{\theta}$$

MLP… Series of **_learnable_** affine transformations of $\vec{x}$

followed by pointwise non-linear map: $f_\phi(\vec{x}) = \hat{\theta}$

Trained by minizing $||\vec{\theta} - \hat{\theta}||_2$

# Typical ML regression

$$\vec{x} \longrightarrow \text{MLP} \longrightarrow \vec{\theta}$$

Would like to have $p(\vec{\theta} \mid \vec{x})$

# Typical ML regression

$$\vec{x} \longrightarrow \text{MLP} \longrightarrow \vec{\theta}$$
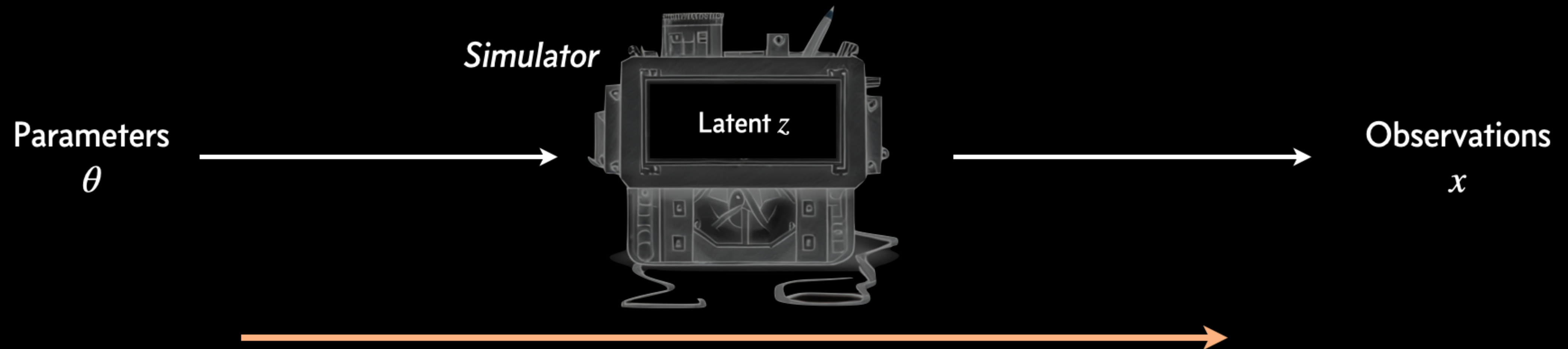
Would like to have $p(\vec{\theta} \mid \vec{x})$

However:  — $p(\vec{x} \mid \vec{\theta})$ might not be tractable

— $p(\vec{\theta} \mid \vec{x})$ might not scale to millions - billions of "runs"

# Typical ML regression

$$\vec{x} \longrightarrow \text{MLP} \longrightarrow \vec{\theta}$$

Would like to have $p(\vec{\theta} \mid \vec{x})$

However:  — $p(\vec{x} \mid \vec{\theta})$ might not be tractable

— $p(\vec{\theta} \mid \vec{x})$ might not scale to millions - billions of "runs"

BUT: if we have access to a simulator, we can approximate $p(\vec{\theta} \mid \vec{x})$

