

# Towards a Complete Census of Young Stars in the local Milky Way

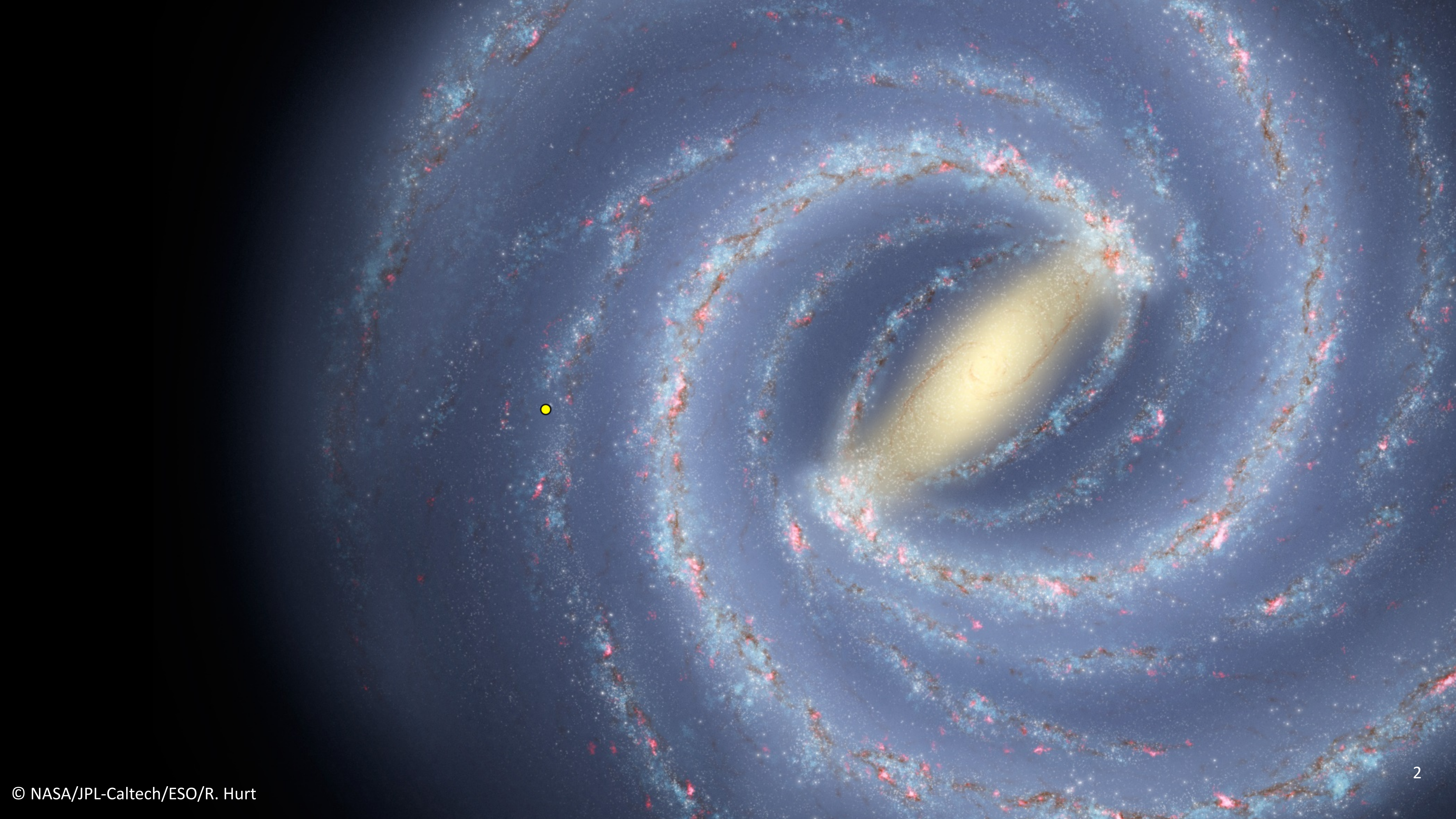
A Self-Consistent Simulation-Based Inference Framework for Incomplete Multi-Survey Data

**Sebastian Ratzenböck @CfA**

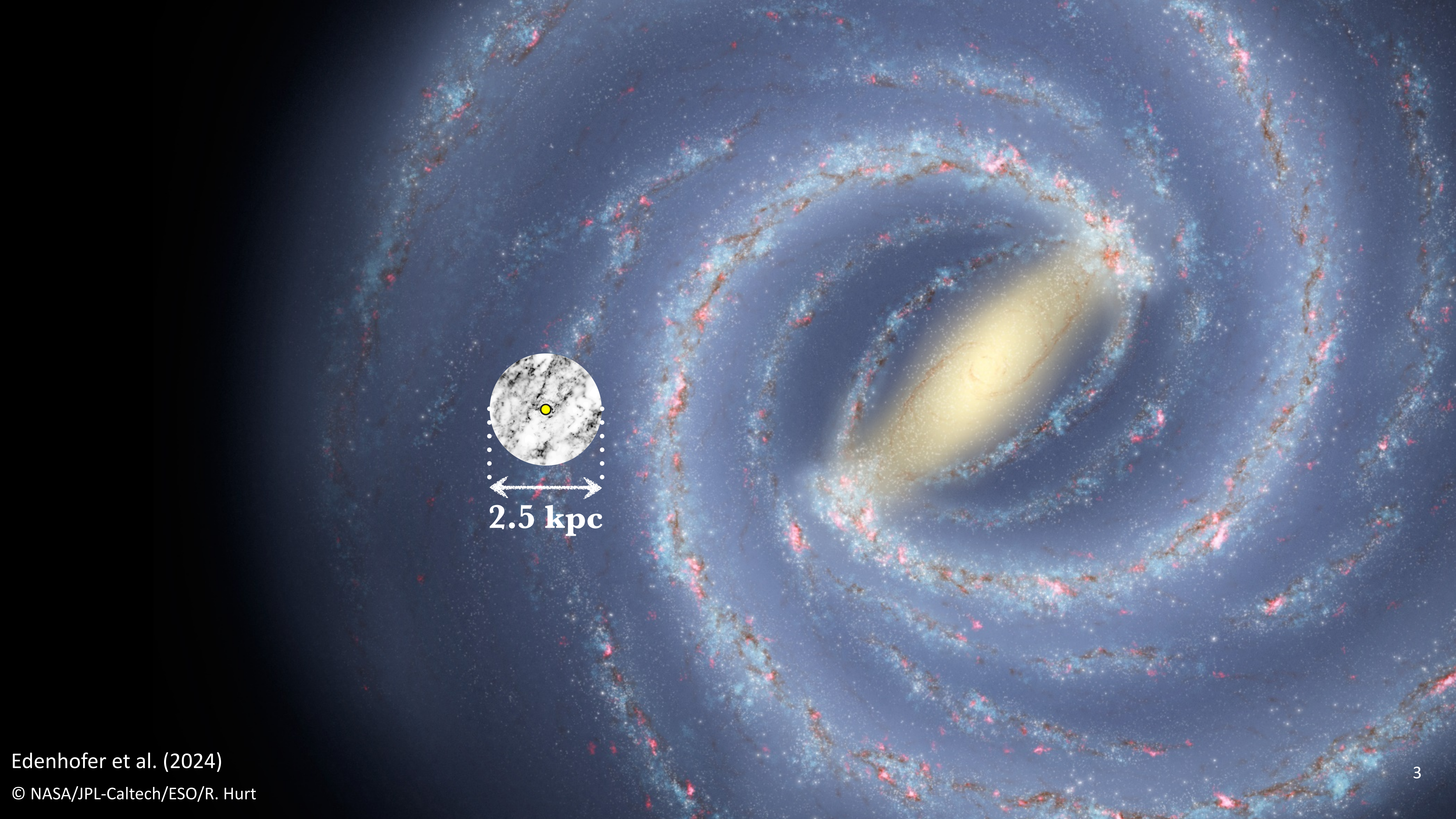
In collaboration with

*Catherine Zucker, Joshua Speagle, Phillip Cargile, Philipp Frank, Andrew Saydjari*







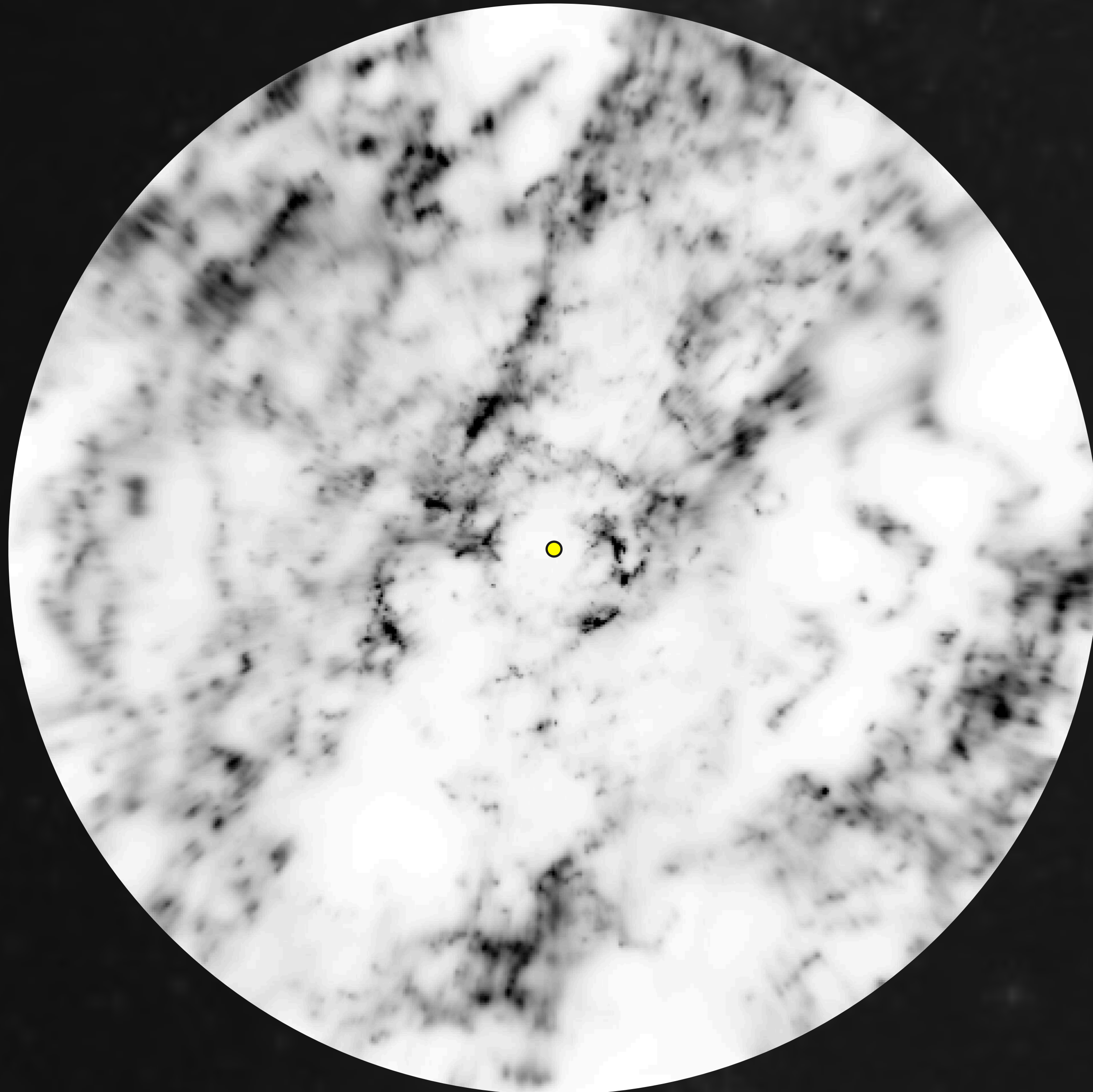




Galactic rotation



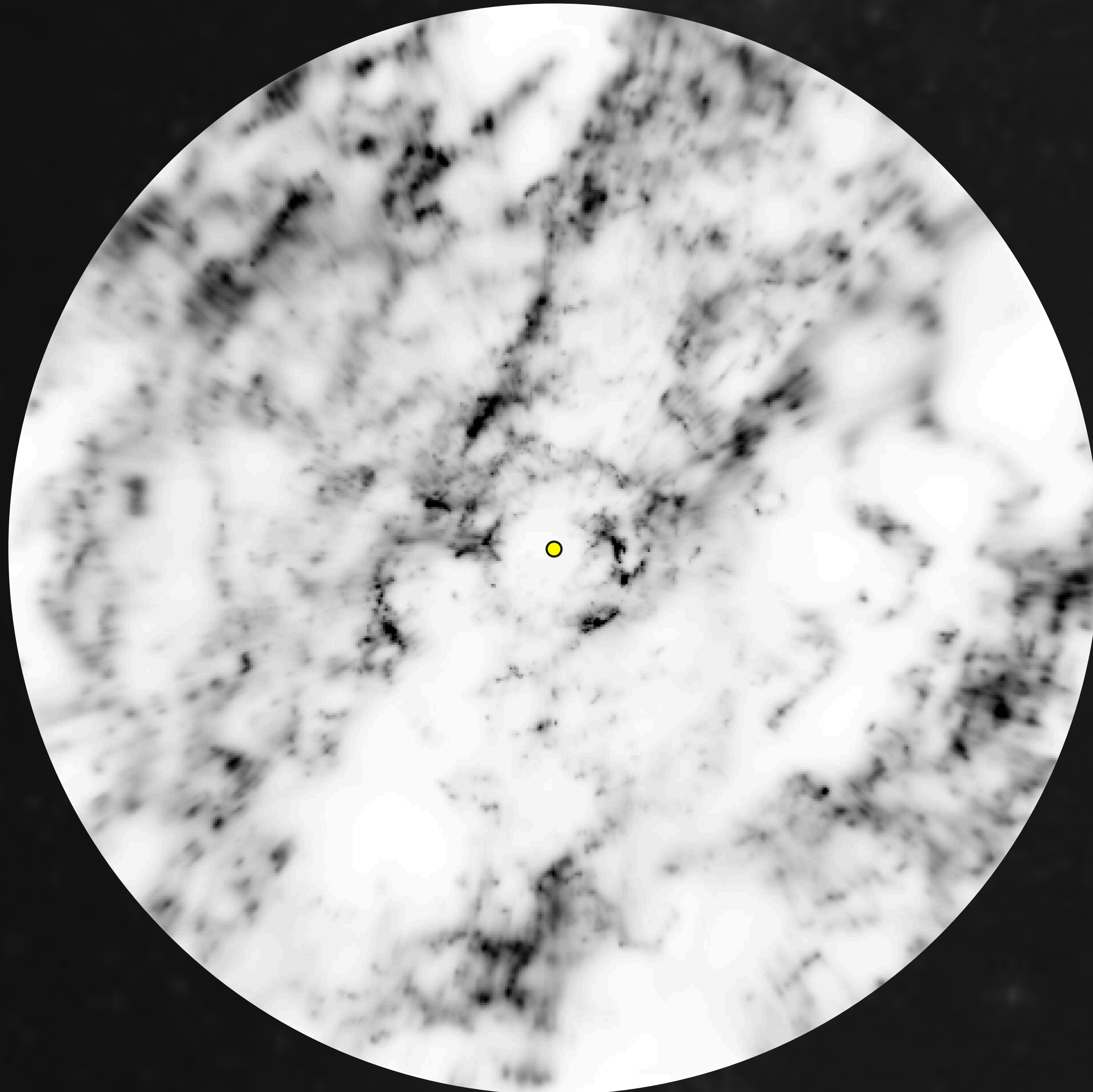
Galactic center



Edenhofer et al. (2024)

© NASA/JPL-Caltech/ESO/R. Hurt





Where are the  
YSOs?



Galactic rotation



Galactic center

Edenhofer et al. (2024)

© NASA/JPL-Caltech/ESO/R. Hurt



# Understanding Galactic baryon cycle

- YSOs connect *cloud*  $\leftrightarrow$  *stars*  $\leftrightarrow$  *feedback*



# Understanding Galactic baryon cycle

- YSOs connect *cloud*  $\leftrightarrow$  *stars*  $\leftrightarrow$  *feedback*
- How does the Milky Way convert gas into stars?



# Understanding Galactic baryon cycle

- YSOs connect *cloud*  $\leftrightarrow$  *stars*  $\leftrightarrow$  *feedback*
- How does the Milky Way convert gas into stars?
- How do stars leave their birth clouds and shape the ISM?



# Understanding Galactic baryon cycle

- YSOs connect *cloud*  $\leftrightarrow$  *stars*  $\leftrightarrow$  *feedback*
- How does the Milky Way convert gas into stars?
- How do stars leave their birth clouds and shape the ISM?
- How do supernovae regulate, trigger, or suppress new generations of stars?



# Understanding Galactic baryon cycle

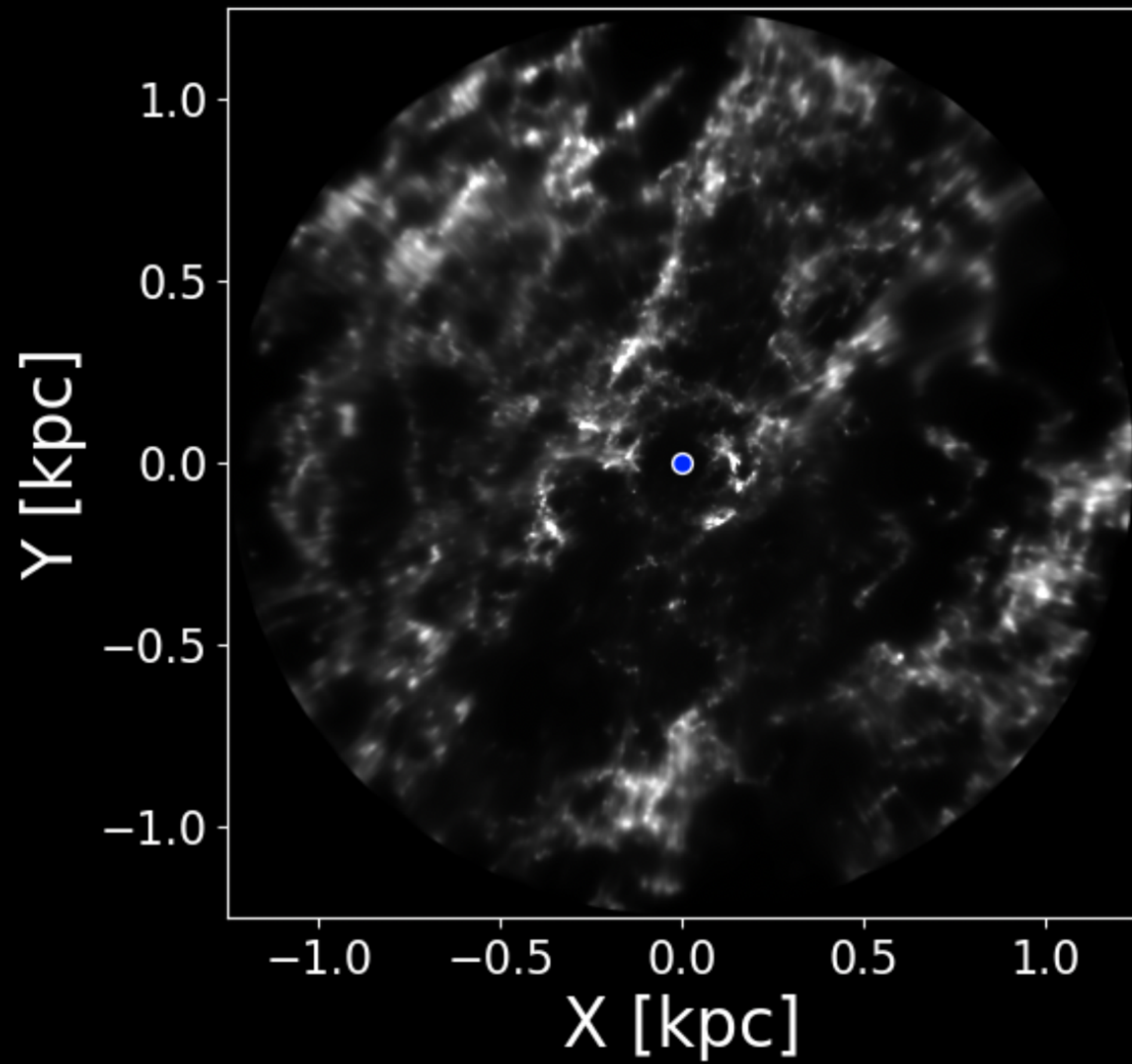
- YSOs connect *cloud*  $\leftrightarrow$  *stars*  $\leftrightarrow$  *feedback*
- How does the Milky Way convert gas into stars?
- How do stars leave their birth clouds and shape the ISM?
- How do supernovae regulate, trigger, or suppress new generations of stars?

**Goal: study YSOs as function of age in relation to dust**



# YSOs in the literature

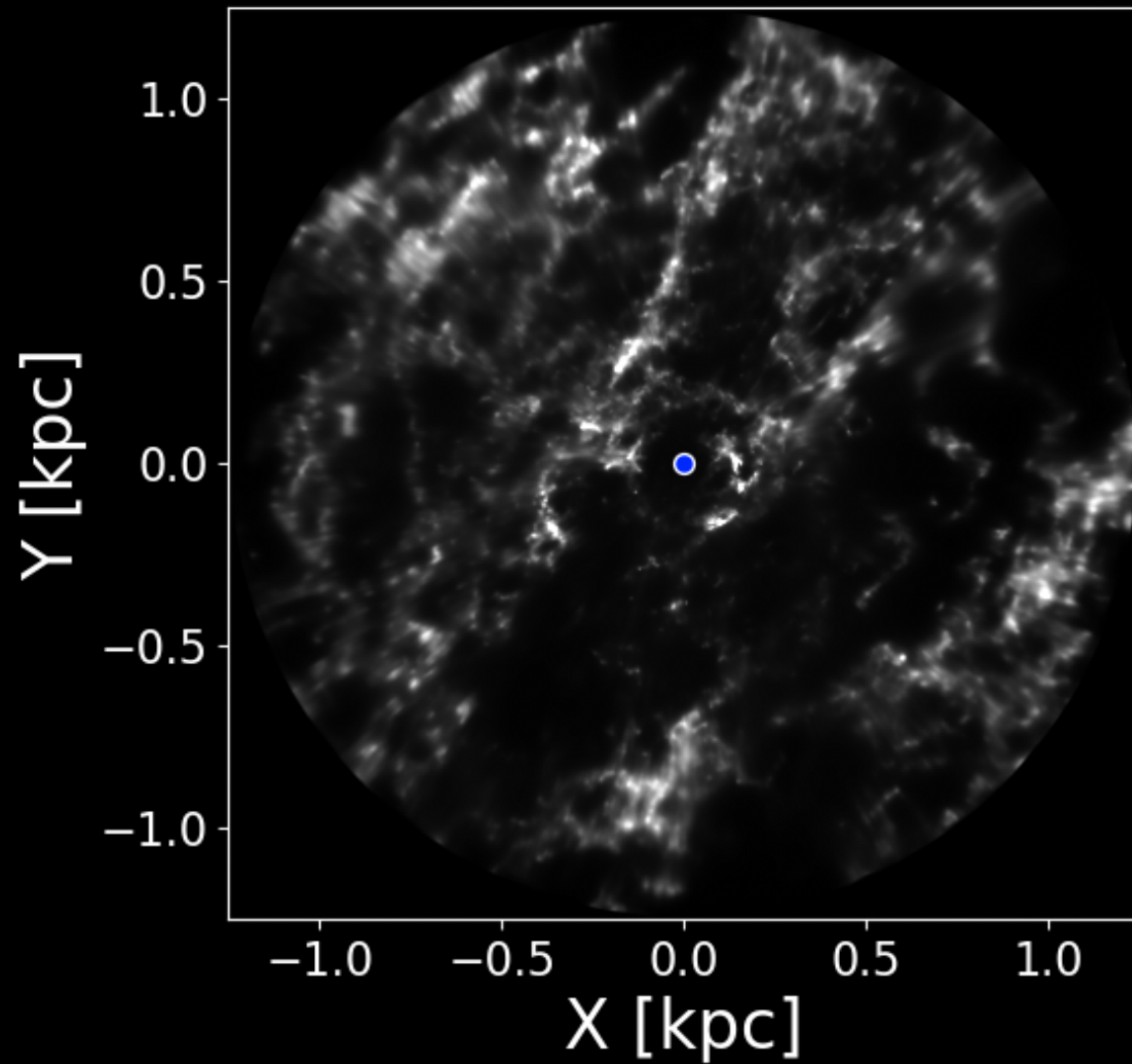
**3D Dust**



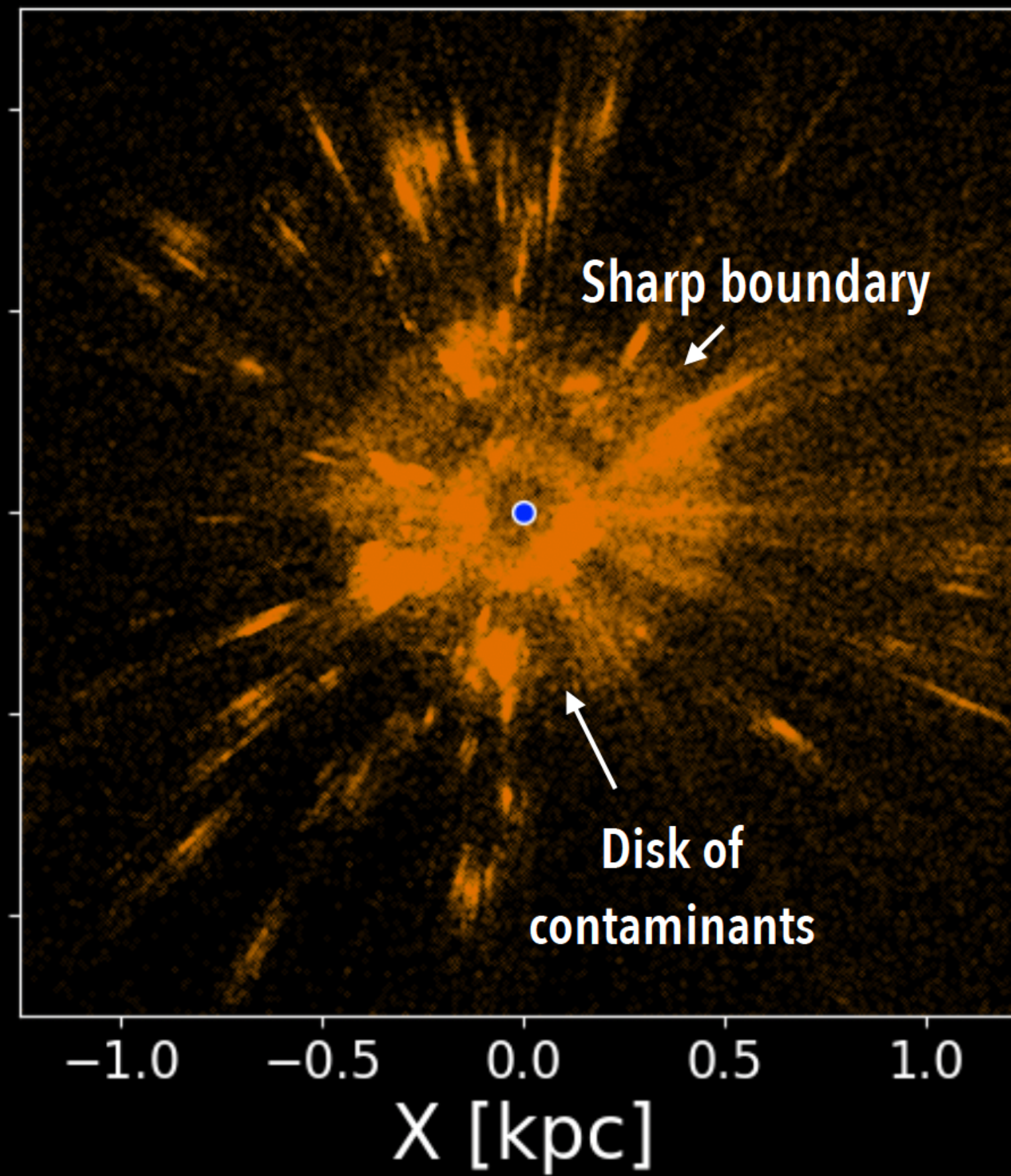


# YSOs in the literature

**3D Dust**



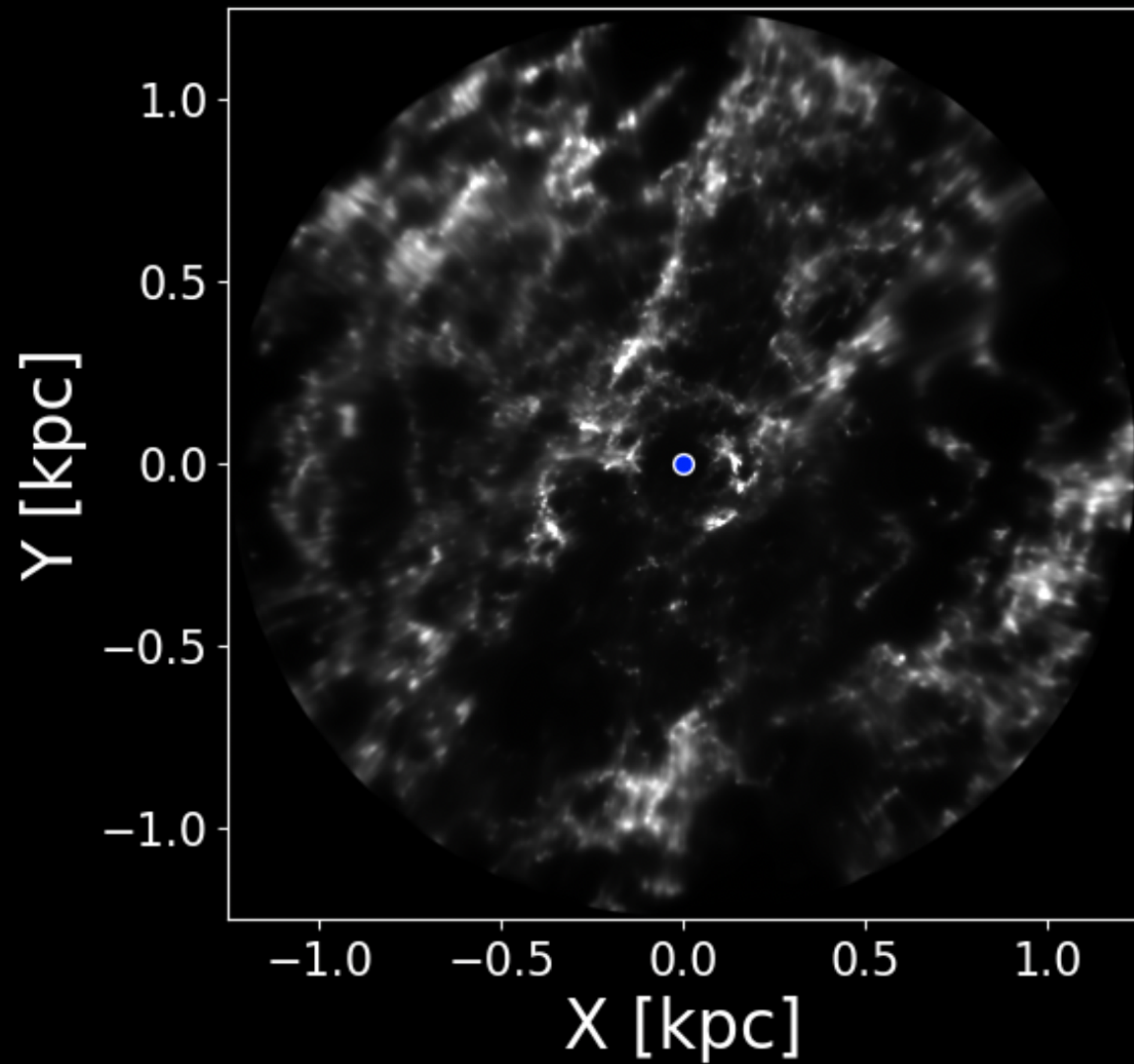
**Young Star (Literature Compilation)**



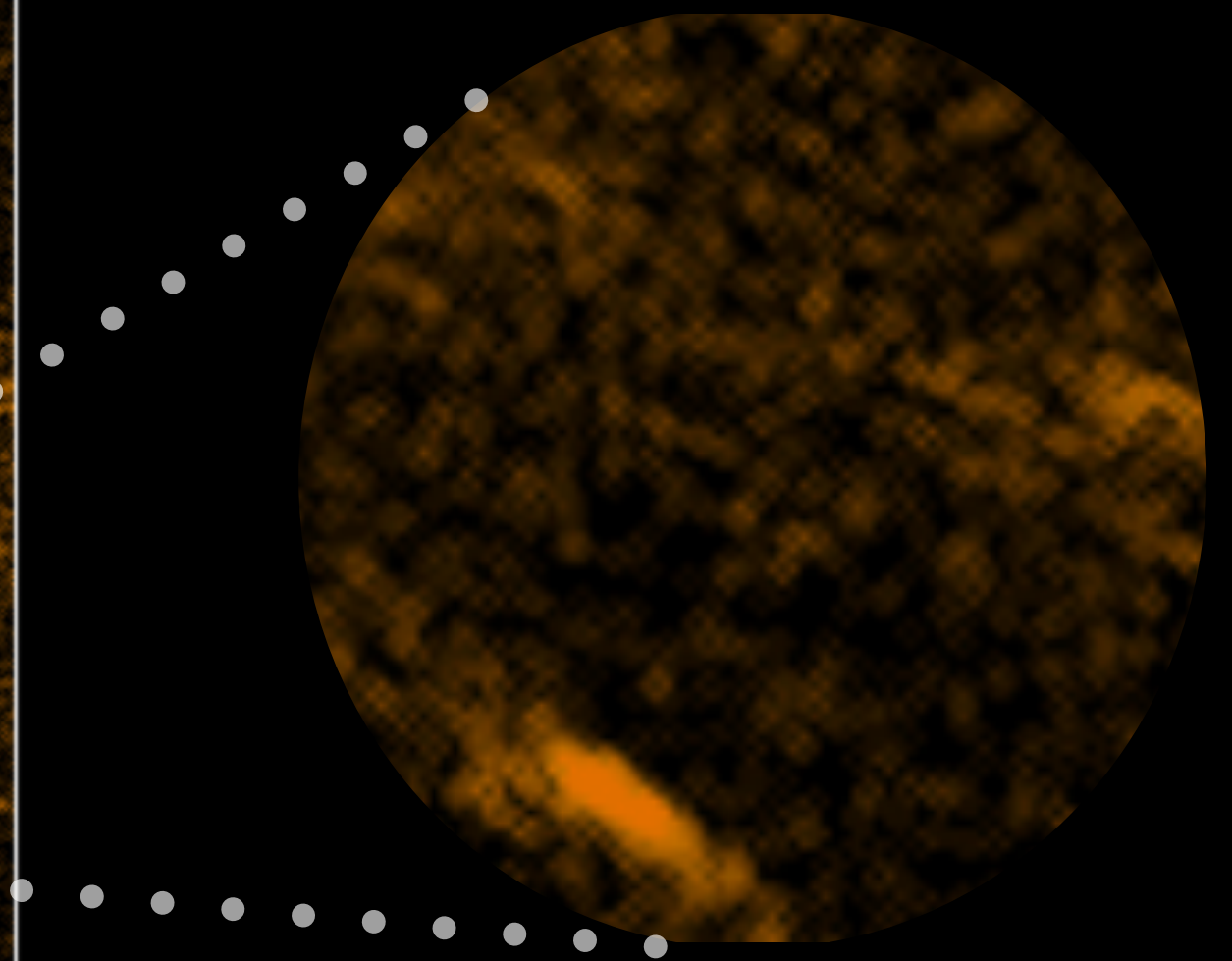
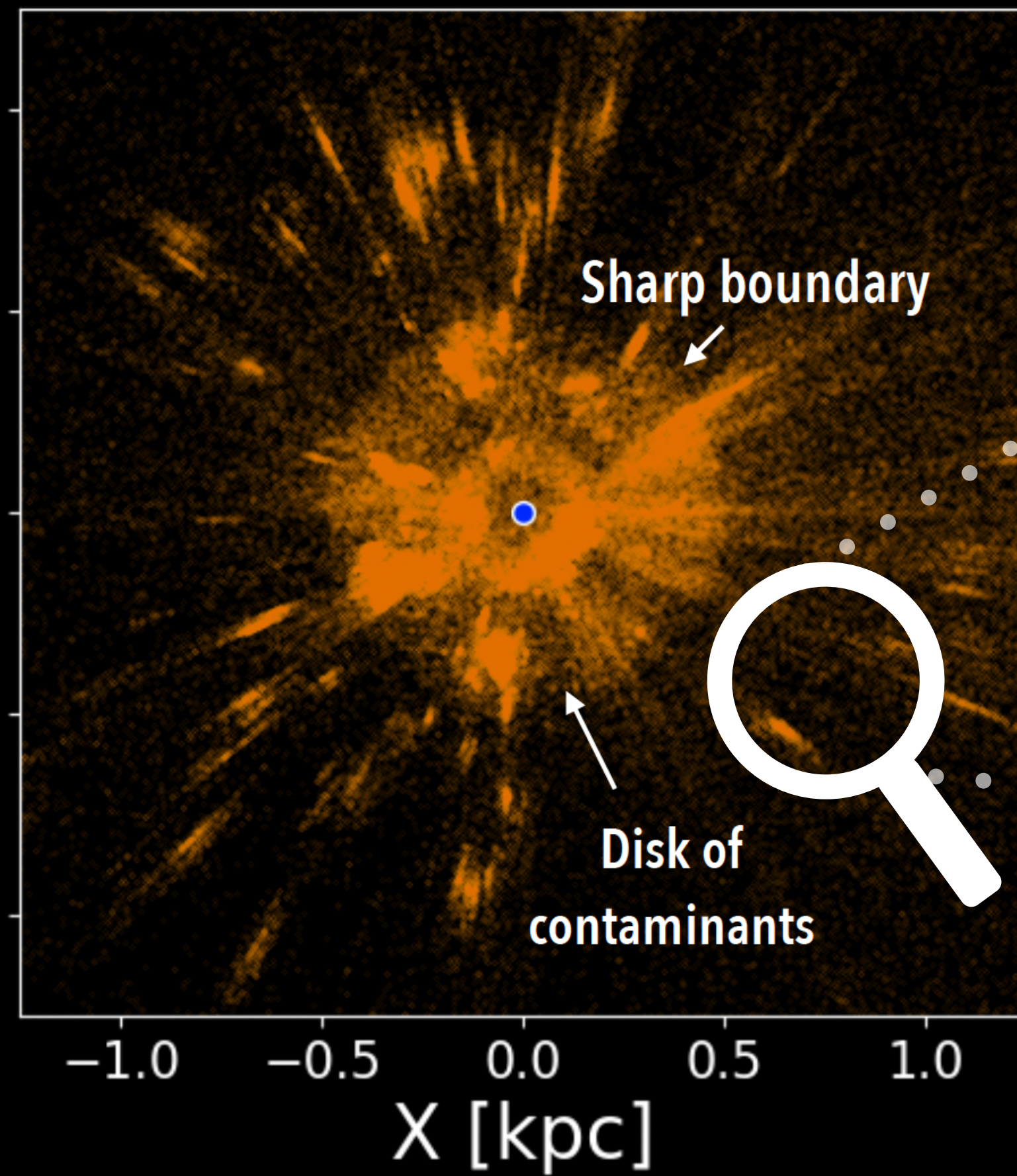