# Simulation based inference (SBI) setup



Simulator

Latent $z$

Parameters
$\theta$

Observations
$x$

**Prediction:**
- Mechanistic forward model
- We can generate samples from a simulator $x \sim p(x \,|\, \theta)$

**Inference:**
- Likelihood $p(x \,|\, \theta) = \int \mathrm{d}z \, p(x, z \,|\, \theta)$ is intractable
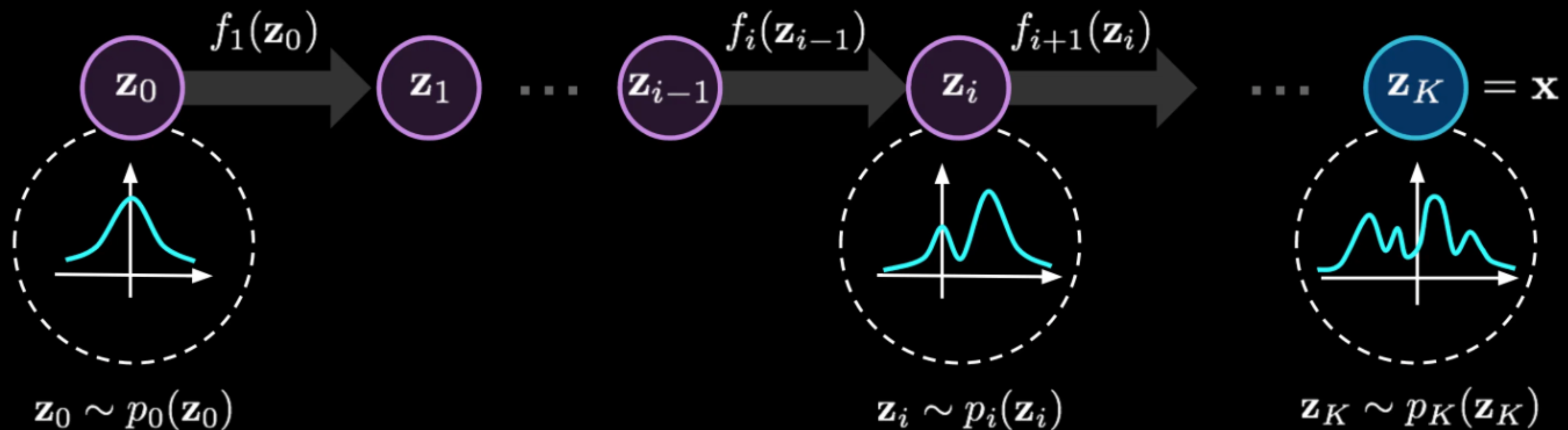- *Inference is challenging*
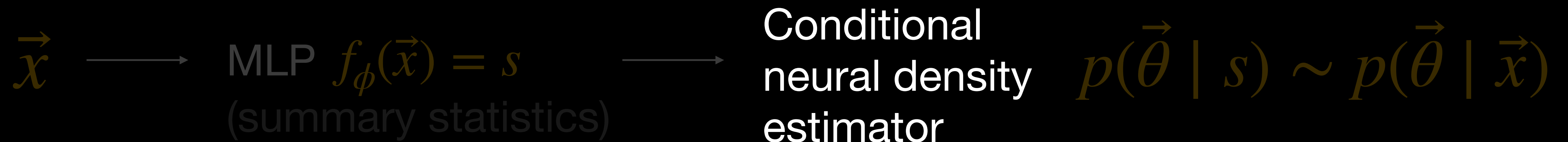
# Neural posterior estimation

$$\vec{x} \longrightarrow \text{MLP} \longrightarrow \vec{\theta}$$

ML regression

# Neural posterior estimation

$$\vec{x} \longrightarrow \text{MLP} \longrightarrow \vec{\theta}$$

ML regression

$$\vec{x} \longrightarrow \text{MLP } f_\phi(\vec{x}) = s \longrightarrow \text{Conditional neural density estimator} \quad p(\vec{\theta} \mid s) \sim p(\vec{\theta} \mid \vec{x})$$

(summary statistics)

Papamakarios & Murray (2016)  & Lueckmann et al. (2017)

# Normalizing flows

$\vec{x} \longrightarrow$ MLP $f_\phi(\vec{x}) = s$ $\longrightarrow$ Conditional neural density estimator $p(\vec{\theta} \mid s) \sim p(\vec{\theta} \mid \vec{x})$
(summary statistics)



Parameterized, invertible maps $f_i$ that transform Gaussian into target distribution

Training objective: ***<u>maximum likelihood</u>***

# Neural posterior estimation

$$\vec{x} \longrightarrow \text{MLP } f_\phi(\vec{x}) = s \longrightarrow \text{Conditional neural density estimator} \quad p(\vec{\theta} \mid s) \sim p(\vec{\theta} \mid \vec{x})$$

(summary statistics)

Cannot deal with missing data

# Transformer: learning with incomplete data

$$\vec{x} \longrightarrow \text{Transformer}$$

$$\uparrow$$

Attention masking: $M_E$

Can enforce conditional independence

— Effectively marginalize over missing values

Gloeckler et al. (2024)