# YSOs in the literature

**3D Dust**



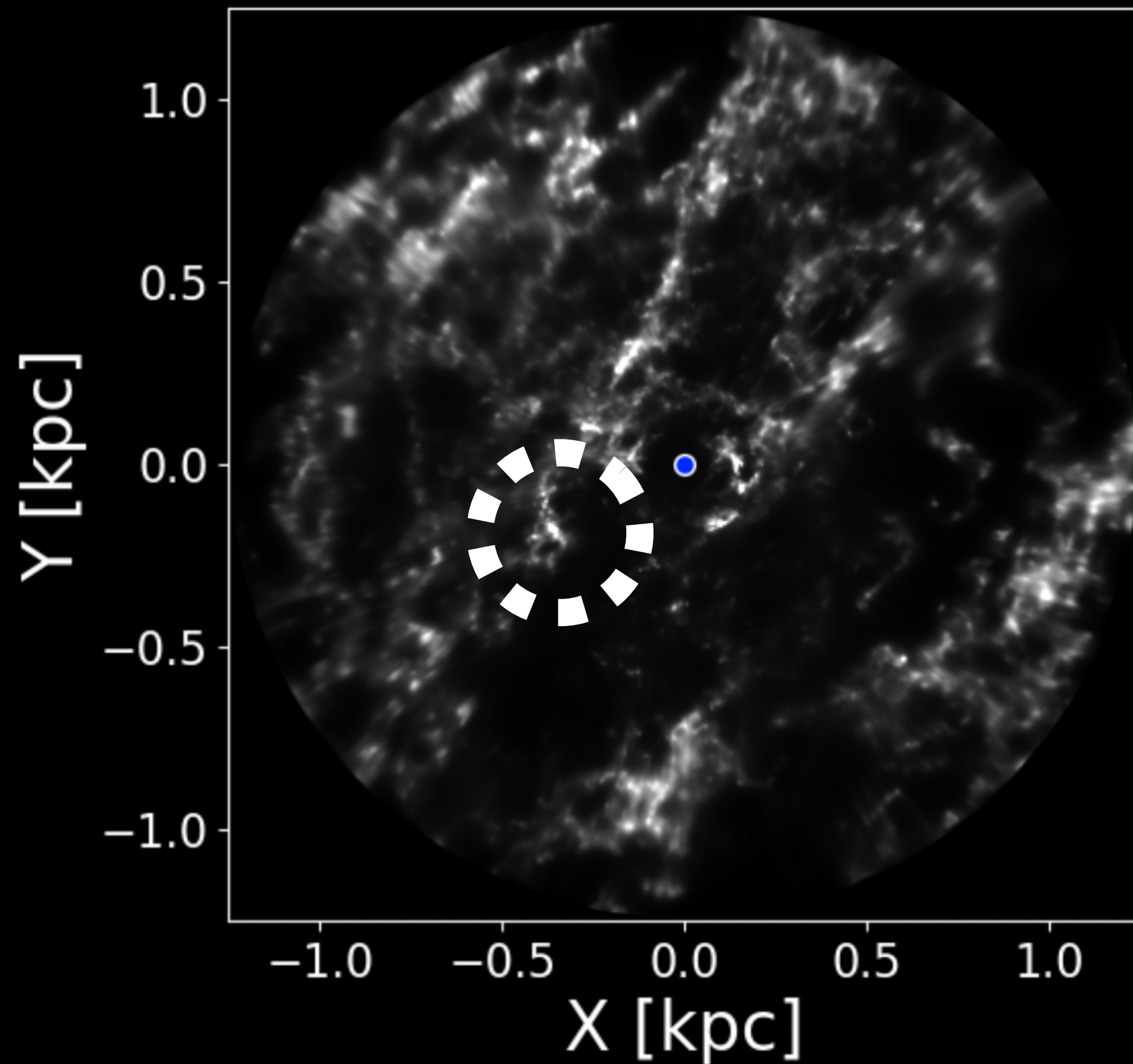
**Young Star (Literature Compilation)**



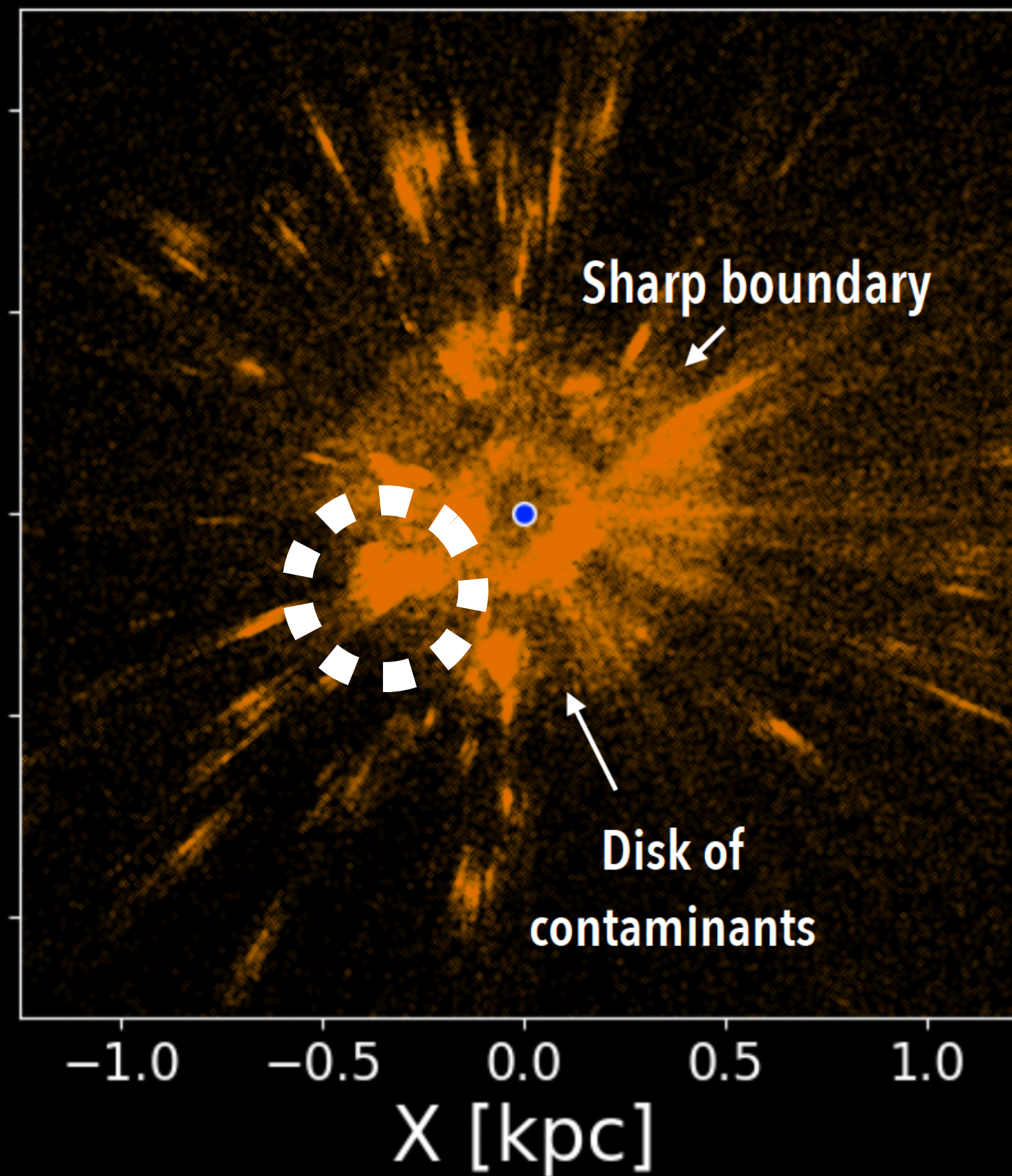


# Sample $\hat{=}$ union of targeted searches

**3D Dust**



**Young Star (Literature Compilation)**





# YSOs in Orion SF region

217 papers on YSOs identified in Orion

Reference	Data type
Struve (1945)	bin
Popper & Plavec (1976)	bin
Hesser et al. (1976)	bin
Marschall & Mathieu (1988)	bin
Bossi et al. (1989)	bin
Mathieu et al. (1991)	bin
Abt et al. (1991)	bin
Popper (1993)	bin
Gagne & Caillault (1994)	X-ray
Kozlova et al. (1995)	SpT, $EW_{\text{acc}}$ , $H\alpha$
Gies et al. (1996)	bin
Hillenbrand (1997)	$T_{\text{eff/bol}}$ , SED
Bolton et al. (1998)	bin
Hillenbrand et al. (1998)	$EW_{\text{acc}}$ disc, SED,
Stassun et al. (1999)	Li, RV, $H\alpha$ , SED
Lloyd & Stickland (1999)	bin
Alcalá et al. (2000)	X-ray, Li, rotation, RV, SpT, $H\alpha$ , bin
Covino et al. (2000)	bin
Herbst et al. (2000)	rotation, $EW_{\text{acc}}$
Rebull et al. (2000)	SpT, SED
Covino et al. (2001)	bin
Rebull (2001)	rotation, SpT, SED
Dolan & Mathieu (2001)	Li, RV, $H\alpha$ , bin, SED
Carpenter et al. (2001)	rotation, SED
Rhode et al. (2001)	rotation, bin
Palla & Stahler (2001)	bin
Béjar et al. (2001)	SpT, SED
Kharchenko (2001)	SpT, bin
Herbst et al. (2002)	rotation
Manset & Bastien (2002)	bin
Feigelson et al. (2002)	X-ray
Maheswar et al. (2003)	Li, SpT, $H\alpha$ , SED
Flaccomio et al. (2003)	X-ray
Skinner et al. (2003)	X-ray
Alcalá et al. (2004)	Li, RV, $T_{\text{eff/bol}}$ , SpT, $H\alpha$ , bin
Scholz & Eislöffel (2004)	rotation $T_{\text{eff/bol}}$ , SpT, logg, $EW_{\text{acc}}$ , $H\alpha$ , SED
Ramírez et al. (2004)	X-ray
Sherry et al. (2004)	SED
Barrado y Navascués et al. (2004)	SpT, $H\alpha$ , SED
Simon et al. (2004)	X-ray
Ali & Noriega-Crespo (2004)	disc, SED

Reference	Data type
Stassun et al. (2004)	Li, bin
Baines et al. (2004)	bin
Scholz & Eislöffel (2005)	rotation, SED
Sicilia-Aguilar et al. (2005)	Li, rotation, RV, SpT, $H\alpha$ , bin
Briceño et al. (2005)	Li, $T_{\text{eff/bol}}$ , SpT, $H\alpha$
Robberto et al. (2005)	bin, SED
Smith et al. (2005)	disc, SED
Getman et al. (2005b)	X-ray
Getman et al. (2005a)	Source list
Kenyon et al. (2005)	Li, RV, $EW_{\text{gravity}}$ , bin
Burningham et al. (2005)	RV, $EW_{\text{gravity}}$ , SED
Franciosini et al. (2006)	X-ray, SpT
Köhler et al. (2006)	bin
Herbig & Griffin (2006)	bin
Stassun et al. (2006)	bin
Caballero (2007)	RV, SpT, SED
Barrado Y Navascués et al. (2007)	bin
Reipurth et al. (2007)	SpT, bin
Lee & Chen (2007)	SpT, $H\alpha$
Briceño et al. (2007)	Li, RV, $H\alpha$
Hernández et al. (2007b)	disc, SED
Barrado y Navascués et al. (2007)	disc, SED
Hernández et al. (2007a)	disc, SED
Irwin et al. (2007)	bin
Caballero (2008)	SED
Stempels et al. (2008)	bin
Caballero & Solano (2008)	SpT, bin, SED
Sacco et al. (2008)	Li, rotation, RV, SpT, $EW_{\text{acc}}$ , $H\alpha$ , bin
López-Santiago & Caballero (2008)	X-ray
Bayo et al. (2008)	$T_{\text{eff/bol}}$
Cargile et al. (2008)	bin
Fűrész et al. (2008)	RV, $H\alpha$
Luhman et al. (2008b)	SED
Maxted et al. (2008)	RV, $EW_{\text{gravity}}$ , bin, SED,
Rodríguez-Ledesma et al. (2009)	rotation, SED,
Fang et al. (2009)	Li, SpT, $EW_{\text{acc}}$ , $H\alpha$ , disc, SED,
Biazzo et al. (2009)	rotation, RV, SpT, $H\alpha$ ,
Frasca et al. (2009)	rotation, $T_{\text{eff/bol}}$ , SpT, SED,
Mace et al. (2009)	bin,
Baldovin-Saavedra et al. (2009)	bin,
Mohanty et al. (2009)	RV, bin, SED,
Tobin et al. (2009)	disc, SED,
Mookerjee & Sandell (2009)	disc, SED,
Hernández et al. (2009)	SpT, $H\alpha$ , SED,
Da Rio et al. (2009)	bin,
Alecian et al. (2009)	rotation, $H\alpha$ , SED,
Parihar et al. (2009)	SpT, $EW_{\text{acc}}$ , veiling, bin, disc, SED,
Connelley & Greene (2010)	bin,
Grinin et al. (2010)	$T_{\text{eff/bol}}$ , SpT,
Da Rio et al. (2010)	rotation, SED,
Cody & Hillenbrand (2010)	bin,
Leone et al. (2010)	$T_{\text{eff/bol}}$ , $H\alpha$ , disc, SED,
Rigliaco et al. (2011)	X-ray, disc, SED,
Barrado et al. (2011)	X-ray, SpT, disc,
Franciosini & Sacco (2011)	Li, $T_{\text{eff/bol}}$ , SpT, $EW_{\text{gravity}}$ ,
Bayo et al. (2011)	X-ray, SpT, $H\alpha$ ,
Ingleby et al. (2011)	bin,
van Eyken et al. (2011)	rotation,
Morales-Calderón et al. (2011)	disc, SED,
Béjar et al. (2011)	bin,
Daemgen et al. (2012)	bin,

Reference	Data type
Caballero et al. (2012)	Li, SpT, $EW_{\text{acc}}$ , $H\alpha$ , disc,
Bayo et al. (2012)	rotation, SpT, $EW_{\text{acc}}$ , $H\alpha$ , disc,
Rigliaco et al. (2012)	SpT, $EW_{\text{acc}}$ , $H\alpha$ ,
Kryukova et al. (2012)	disc, SED,
Megeath et al. (2012)	disc, SED,
Gómez Maqueo Chew et al. (2012)	bin,
Da Rio et al. (2012)	$T_{\text{eff/bol}}$ ,
Hsu et al. (2012)	Li, SpT, $EW_{\text{acc}}$ , $H\alpha$ , SED,
Morales-Calderón et al. (2012)	rotation, SpT, bin,
Peña Ramírez et al. (2012)	SED,
Manara et al. (2012)	$H\alpha$ , disc,
Correia et al. (2013)	bin,
Hillenbrand et al. (2013)	SpT, bin,
Hsu et al. (2013)	Li, RV, SpT, disc, SED,
Stutz et al. (2013)	$T_{\text{eff/bol}}$ , bin, disc, SED,
Windemuth et al. (2013)	bin,
Pillitteri et al. (2013)	X-ray,
Kim et al. (2013)	X-ray, $T_{\text{eff/bol}}$ , SpT, bin, disc,
Fang et al. (2013)	Li, SpT, $EW_{\text{acc}}$ , $H\alpha$ , disc, SED,
Szegedi-Elek et al. (2013)	$H\alpha$ , SED,
Kuhn et al. (2013)	SED,
Povich et al. (2013)	Source list,
Broos et al. (2013)	X-ray, disc,
López-García et al. (2013)	X-ray,
Sergison et al. (2013)	Li, RV, veiling, $H\alpha$ ,
Bouy et al. (2014)	SED,
Pettersson et al. (2014)	$H\alpha$ , bin, SED,
Fang et al. (2014)	bin,
Sanchez et al. (2014)	$T_{\text{eff/bol}}$ , SED,
Ingraham et al. (2014)	$T_{\text{eff/bol}}$ , SpT,
Kuhn et al. (2014)	SpT, bin,
Hernández et al. (2014)	Li, RV, SpT, $EW_{\text{gravity}}$ , $H\alpha$ , bin, disc, SED,
Theissen & West (2014)	RV, $T_{\text{eff/bol}}$ , SpT, $EW_{\text{gravity}}$ , $EW_{\text{acc}}$ , $H\alpha$ ,
Downes et al. (2014)	SpT, $H\alpha$ ,
Rice et al. (2015)	disc, SED,
Downes et al. (2015)	SpT, $EW_{\text{acc}}$ , $H\alpha$ , disc, SED,
Meingast et al. (2016)	rotation, bin, SED,
Hasenberger et al. (2016)	$T_{\text{eff/bol}}$ , SpT, $H\alpha$ , disc, SED,
Megeath et al. (2016)	SED,
Karim et al. (2016)	SpT, disc,
Da Rio et al. (2016)	disc, SED,
Kounkel et al. (2016a)	rotation, $T_{\text{eff/bol}}$ , logg, SED,
Kounkel et al. (2016b)	Li, RV, $T_{\text{eff/bol}}$ , bin,
Lewis & Lada (2016)	bin, disc,
Furlan et al. (2016)	bin, disc,
Kim et al. (2016)	$T_{\text{eff/bol}}$ , disc, SED,
Messina et al. (2016)	Li, $T_{\text{eff/bol}}$ , SpT, $H\alpha$ , bin, disc,
Pillitteri et al. (2017)	bin,
Fang et al. (2017)	X-ray,
Suárez et al. (2017)	Li, SpT, $H\alpha$ , disc,
Kounkel et al. (2017a)	Li, $T_{\text{eff/bol}}$ , SpT, $H\alpha$ , disc,
Kounkel et al. (2017b)	Li, SpT, $H\alpha$ ,
Simon & Toraskar (2017)	bin,
Kounkel et al. (2017c)	bin,
Jaehnig et al. (2017)	Li, RV, $T_{\text{eff/bol}}$ , $H\alpha$ , disc,
Getman et al. (2017)	bin,
Messina et al. (2017)	X-ray, SED,
Fernandez et al. (2017)	SED,
GRAVITY Collaboration et al. (2018)	bin,

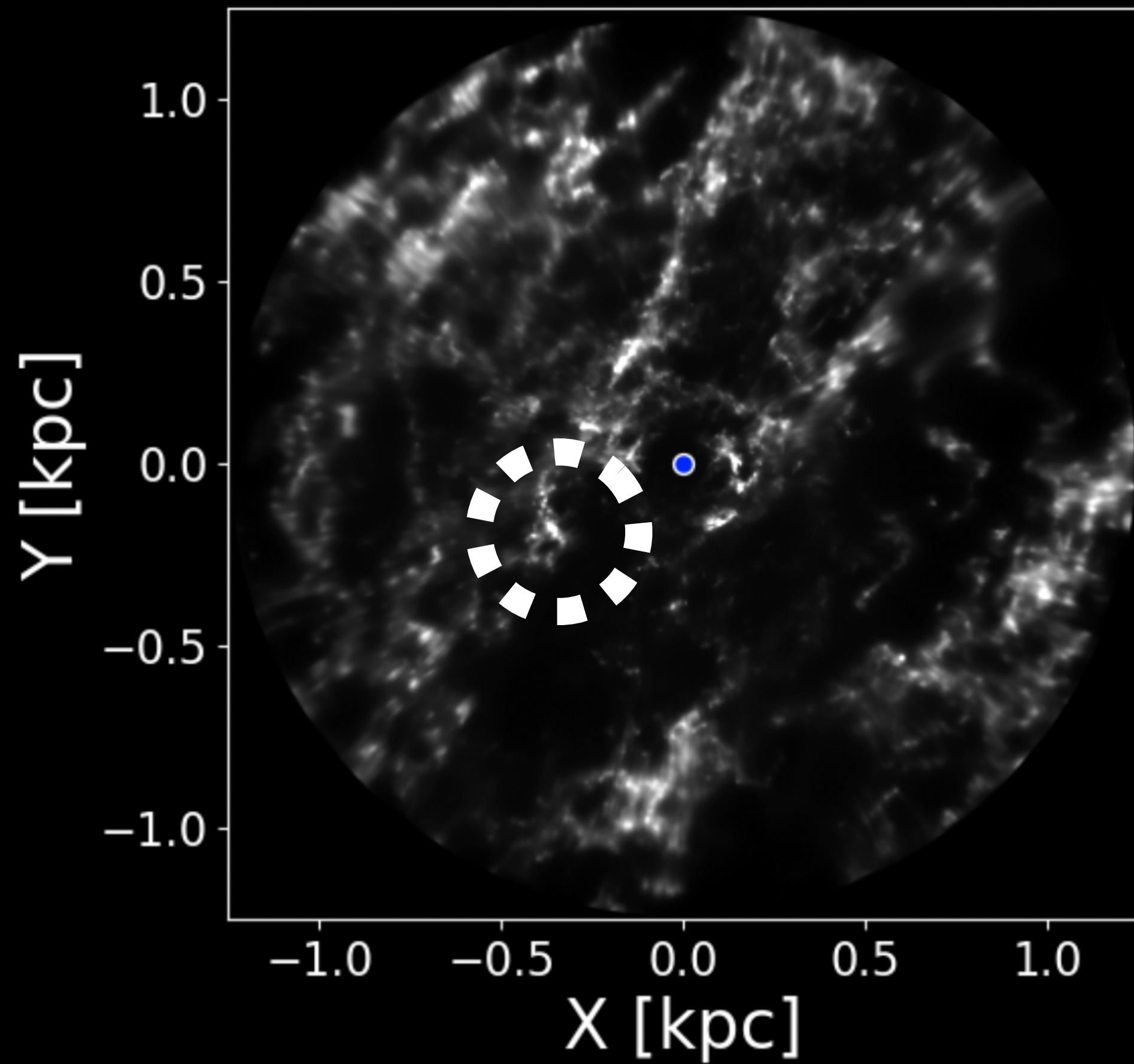
Reference	Data type
Vioque et al. (2018)	$T_{\text{eff/bol}}$ , $H\alpha$ , bin, disc, SED,
Zari et al. (2018)	RV,
Kounkel et al. (2018)	RV, $T_{\text{eff/bol}}$ , logg,
Caballero et al. (2018)	bin,
Grant et al. (2018)	SpT, disc, SED,
Yao et al. (2018)	$T_{\text{eff/bol}}$ , logg, $EW_{\text{acc}}$ , disc, SED,
Cottle et al. (2018)	disc, SED,
Davies et al. (2018)	bin,
Getman et al. (2018)	bin,
Getman et al. (2018)	disc,
Jayasinghe et al. (2018)	rotation,
Duchêne et al. (2018)	bin,
Großschedl et al. (2019)	disc, SED,
Jerabkova et al. (2019)	bin,
Caballero et al. (2019)	Li, SpT, $EW_{\text{acc}}$ , $H\alpha$ ,
Briceño et al. (2019)	Li, $T_{\text{eff/bol}}$ , SpT, $EW_{\text{gravity}}$ , $H\alpha$ , SED,
Kim et al. (2019)	bin, SED,
Kounkel et al. (2019)	RV, $T_{\text{eff/bol}}$ , logg, veiling, bin, disc,
Mairs et al. (2019)	bin,
McBride & Kounkel (2019)	RV,
Tobin et al. (2019)	bin,
De Furio et al. (2019)	bin,
Suárez et al. (2019)	SED,
Bouma et al. (2020)	bin,
Tokovinin et al. (2020)	bin,
Kounkel et al. (2020)	Source list,
Manzo-Martínez et al. (2020)	SpT, $H\alpha$ ,
Robberto et al. (2020)	$T_{\text{eff/bol}}$ ,
Strampelli et al. (2020)	bin,
Fischer et al. (2020)	SED,
Jackson et al. (2020)	RV, $T_{\text{eff/bol}}$ , logg, SED,
Lavail et al. (2020)	bin,
Frasca et al. (2021)	Source list,
Pinzón et al. (2021)	rotation, RV, $T_{\text{eff/bol}}$ , SpT, bin, disc,
Fang et al. (2021)	SpT, veiling, bin,
Habel et al. (2021)	disc,
Laos et al. (2021)	$EW_{\text{acc}}$ ,
Serna et al. (2021)	rotation, $T_{\text{eff/bol}}$ , bin, SED,
Arun et al. (2021)	bin,
Kos et al. (2021)	Li, rotation, RV, $T_{\text{eff/bol}}$ , logg,
Proffitt et al. (2021)	bin,
Franciosini et al. (2022)	RV, $T_{\text{eff/bol}}$ ,
Flaischlen et al. (2022)	$T_{\text{eff/bol}}$ ,
Thanathibodee et al. (2022)	$T_{\text{eff/bol}}$ , SpT, $EW_{\text{acc}}$ , $H\alpha$ , bin,
Kounkel et al. (2022)	rotation,
Pittman et al. (2022)	rotation,
Cao et al. (2022)	$T_{\text{eff/bol}}$ , disc,
De Furio et al. (2022b)	bin,
Theissen et al. (2022)	rotation, RV, $T_{\text{eff/bol}}$ , veiling, bin,
De Furio et al. (2022a)	bin,
Hernández et al. (2023)	Li, SpT, $H\alpha$ ,
Federman et al. (2023)	$T_{\text{eff/bol}}$ , disc,
Damian et al. (2023)	$T_{\text{eff/bol}}$ , SpT, SED,
Smith et al. (2023)	rotation, $T_{\text{eff/bol}}$ , disc,