# Transformer: learning with incomplete data

$$\vec{x} \longrightarrow \boxed{\text{Transformer}} \longrightarrow \text{Conditional density estimator (flow matching)} \quad p(\vec{\theta} \mid \vec{x})$$

Attention masking: $M_E$

# Transformer: learning with incomplete data

$$\vec{x} \longrightarrow \text{Transformer} \longrightarrow \text{Conditional density estimator (flow matching)} \quad p(\vec{\theta} \mid \vec{x})$$

Attention masking: $M_E$

Gaia
2MASS
WISE
LAMOST

WISE
Spitzer
APOGEE

# Model implementation
## II. Dealing with model misspecification

# Input split into *simulated*, *real & paired* data

$$x \longrightarrow \begin{matrix} x_{\mathrm{sim}} \\ x_{\mathrm{real}} \\ x_{\mathrm{sim-real-pairs}} \end{matrix}$$

# Modality encoders: split into indiv. spectra

$x_{\text{sim}}$

Apogee

Boss

Lamost

XP


Photometry

Gaia, 2Mass, …

# Modality encoders: encode

$x_{\mathrm{sim}}$

Apogee     *MLP encoder*

$z_{\mathrm{s}}^{\mathrm{A}} = f_{\phi}(x_{\mathrm{s}}^{\mathrm{A}}) \;\; \hat{=}$ summary statistics

Boss               $z_{\mathrm{s}}^{\mathrm{B}}$

Lamost           $z_{\mathrm{s}}^{\mathrm{L}}$

XP                  $z_{\mathrm{s}}^{\mathrm{xp}}$

Photometry

Gaia, 2MASS, … $\longrightarrow$ $z_{\mathrm{s}}^{\mathrm{G}}$

$\vdots$                                    $\vdots$

# Modality encoders: alignment loss



Goal: force latent representations of simulated and real data to "look the same"

# Sim-real alignment via optimal transport

$$\vec{z}_{\text{sim}} = f_\phi(\vec{x}_{\text{sim}})$$

$$\vec{z}_{\text{sim}-\text{real}-\text{pairs}}$$



$$\vec{z}_{\text{real}}$$

(a) Samples      (b) Exact      (c) Gromov-Wasserstein      Gu et al. (2022)

Goal: force latent representations of simulated and real data to "look the same"

Gu et al. (2022)

# Let's test this

# Galaxy A … "sim"

- Dust according to Bayestar2019 (Green+2018)

# Galaxy B … "real"

- Dust according to Edenhofer+2024

## Galaxy A … "sim"

Synthetic spectra:

- BaSeL 2.2, ~ Atlas 9 empirically recalibrated (Leujeune+1998)

## Galaxy B … "real"

Synthetic spectra:

- BT-Settl Library (Allard, Hauschildt and Schweitzer 2000)

# Model differences (BaSeL)



logTeff = 3.45 | Av = 0.96 mag

logTeff = 3.47 | Av = 0.97 mag

F_lambda
dusty
GAIA_GAIA3.Gbp
GAIA_GAIA3.G
GAIA_GAIA3.Grp
2MASS_H
2MASS_J
2MASS_Ks
SPITZER_IRAC_36
SPITZER_IRAC_45
SPITZER_IRAC_58
SPITZER_IRAC_80