Roquette+2025

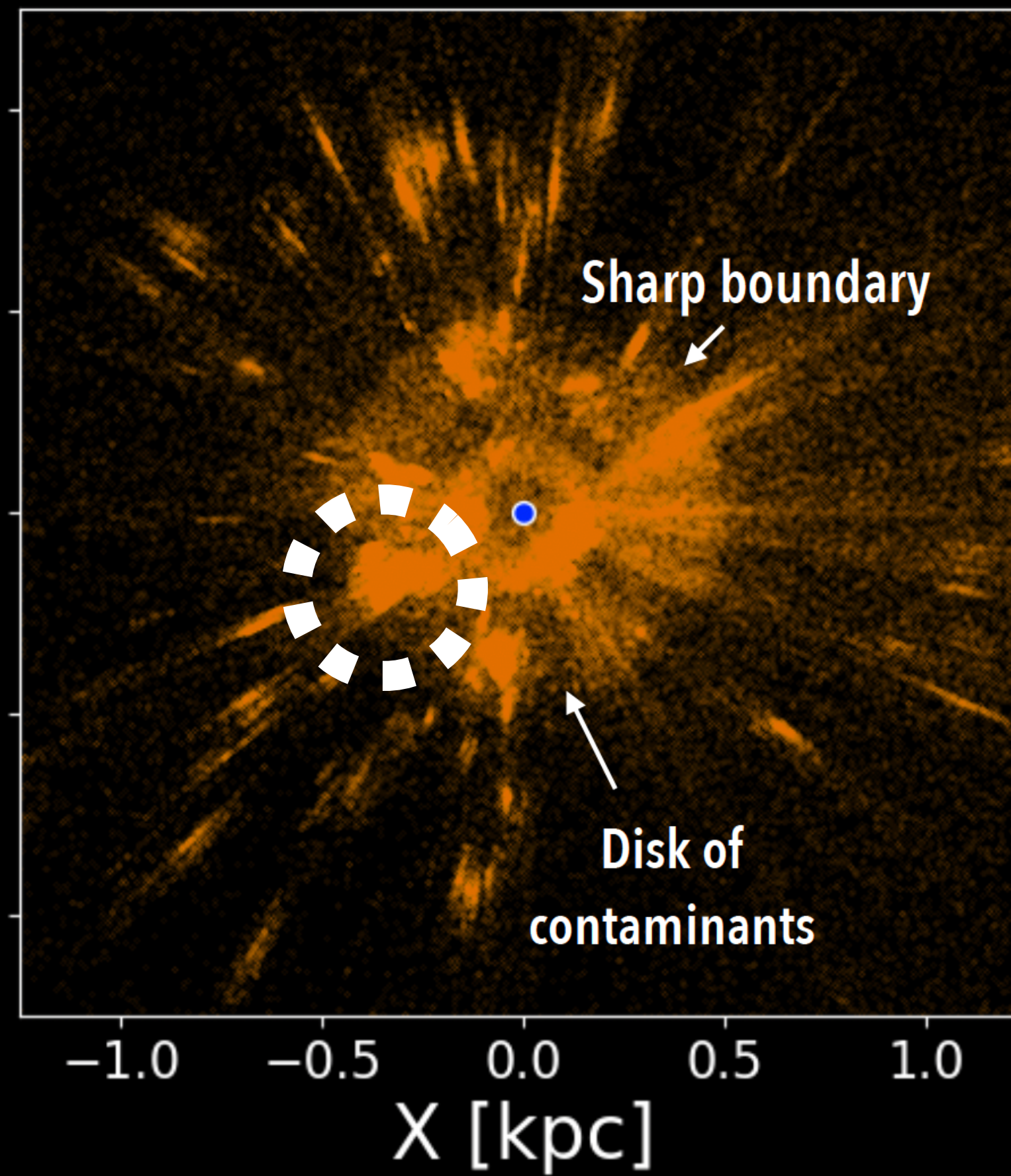


# Inhomogeneous selection function

**3D Dust**



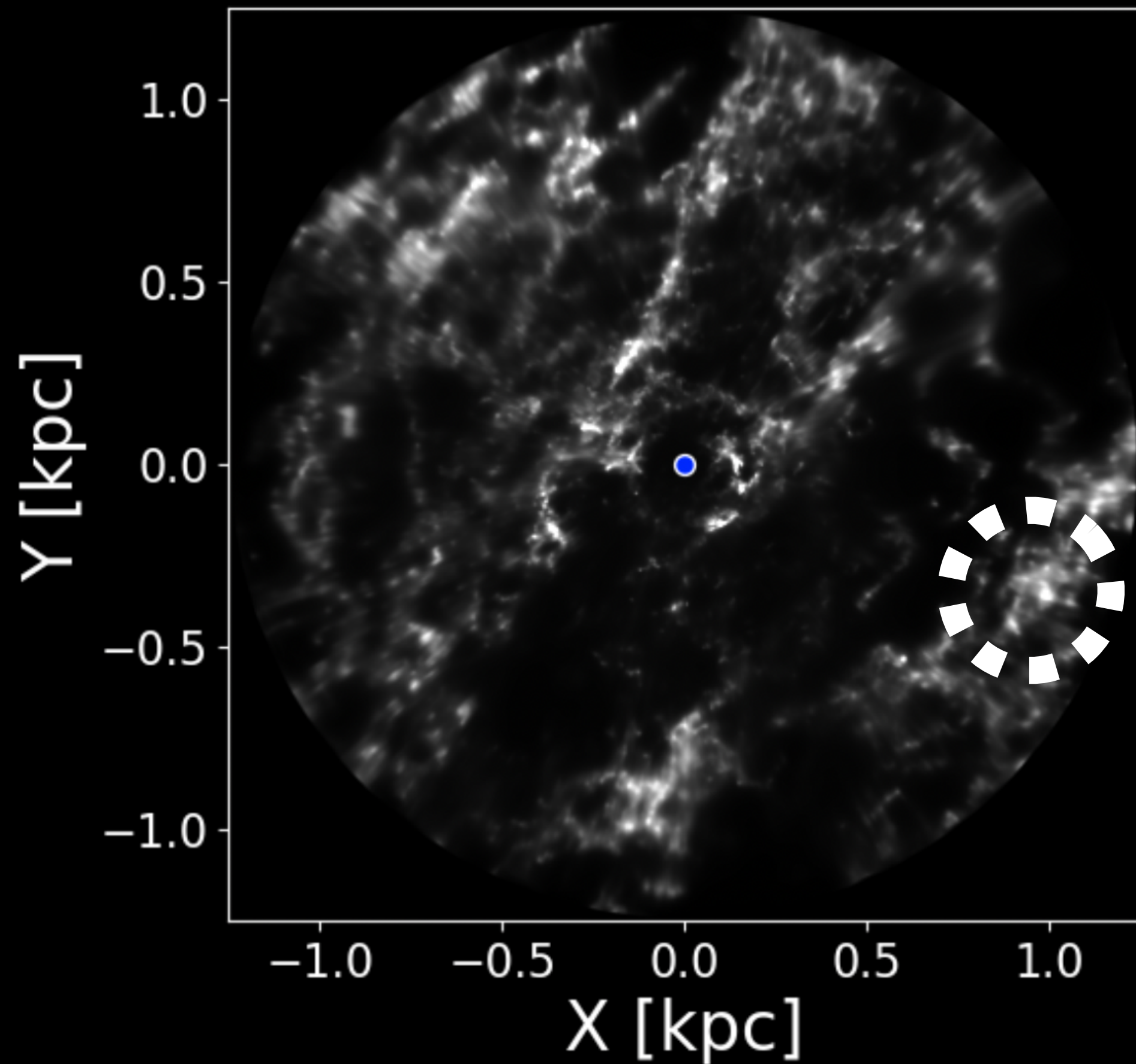
**Young Star (Literature Compilation)**



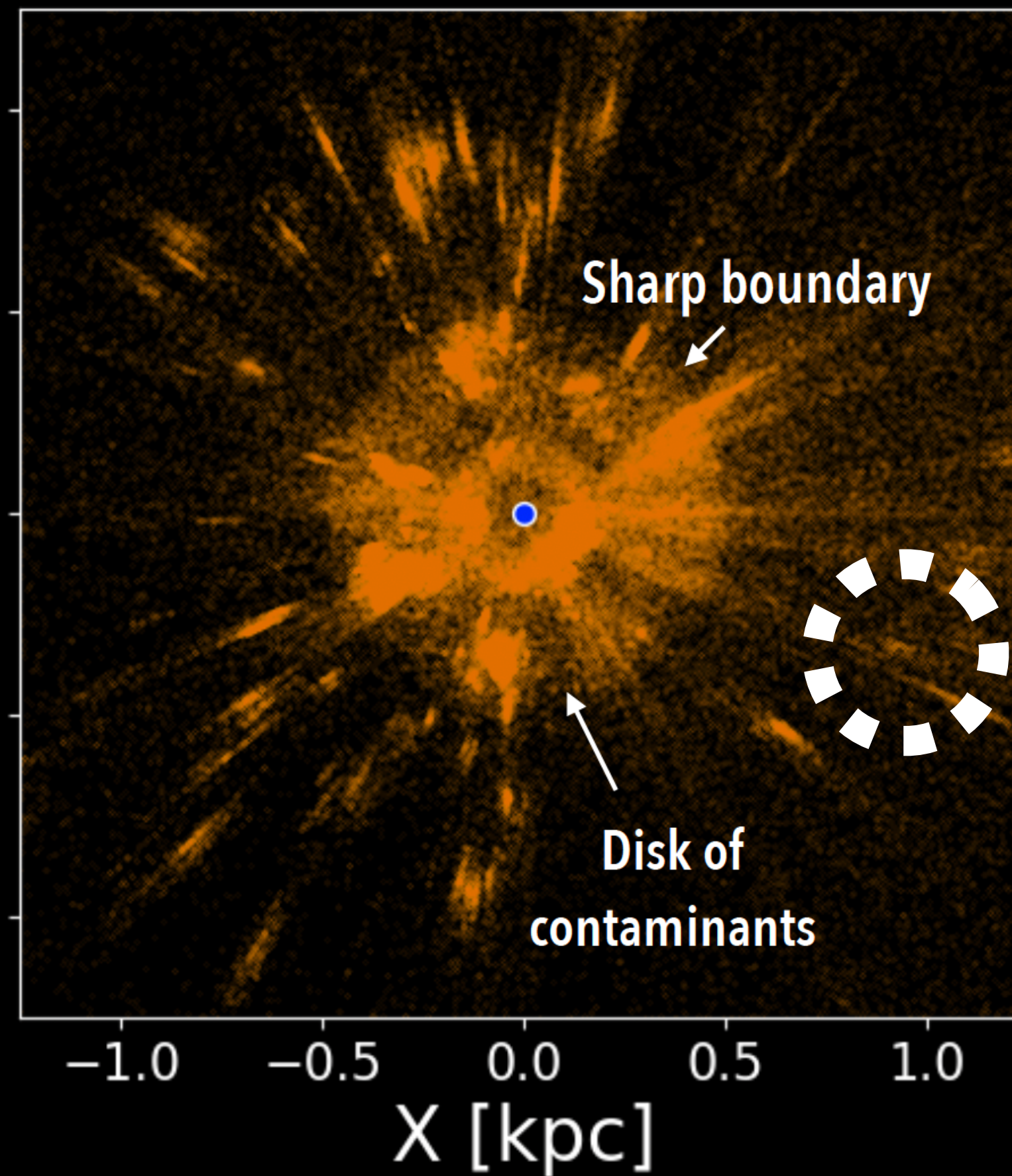


# Inhomogeneous selection function

**3D Dust**

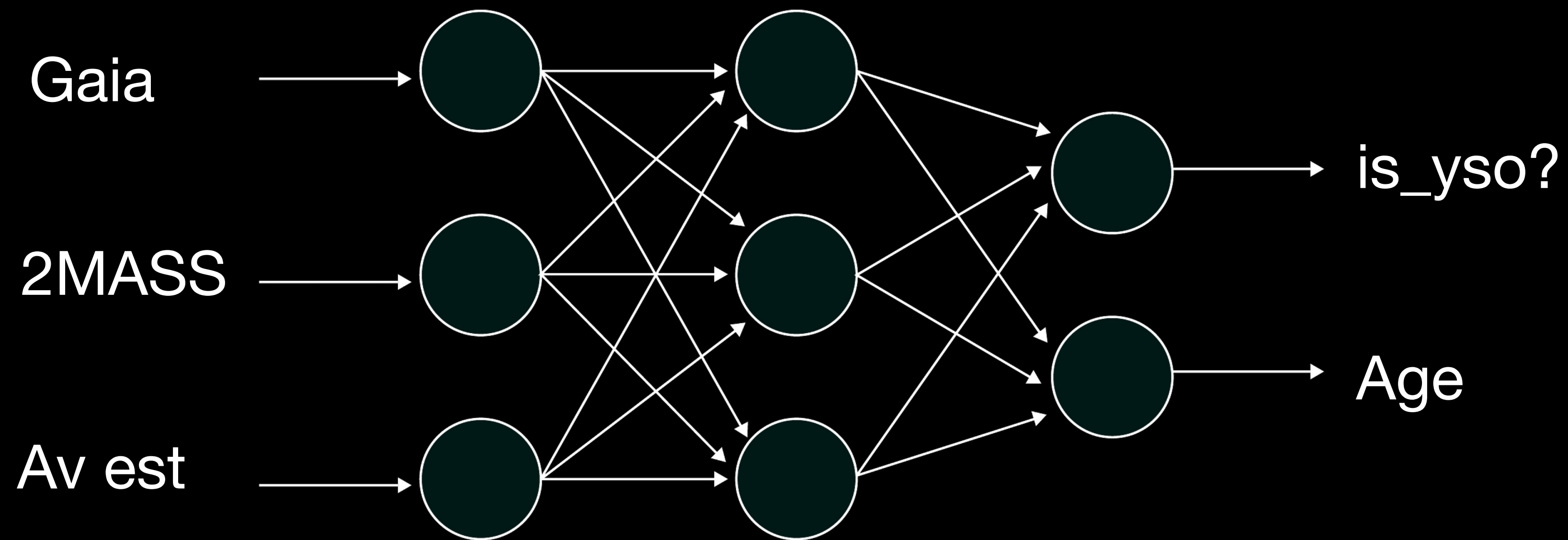


**Young Star (Literature Compilation)**



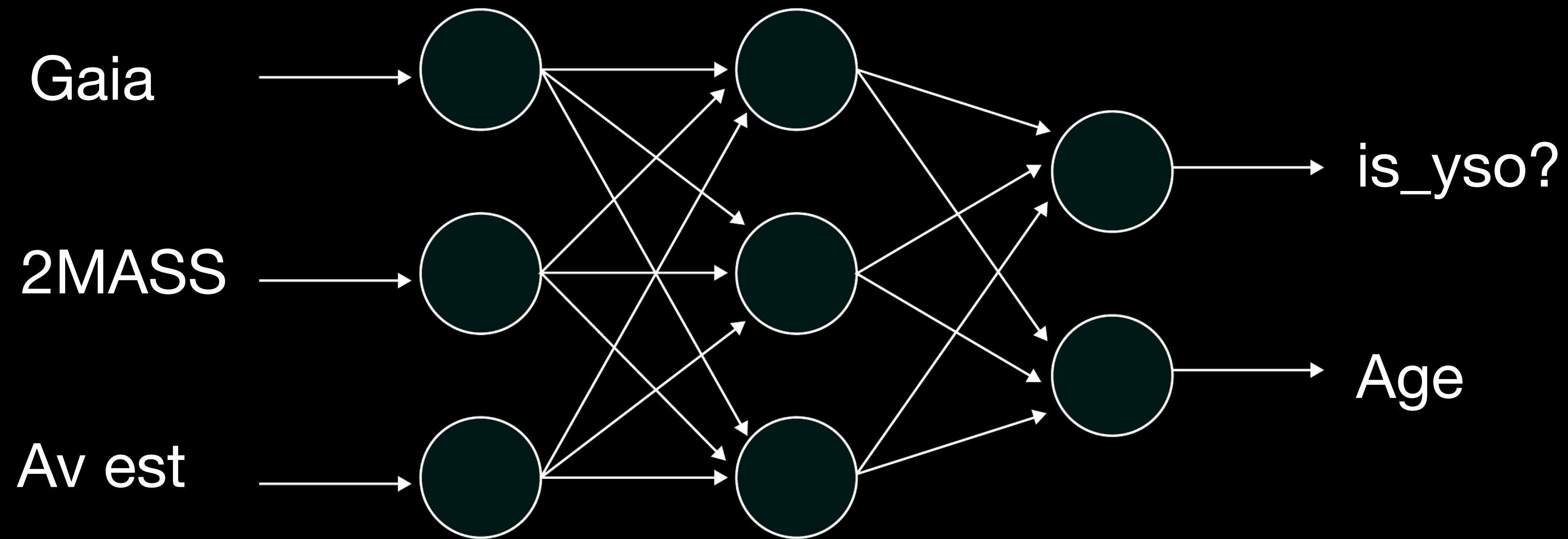


# Can ML help?

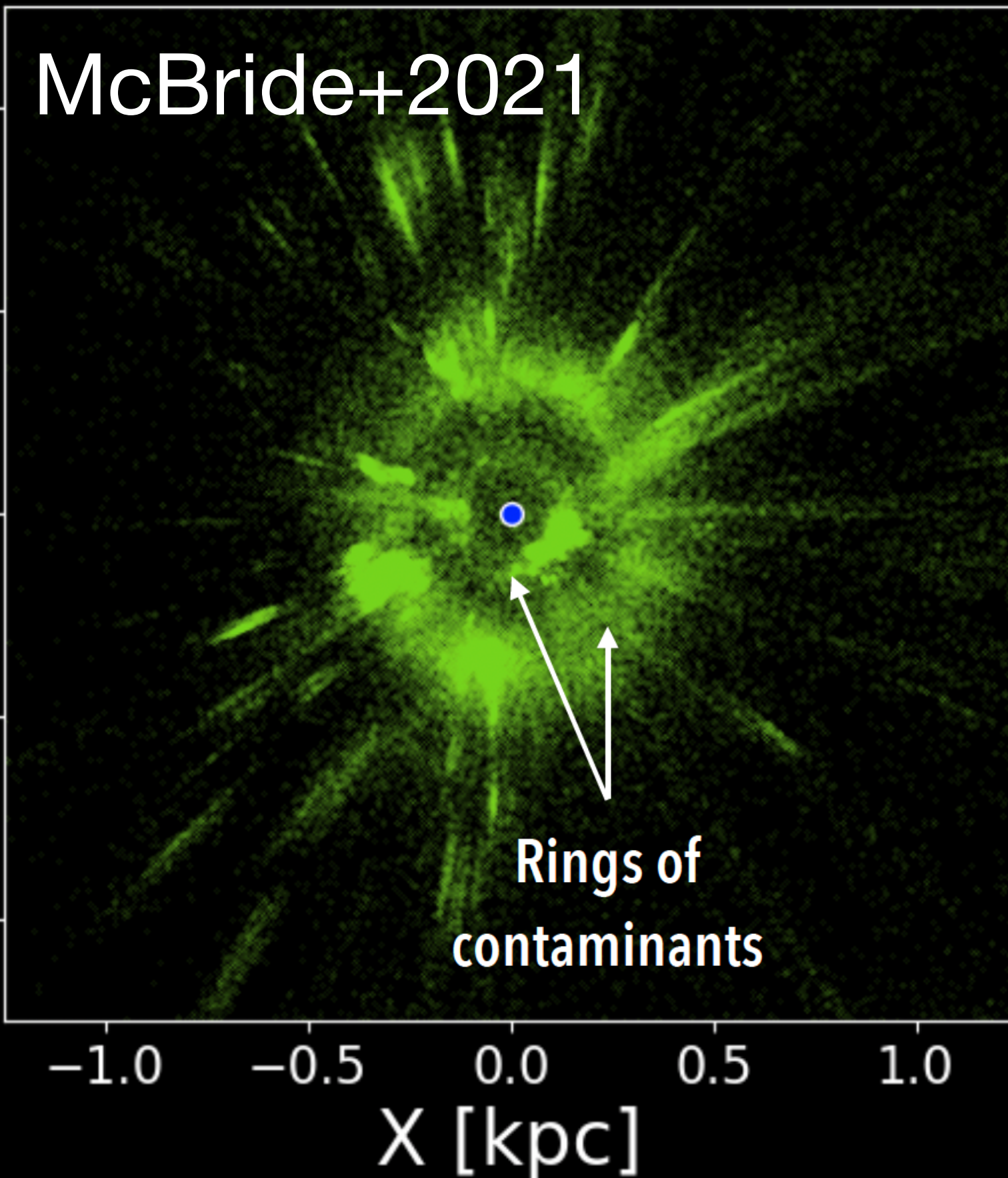




# Can ML help?



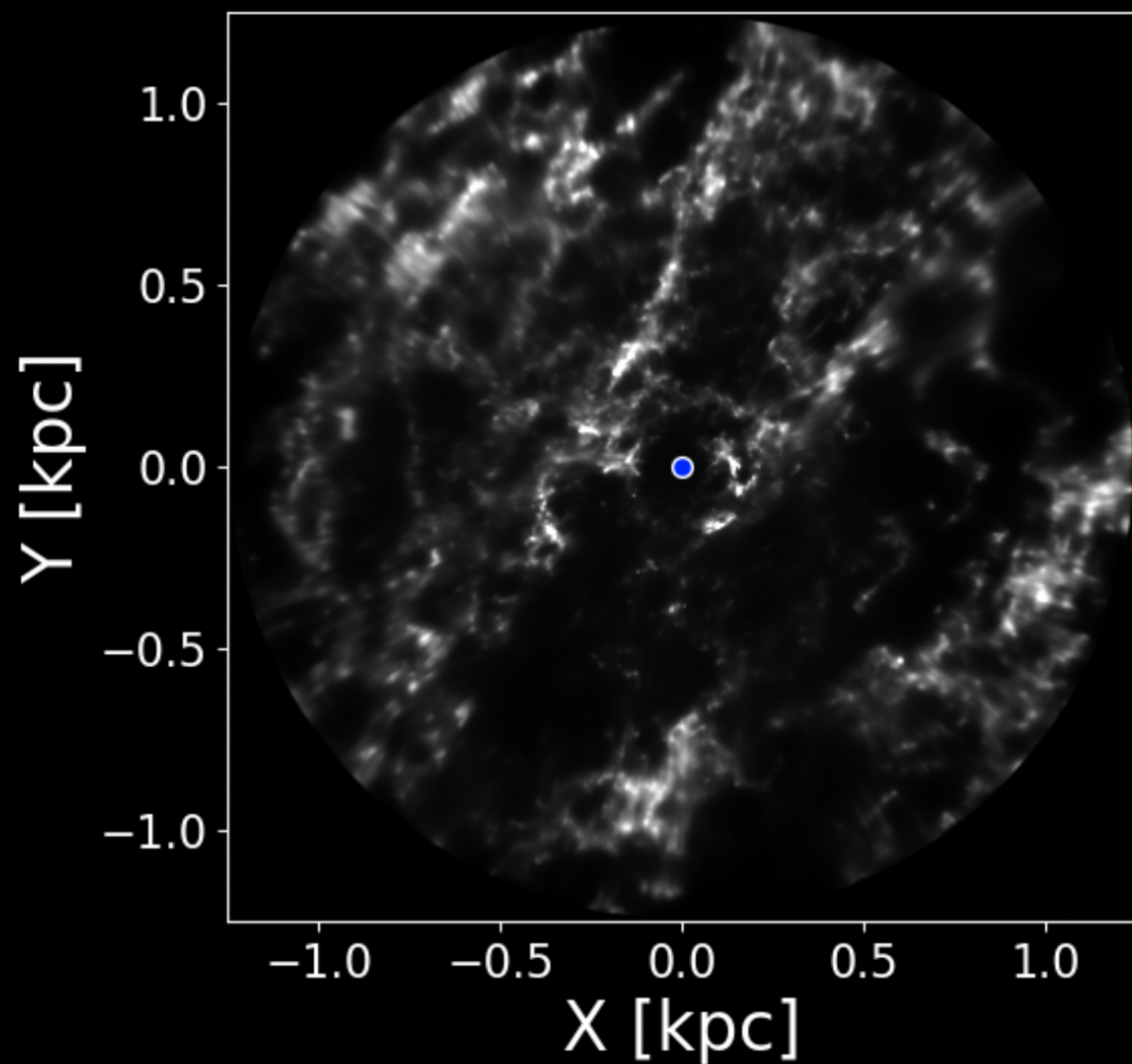
## Young Star (SOTA Machine Learning)



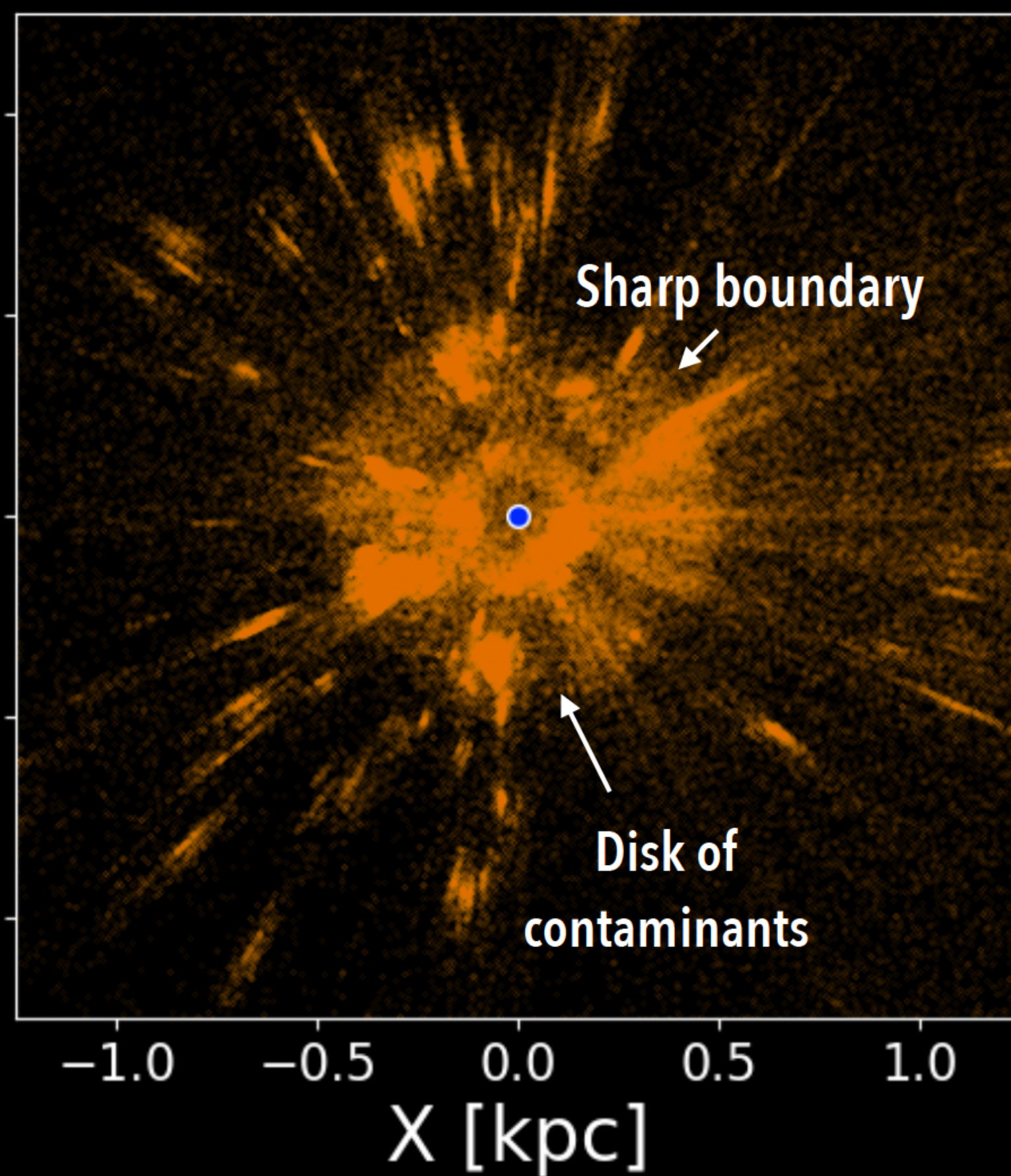


# Can ML help?

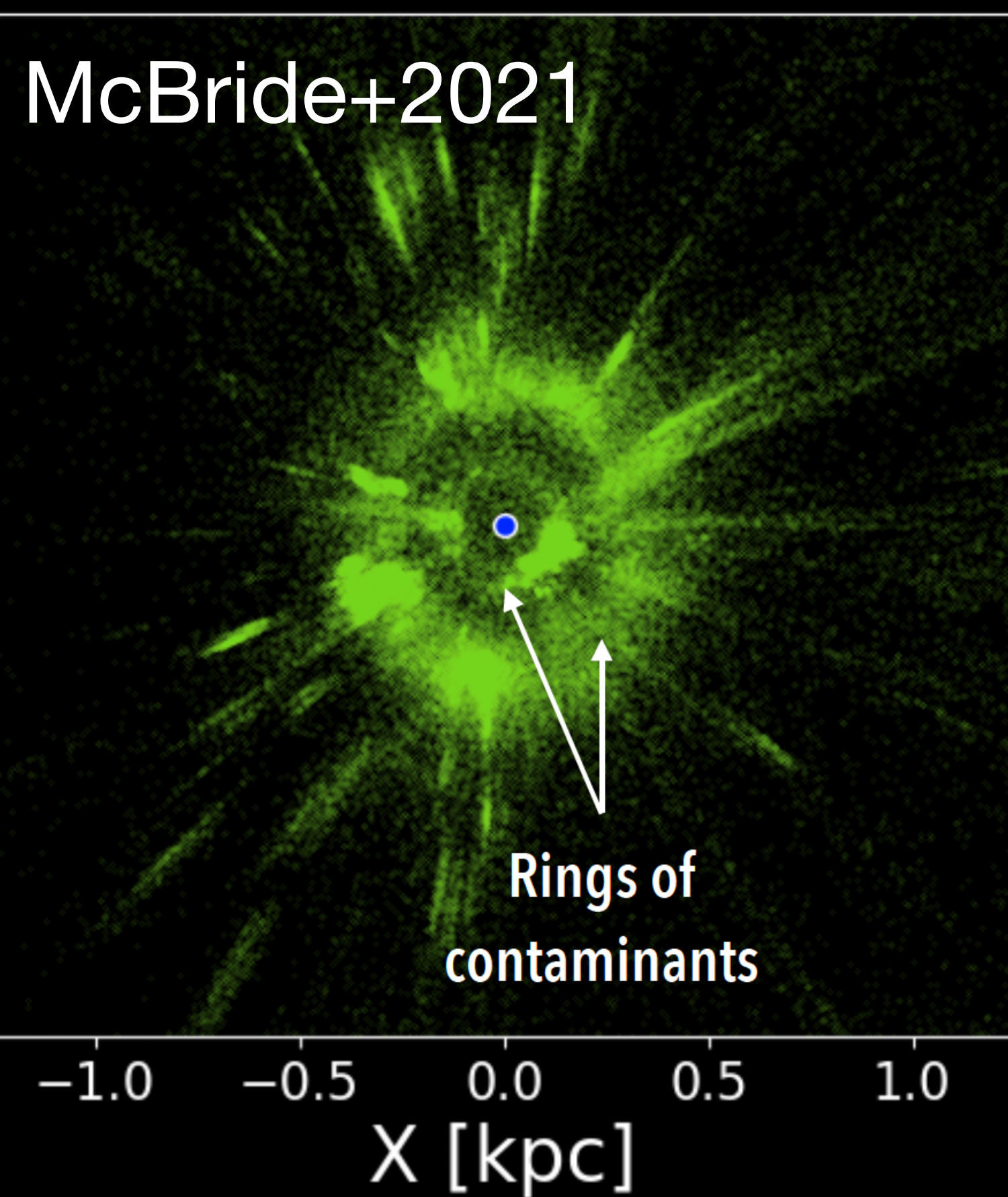
**3D Dust**



**Young Star (Literature Compilation)**



**Young Star (SOTA Machine Learning)**



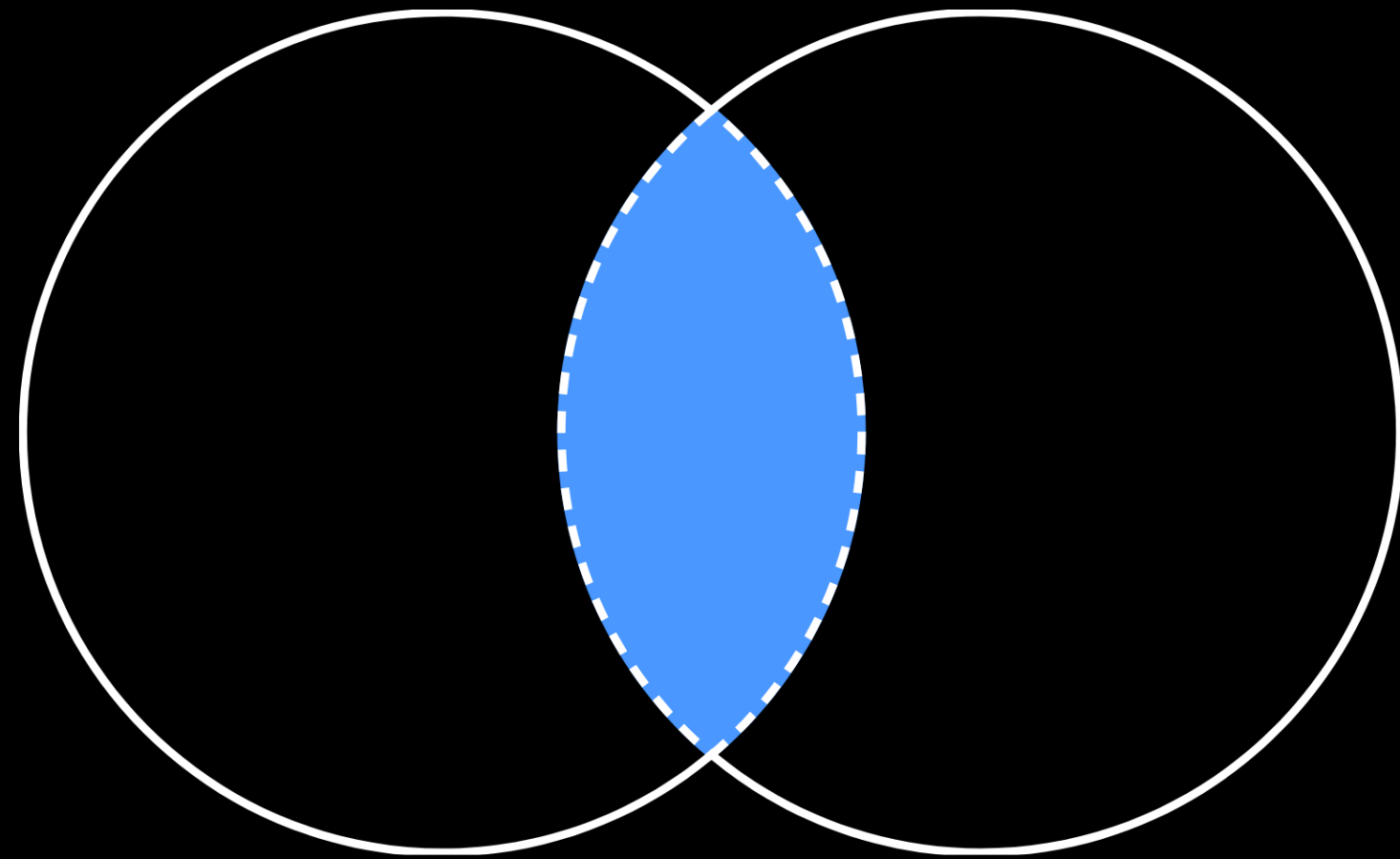


# Ways to improve ML YSO catalog

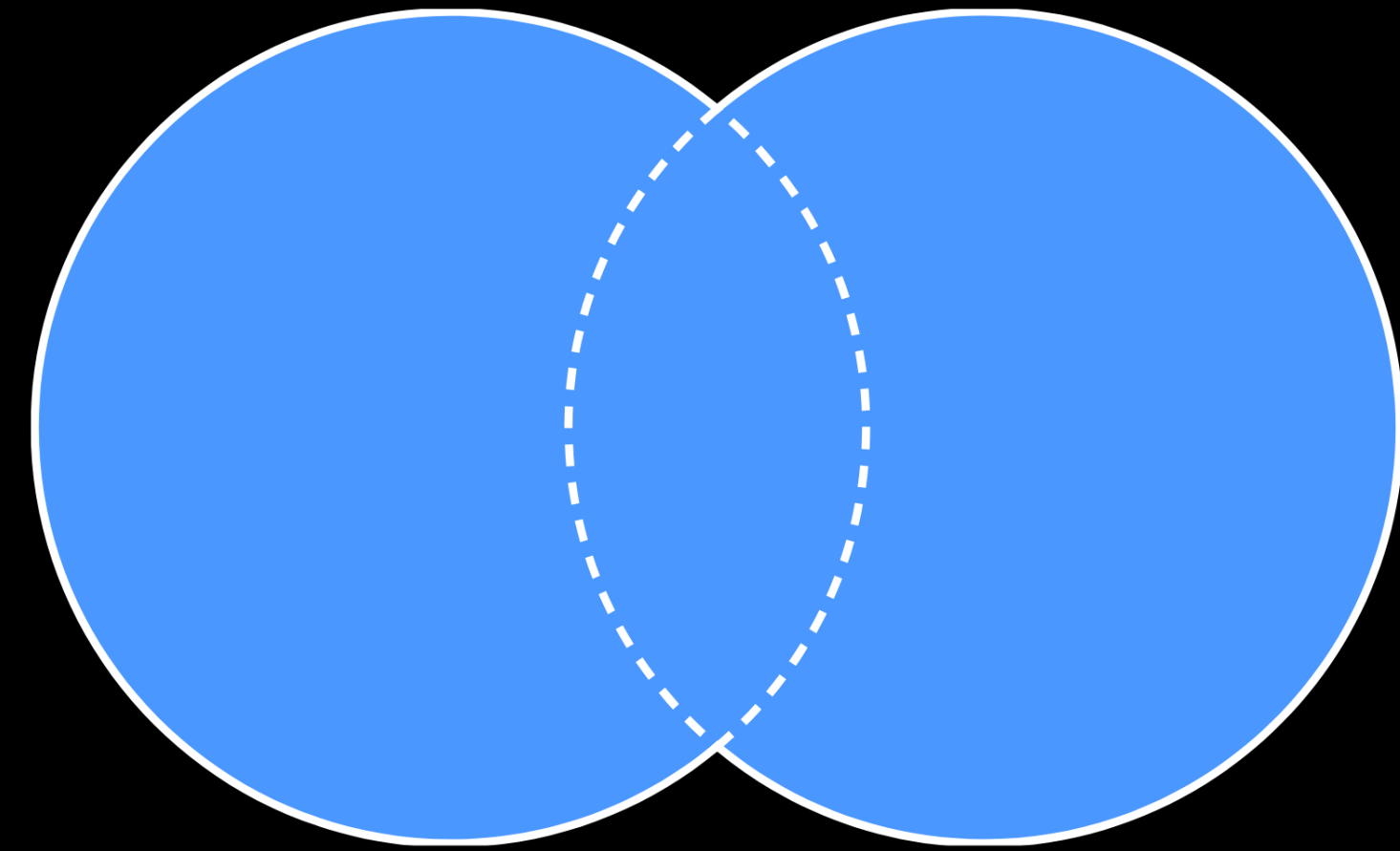


# Ways to improve ML YSO catalog

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



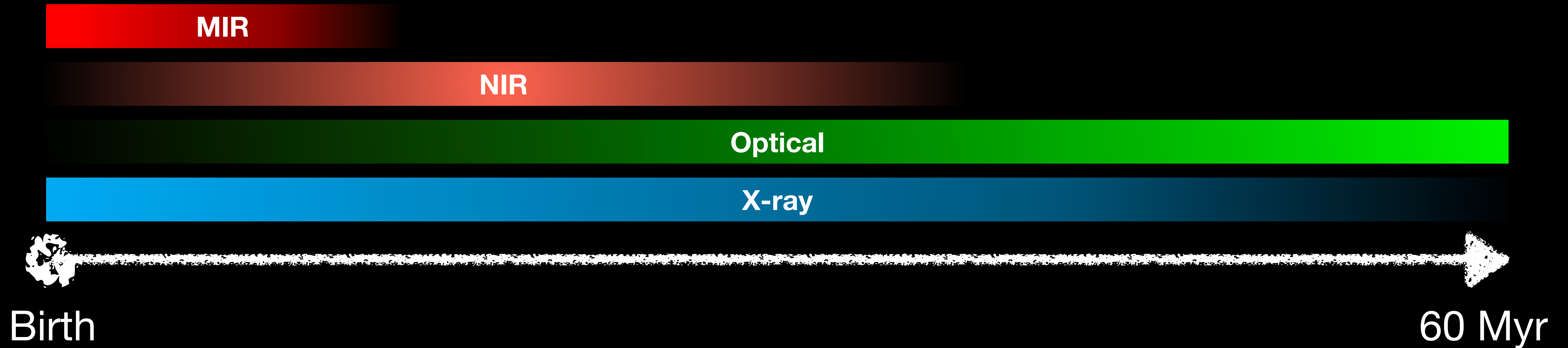
Often works focus  
on intersection



Aim to be truly multi-survey



# Untapped potential of current methods





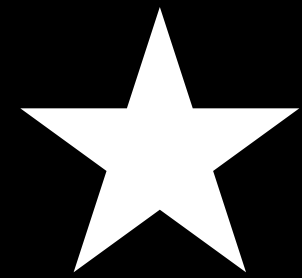
# Many different YSO tracers data sources

- **IR-excess** from disk
- **Lithium Depletion**: Li I ( $\lambda 6708\text{\AA}$ ) line
- **Gravity tracers** (weaker *logg* comp. to MS star)
  - Either med to high-R spectra or EW of absorption lines from *logg* (e.g., *Na I*, *K I*, *TiO*, *CaH 3*)
- **Emission lines from accretion** process
  - **H $\alpha$  line** - indicated also by increased chromospheric activity, or Balmer (Opt), Paschen or Brackett (IR), He I, Ca II near-IR triplet
- **Optical and IR veiling**
  - Excess flux from high-temperature material in inner disc regions
  - dilutes photospheric spectral lines & enhances spectral continuum
- **Variability**
  - Seen e.g. in larger mean flux uncertainties in multi-epoch photometry
- **Increased X-ray emission** & stellar rotation

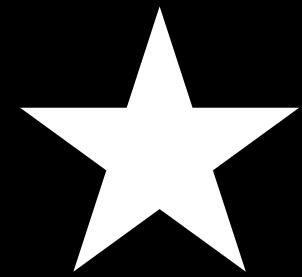


# Ways to improve ML YSO catalog

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



Gaia  
2MASS  
WISE  
LAMOST



WISE  
Spitzer  
APOGEE



# Ways to improve ML YSO catalog

- Data fusion: use as many informative data sets as possible
- Provide well-calibrated posteriors over stellar parameters given spectra & photometric observations



# Ways to improve ML YSO catalog

- Data fusion: use as many informative data sets as possible
  - Provide well-calibrated posteriors over stellar parameters given spectra & photometric observations
  - Scale inference to  $> 1\text{M} - 1\text{B}$  stars
- *Simulation-based inference for **incomplete, multi-survey** data*



# Model implementation

## I. SBI model

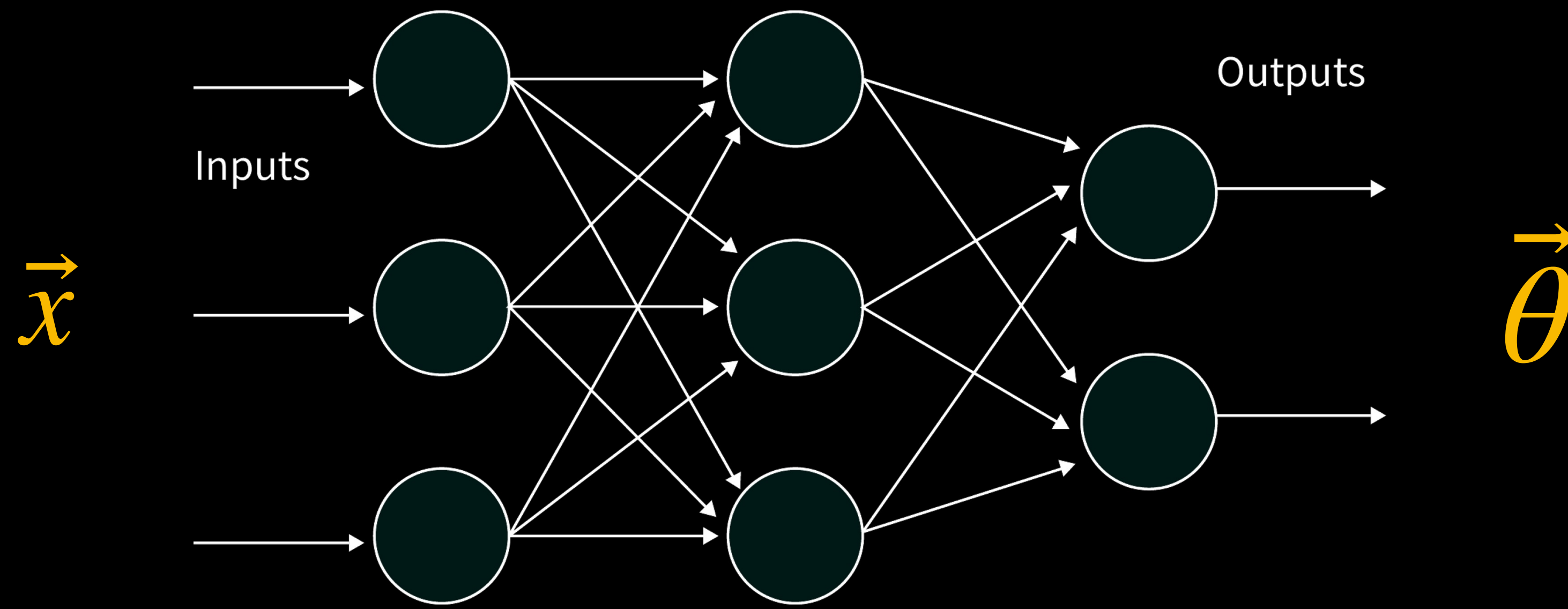


# Typical ML regression





# Typical ML regression

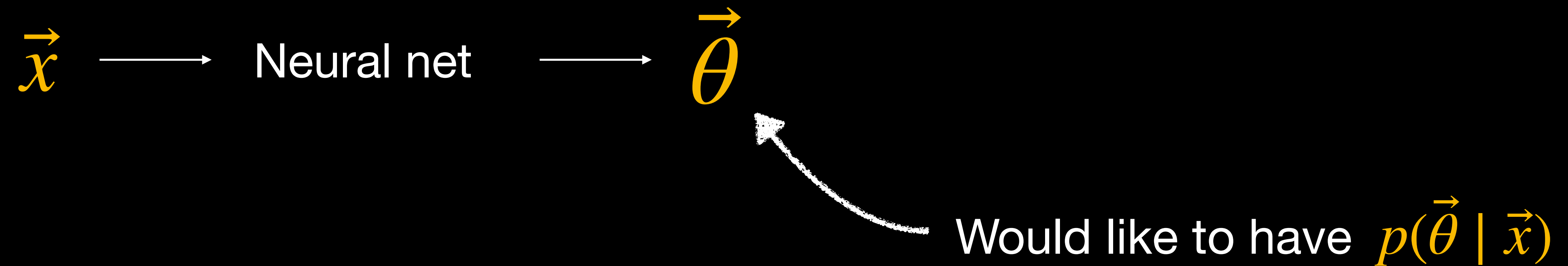


Series of learnable affine transformations of  $\vec{x}$   
followed by pointwise non-linear map:  $f_{\phi}(\vec{x}) = \hat{\theta}$

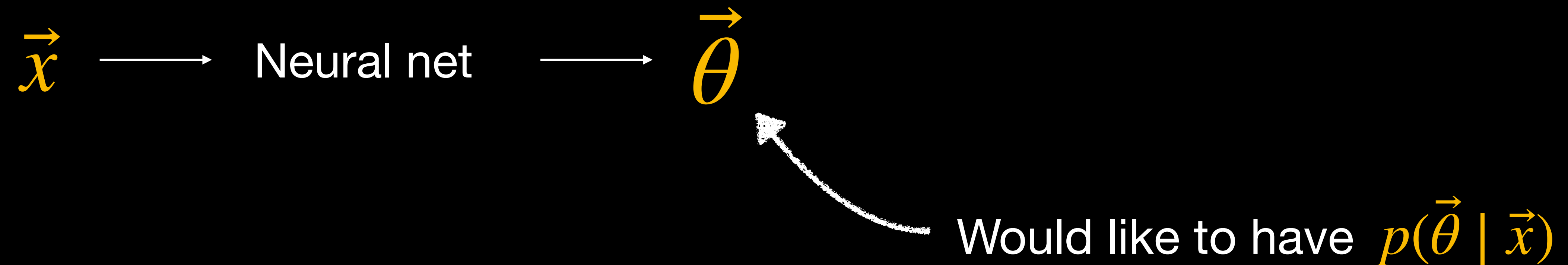
$\phi$ ...learnable parameters



# Typical ML regression



# Typical ML regression

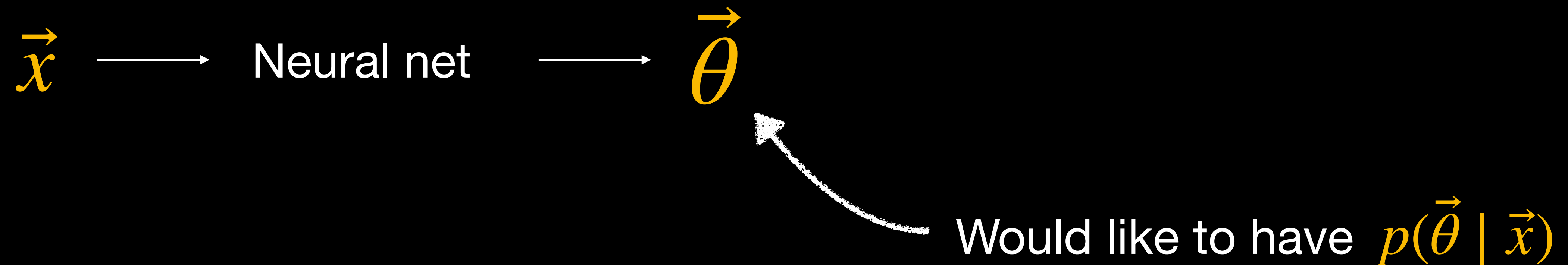


However:

- $p(\vec{x} | \vec{\theta})$  might not be tractable
- $p(\vec{\theta} | \vec{x})$  might not scale to millions - billions of “runs”



# Typical ML regression

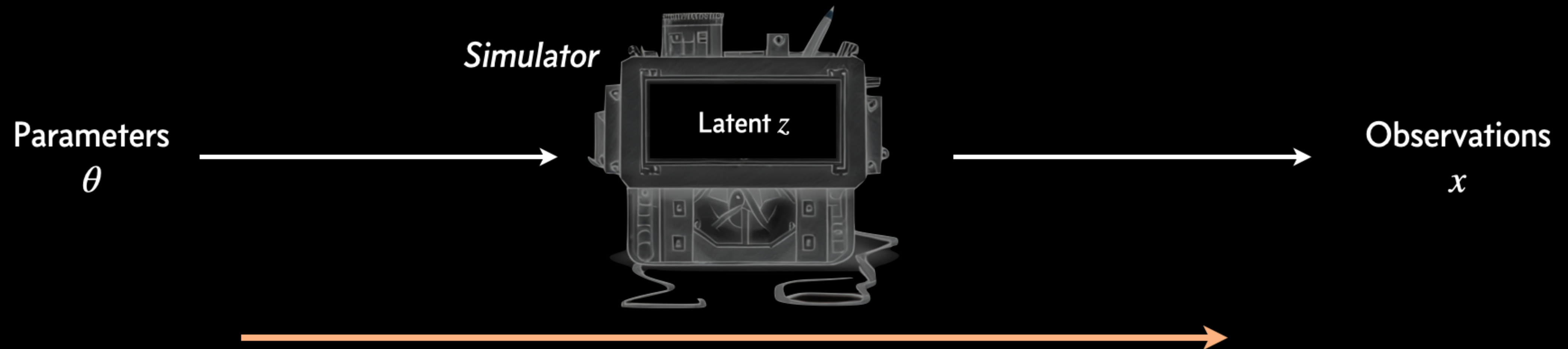


However:

- $p(\vec{x} | \vec{\theta})$  might not be tractable
- $p(\vec{\theta} | \vec{x})$  might not scale to millions - billions of “runs”

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

# Simulation based inference (SBI) setup



Prediction:

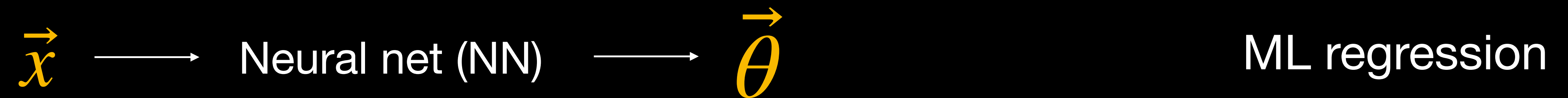
- Mechanistic forward model
- We can generate samples from a simulator  $x \sim p(x | \theta)$

Inference:

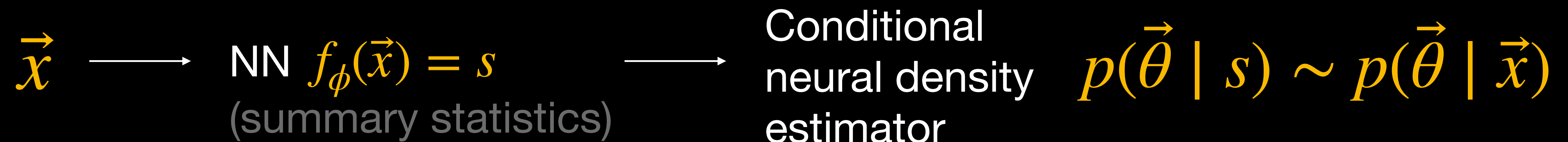
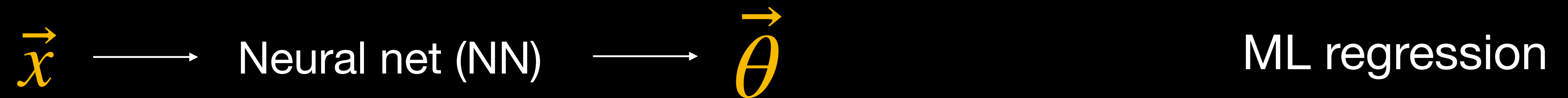
- Likelihood  $p(x | \theta) = \int dz p(x, z | \theta)$  is intractable
- *Inference is challenging*



# Neural posterior estimation

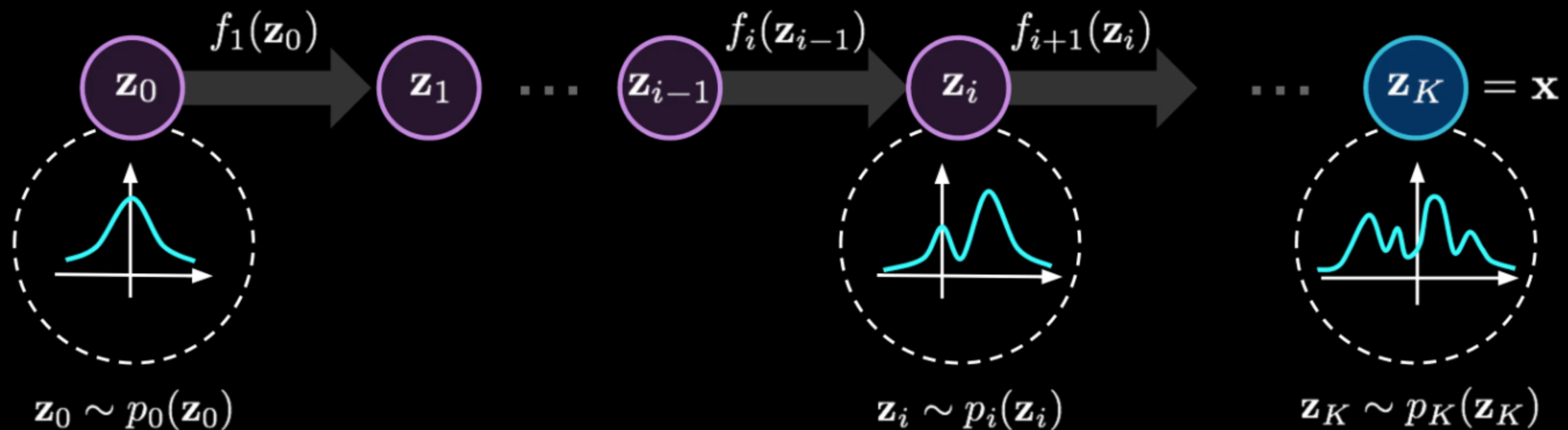
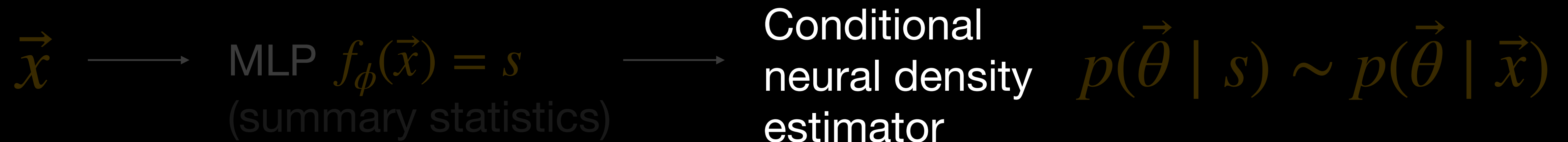