# Model differences (BTSettl)

# Results

# Pilot study: Gaia+2MASS+WISE

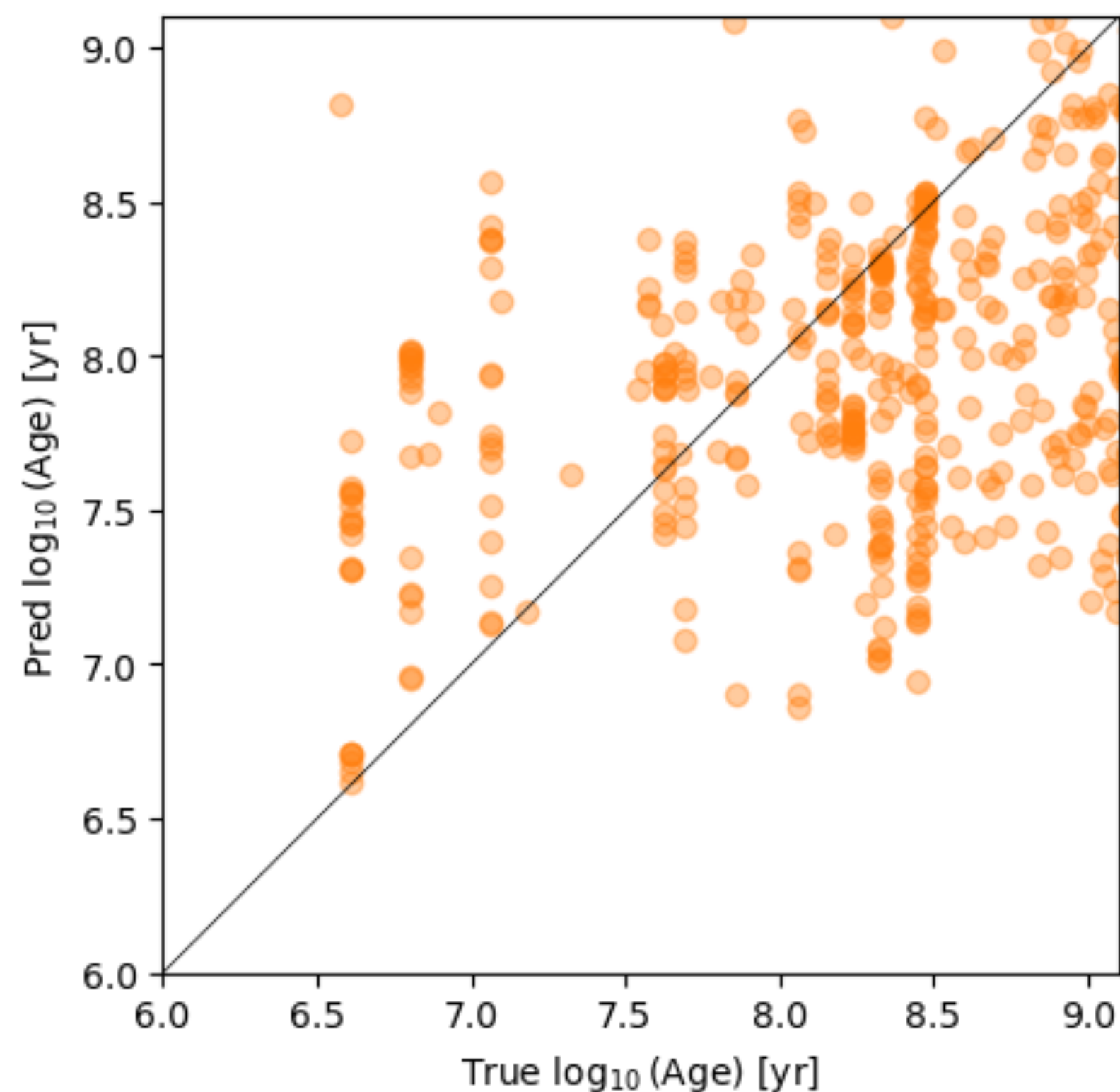# Updated pipeline + XP spectra
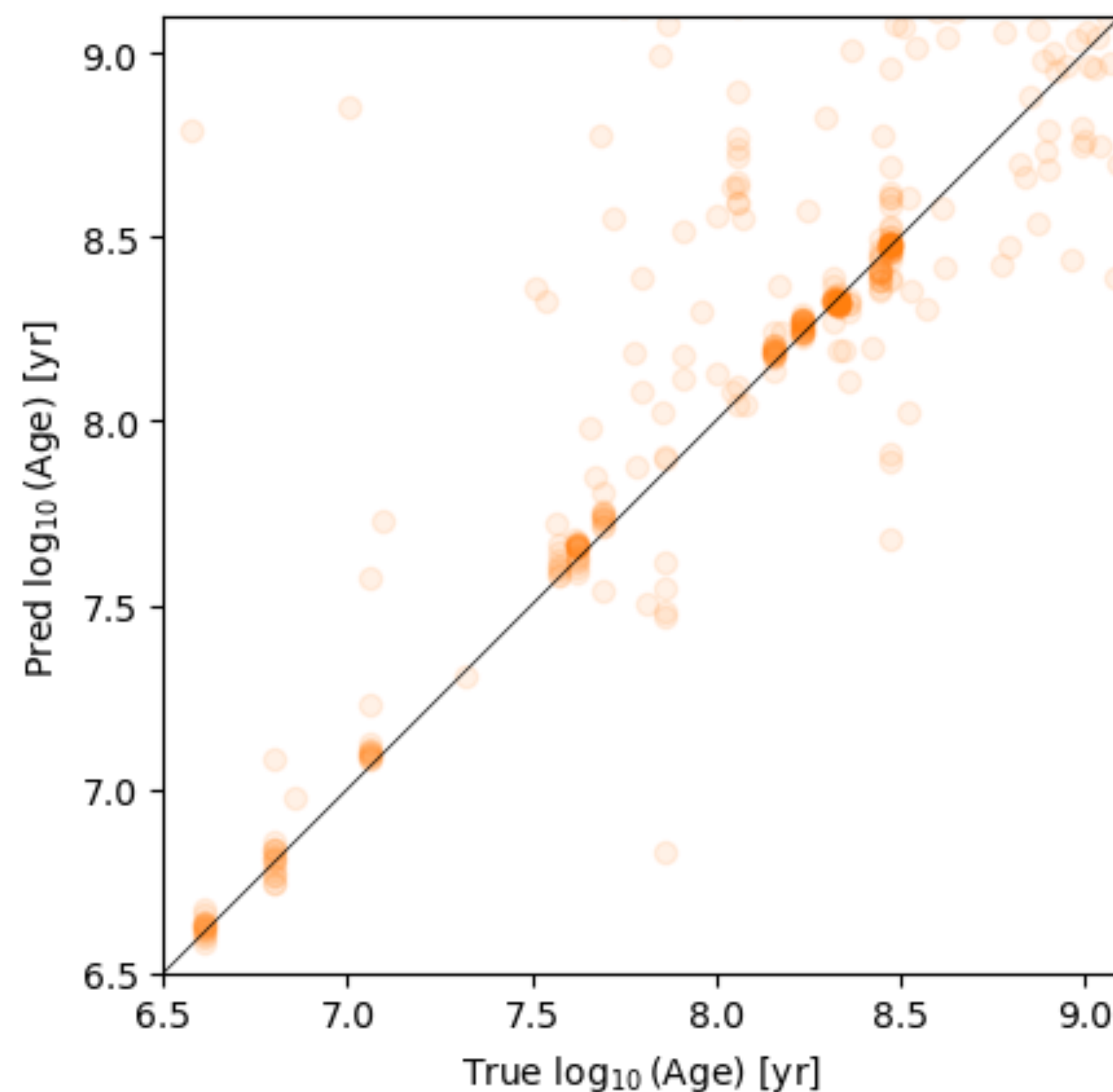