# Neural posterior estimation





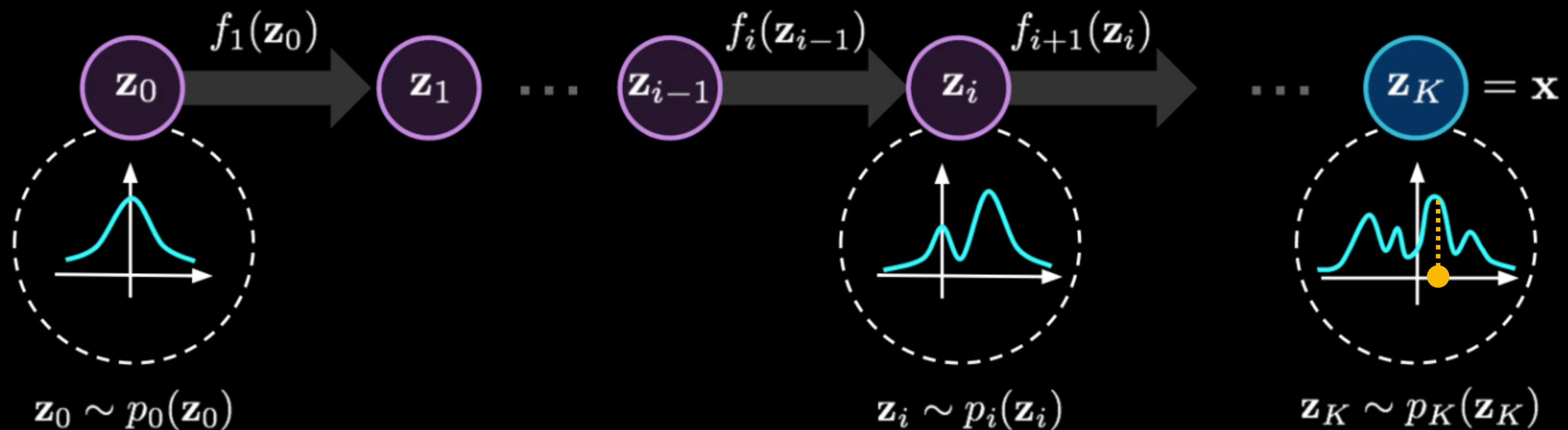
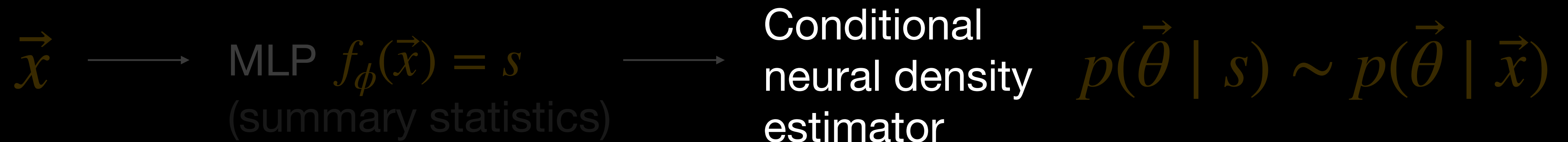
# Normalizing flows



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

Training objective: **maximum likelihood**

# Normalizing flows

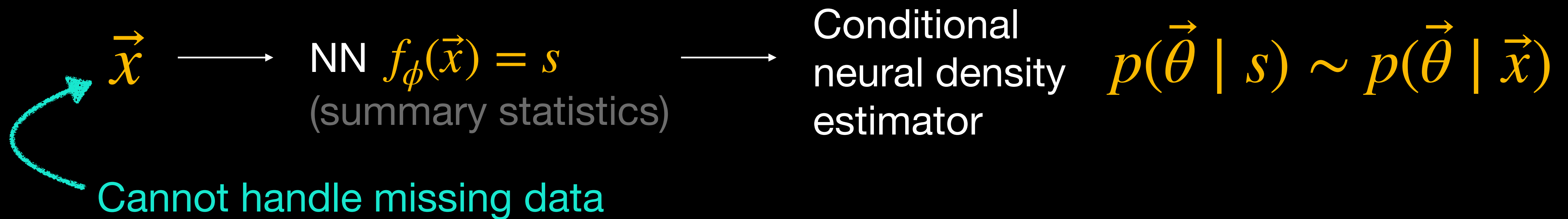
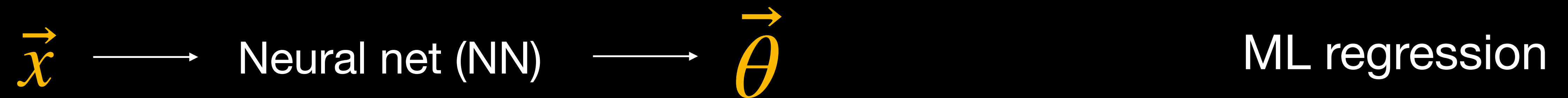


Parameterized, invertible maps  $f_i$  that transform Gaussian into target distribution

Training objective: **maximum likelihood**



# Neural posterior estimation



# Transformer: learning with incomplete data

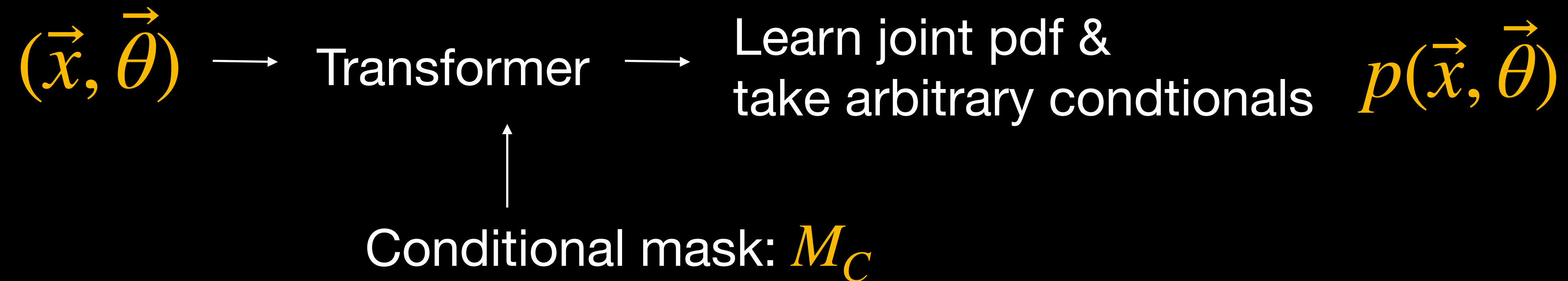
$(\vec{x}, \vec{\theta}) \longrightarrow$  Transformer

↑  
Conditional mask:  $M_C$

Model trained to allow for arbitrary conditioning

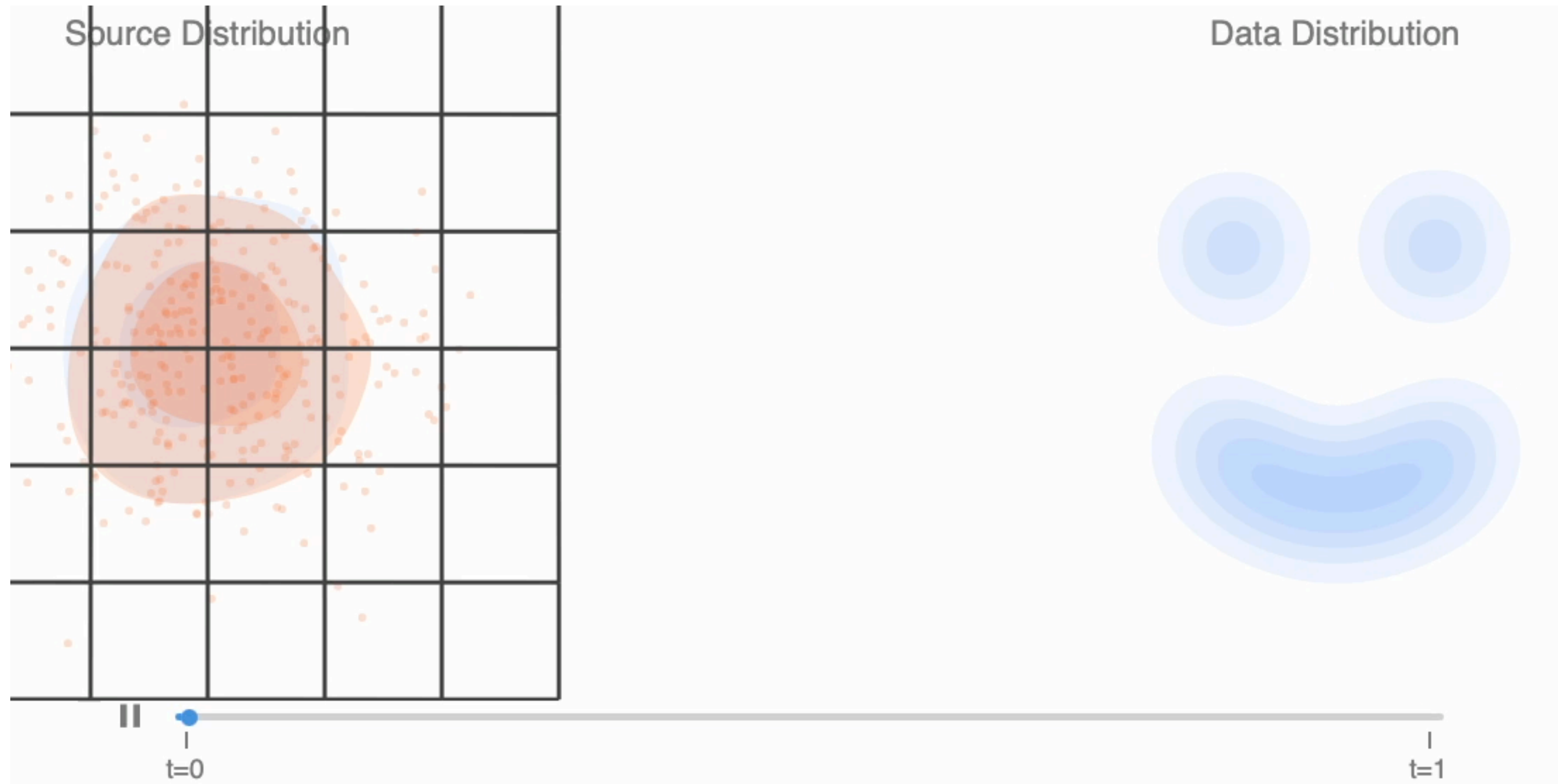


# Transformer: learning with incomplete data



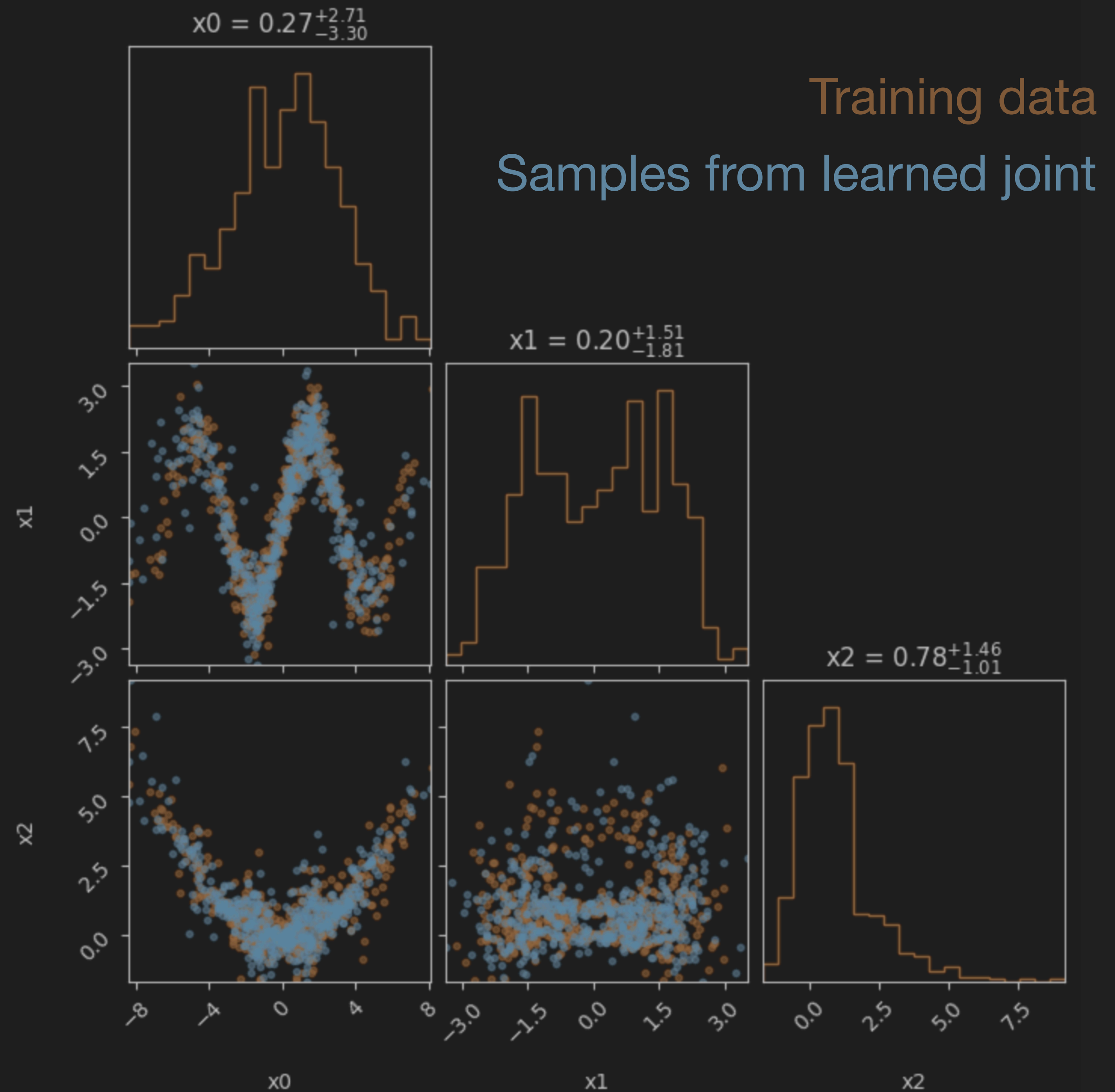
Model trained to allow for arbitrary conditioning  
— condition joint only on actually observed quantities

# Flow matching: time-dependent vector field

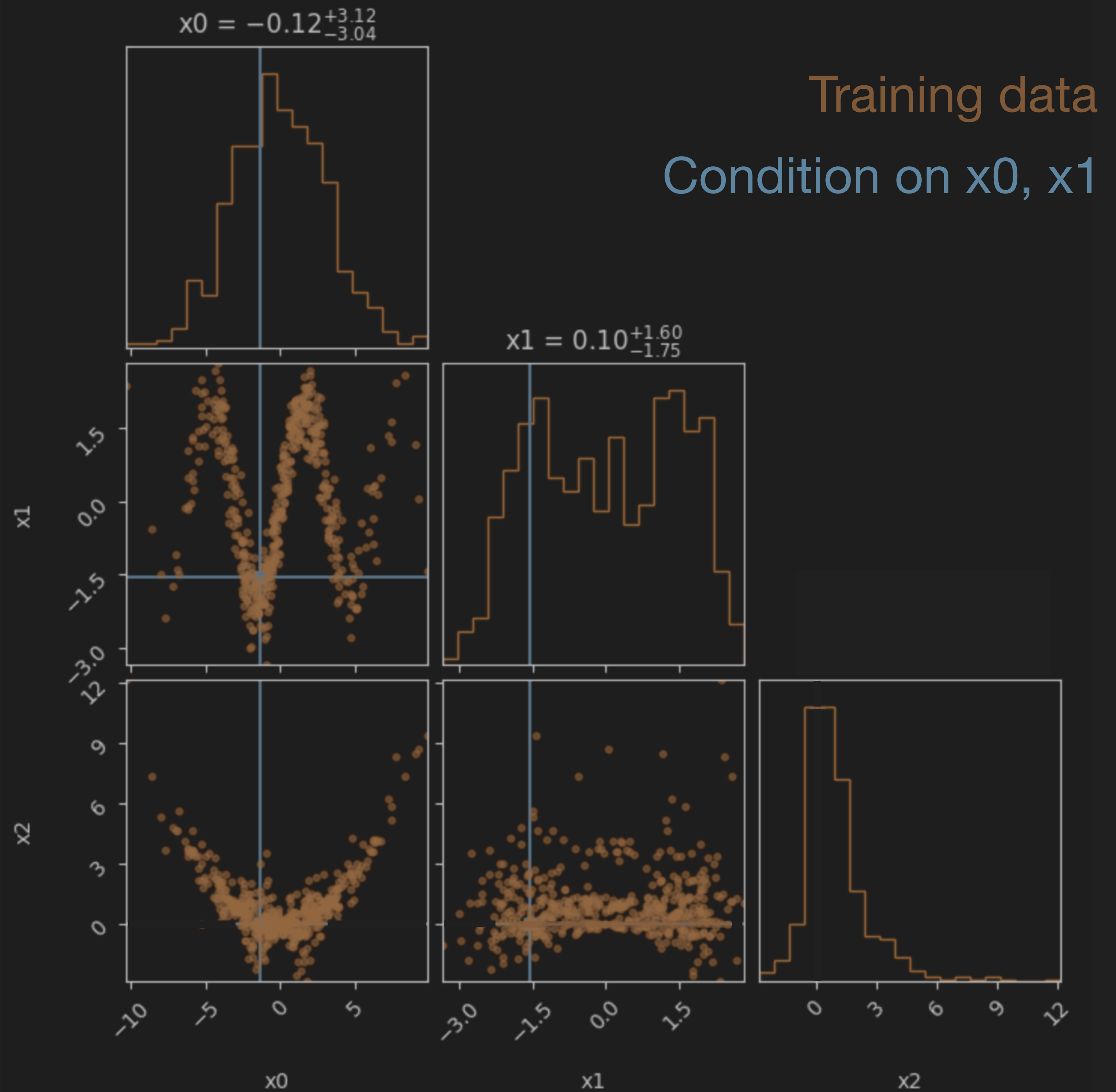
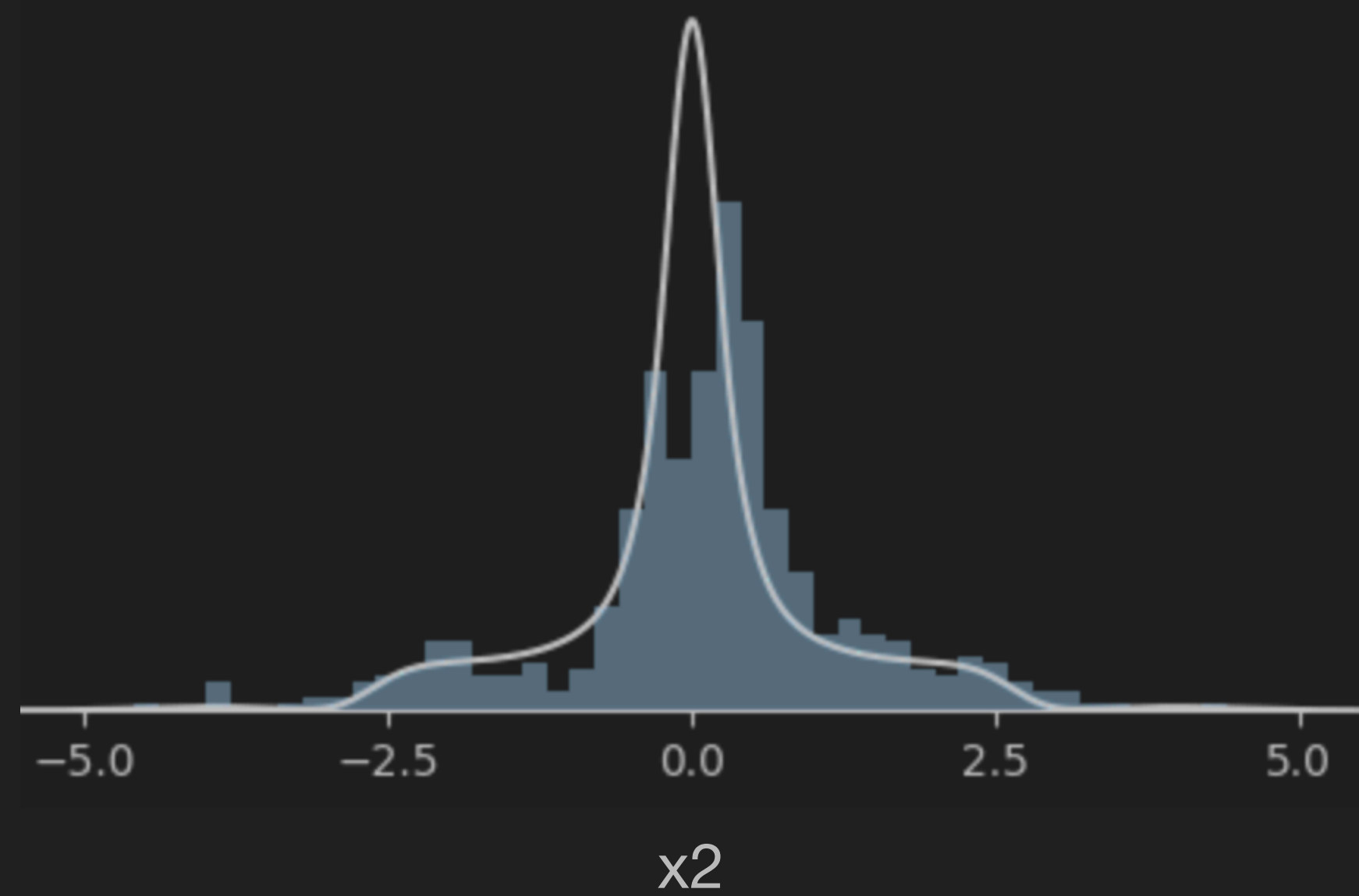