Posterior mean **sim**

# Without DA
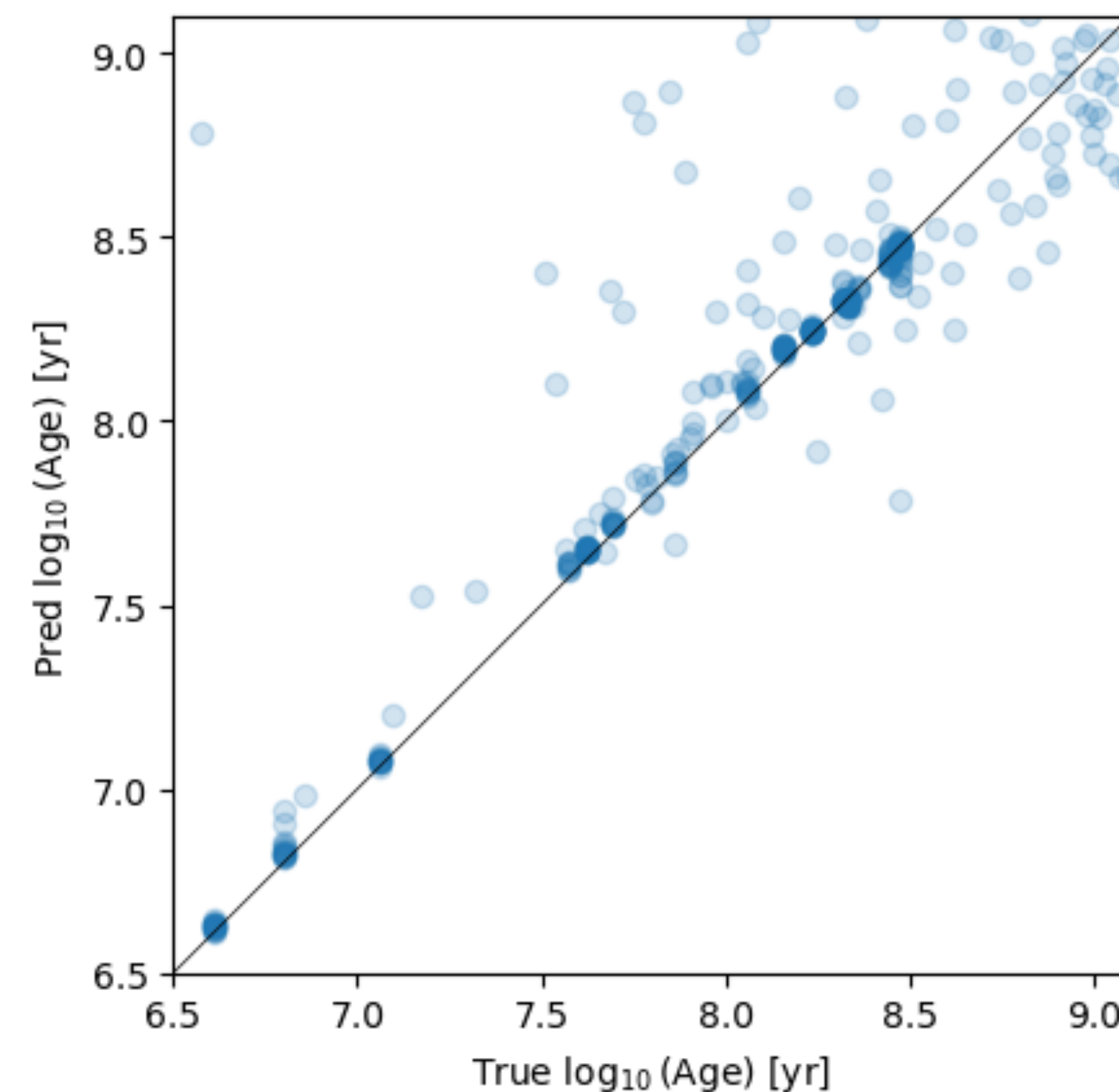
# With DA

# Predictions



Posterior mean "**real**"

Posterior mean "**real**"

Posterior mean **sim**

**no domain adaption**
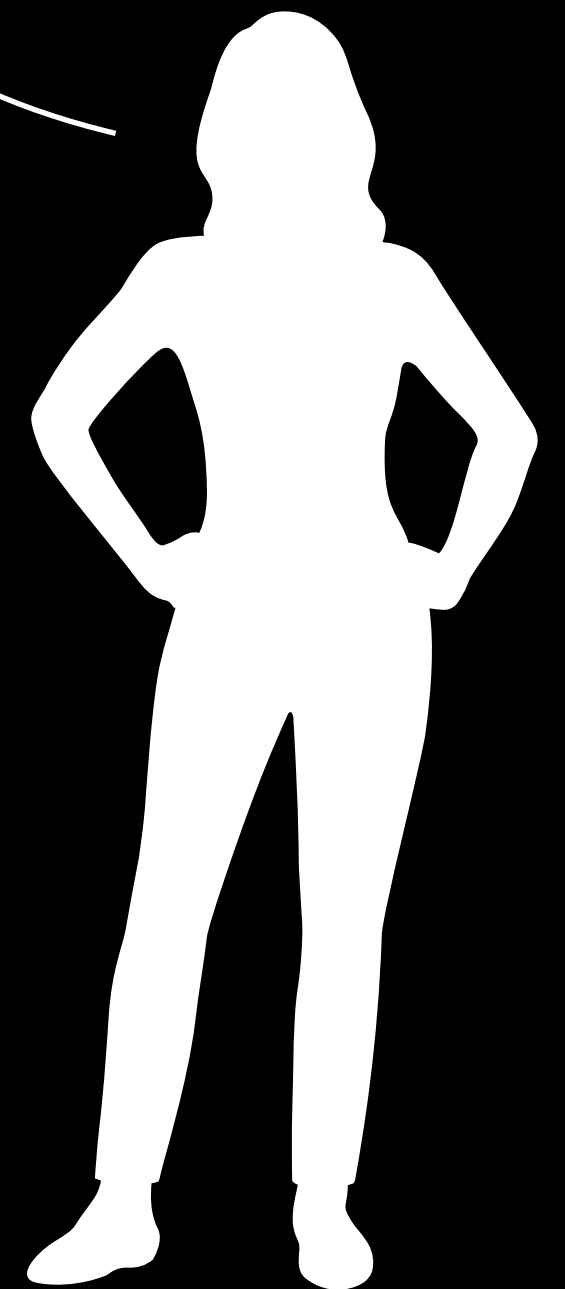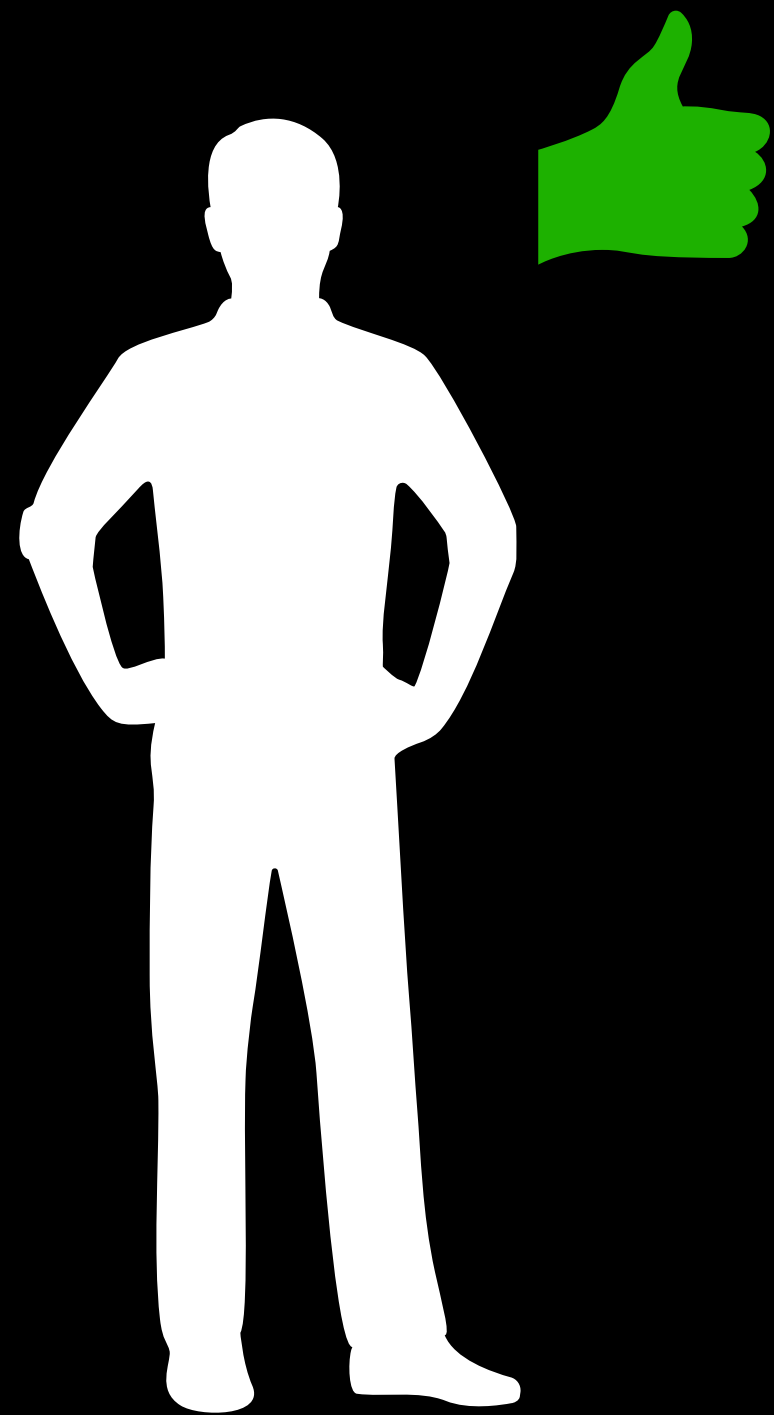
**with domain adaption**

# Summary

- Combine **flow matching models** + **transformer model** to learn arbitrary **conditionals** and **marginals**

- Add OT + pair loss to close domain gap

- Obtain promising results on simulations

# Come find me

I want to know more!

# Come find me

I don't trust ML models!
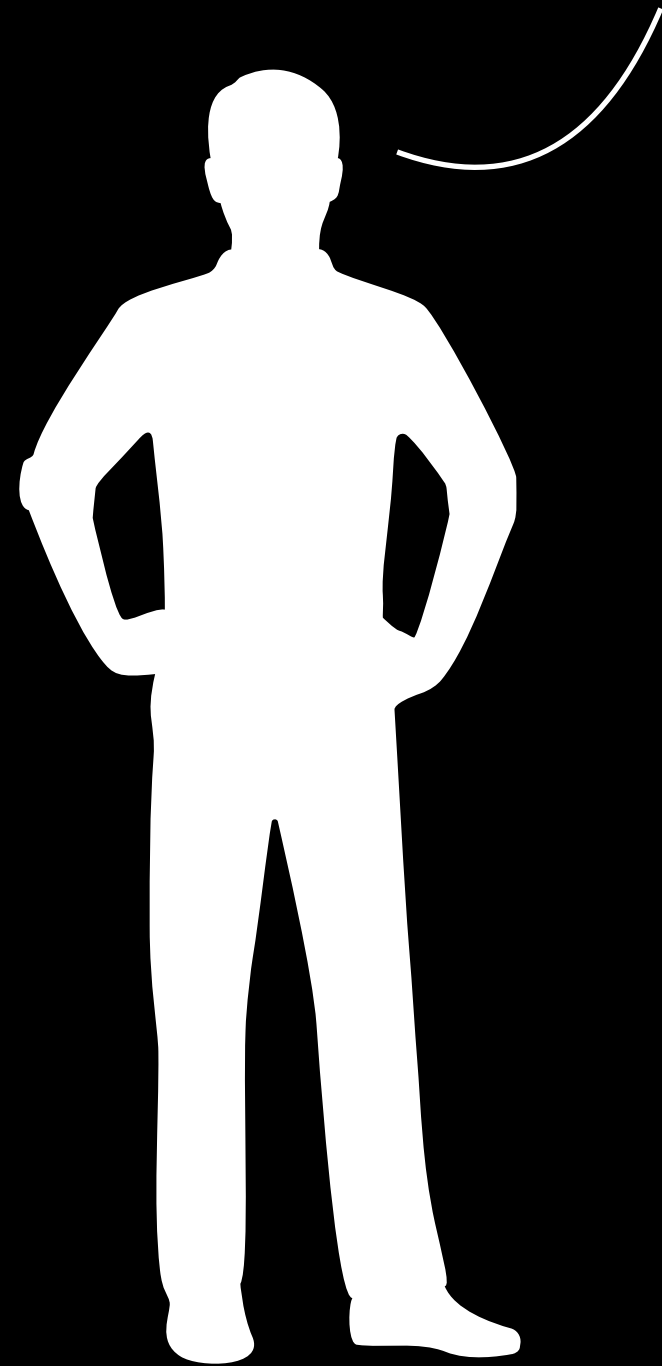
# Come find me

What would you need to see
(on sims) to trust them more?

I don't trust ML models!

Thank you!

# Backup