# 3D toy example



# 3D toy example





# 3D toy example

Flow not well constrained in low data regime



# 3D toy example

Flow not well constrained in low data regime

**Guidance** term towards physical manifold in low data regime





# 3D toy example

Flow not well constrained in low data regime

**Guidance** term towards physical manifold in low data regime

Baldan+2025; Bastek+2025

$$\mathcal{L} = w_{\text{FM}} \mathcal{L}_{\text{FM}} + w_{\text{R}} \mathcal{L}_{\text{R}}$$



# 3D toy example

Flow not well constrained in low data regime

**Guidance** term towards physical manifold in low data regime

Baldan+2025; Bastek+2025

$$\mathcal{L} = w_{\text{FM}} \mathcal{L}_{\text{FM}} + w_{\text{R}} \mathcal{L}_{\text{R}}$$

$$\mathcal{L}_{\text{R}} = f(x_1, \dots, x_n)$$

E.g.:  $f(x_1, \dots, x_n) = x_1^2 + x_2^2 + x_3^2$





# Guidance via measurement operator constraints

$$\mathcal{L}_R = f(x_1, \dots, x_n)$$

The diagram illustrates the measurement operator equation, which relates the observed flux in a specific band to the intrinsic spectral energy distribution (SED) of a source, accounting for various physical and observational factors. The equation is:

$$f_i = a(d, R_\star) \int k_i(\lambda) A(\lambda; A_V, R_V) s(\lambda; T_{\text{eff}}, \log g, Z, t_\star) d \log \lambda$$

The components of the equation are labeled as follows:

- Amplitude**: Points to the term  $a(d, R_\star)$ .
- Attenuation**: Points to the term  $A(\lambda; A_V, R_V)$ .
- Flux in band i**: Points to the variable  $f_i$ .
- Passband filter**: Points to the term  $k_i(\lambda)$ .
- intrinsic surface SED**: Points to the term  $s(\lambda; T_{\text{eff}}, \log g, Z, t_\star)$ .

# Guidance via measurement operator constraints

$$\mathcal{L}_R = f(x_1, \dots, x_n)$$

$$f_i = a(d, R_\star) \int k_i(\lambda) A(\lambda; A_V, R_V) s(\lambda; T_{\text{eff}}, \text{logg}, Z, t_\star) d \log \lambda$$



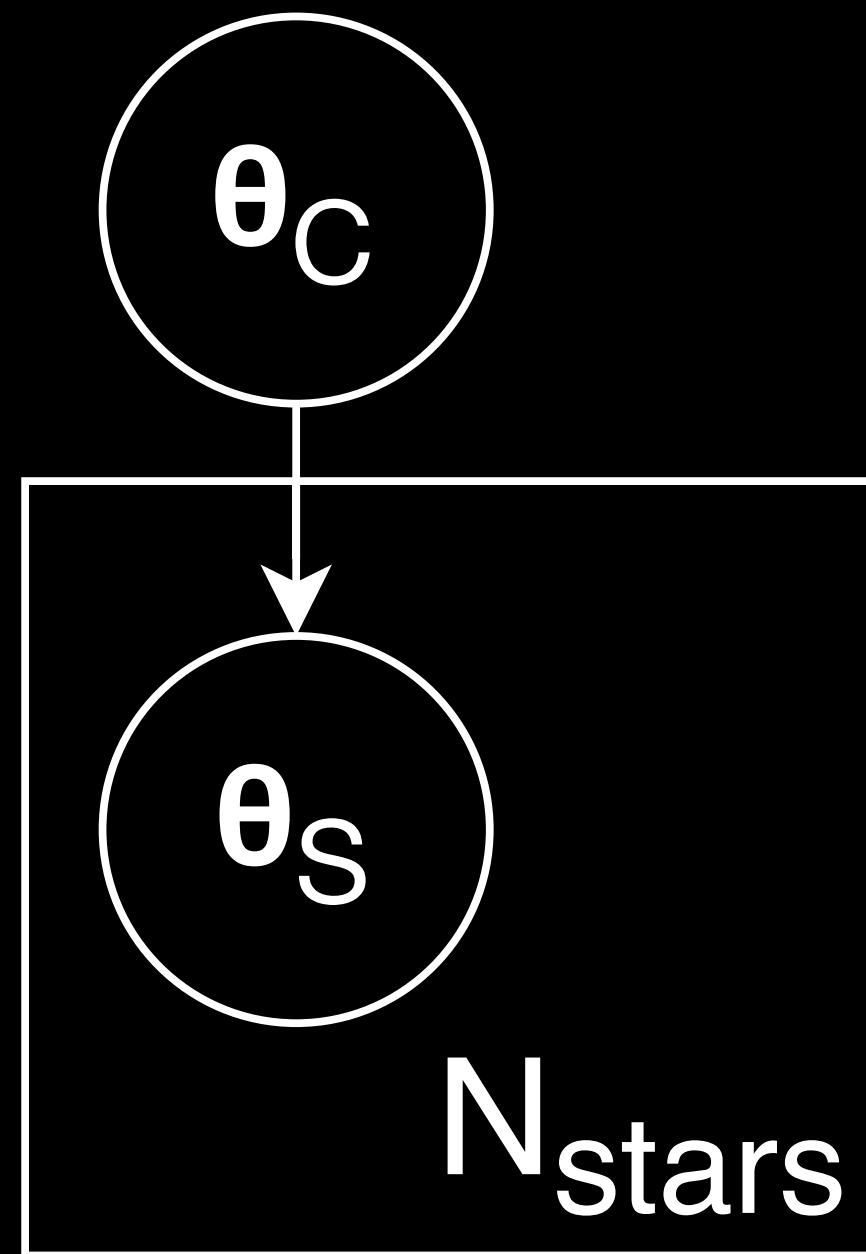
$$[\mu(\lambda) + \sum_m c_m \phi_m(\lambda)]$$

$$p(\vec{x}, \vec{\theta}, \vec{c})$$



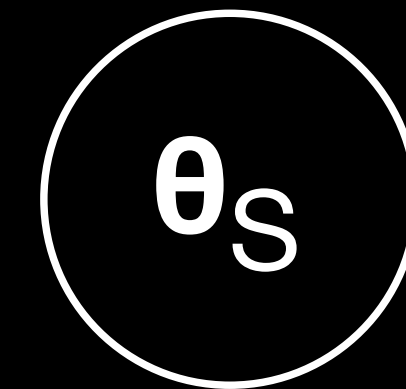
**Pilot study**

## Clusters/young stars



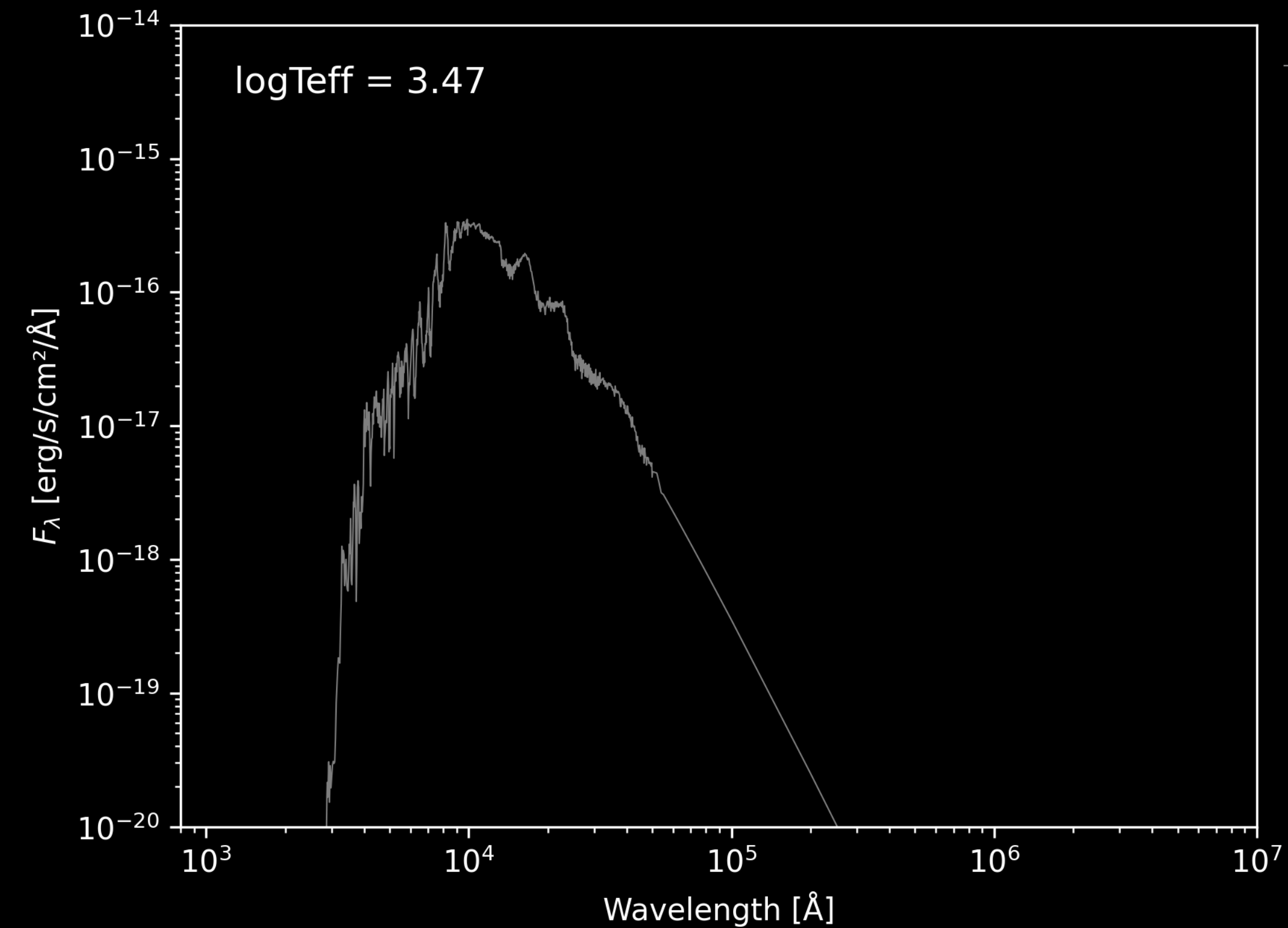
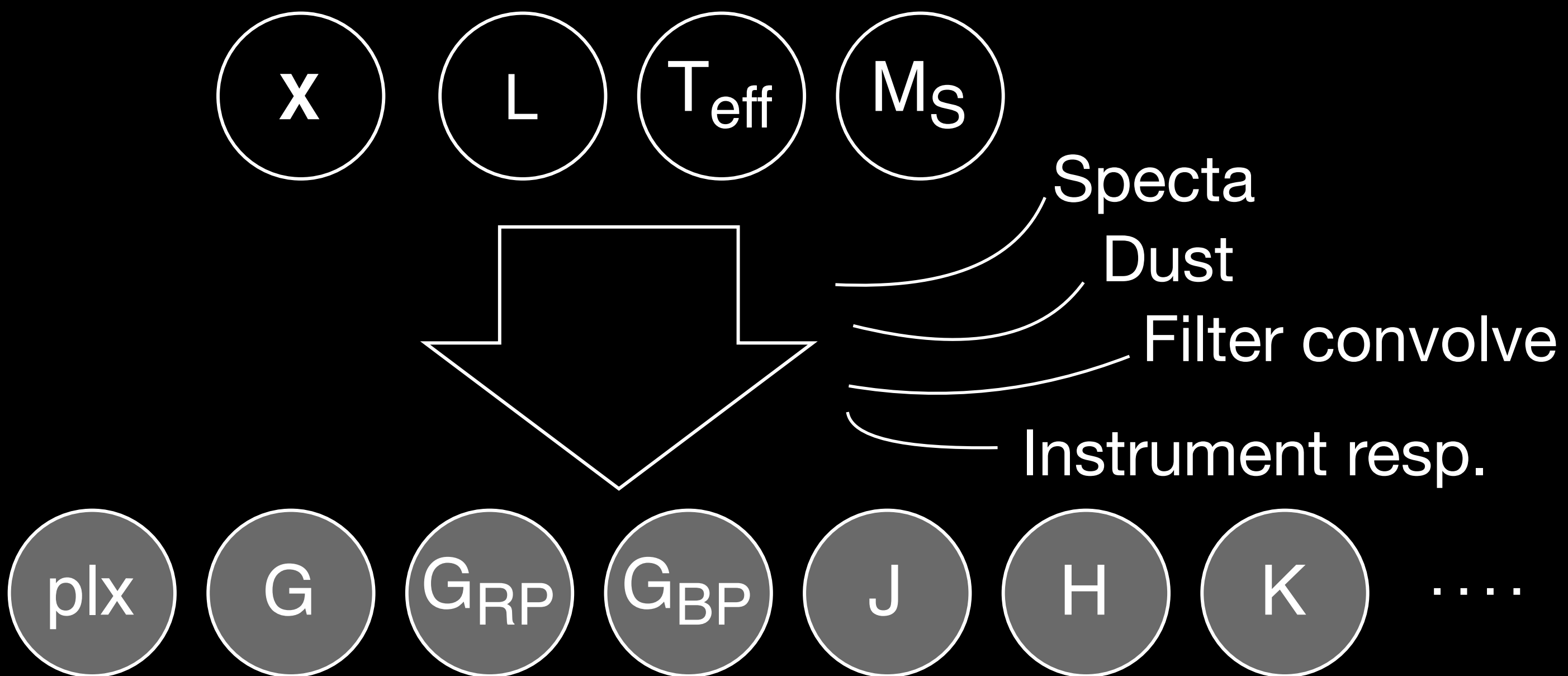
## Milky Way model

- *Galaxia* code
  - Galactic structure + kinematics, SFH + chemistry evolution
- thin+thick+halo+bulge

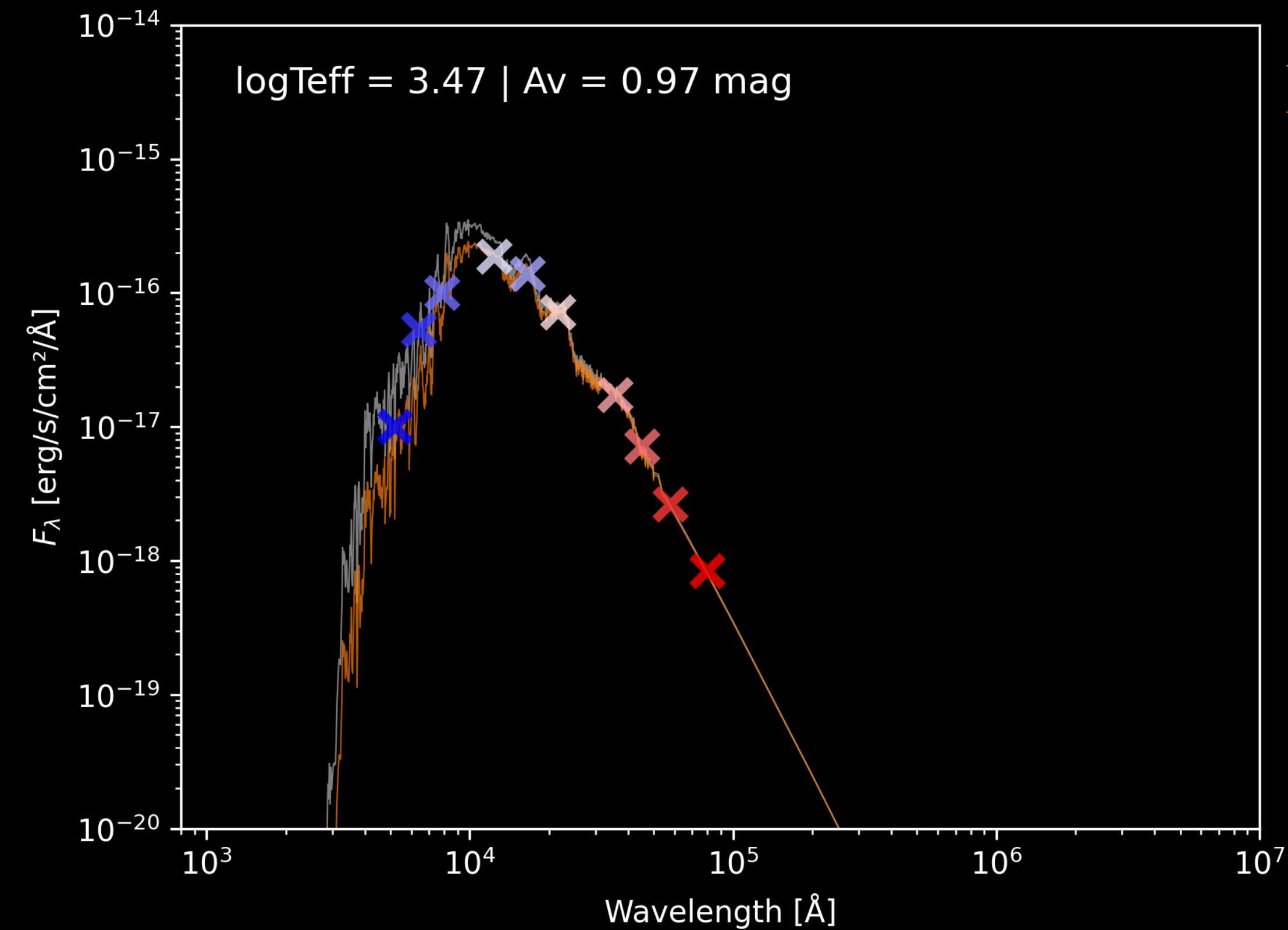
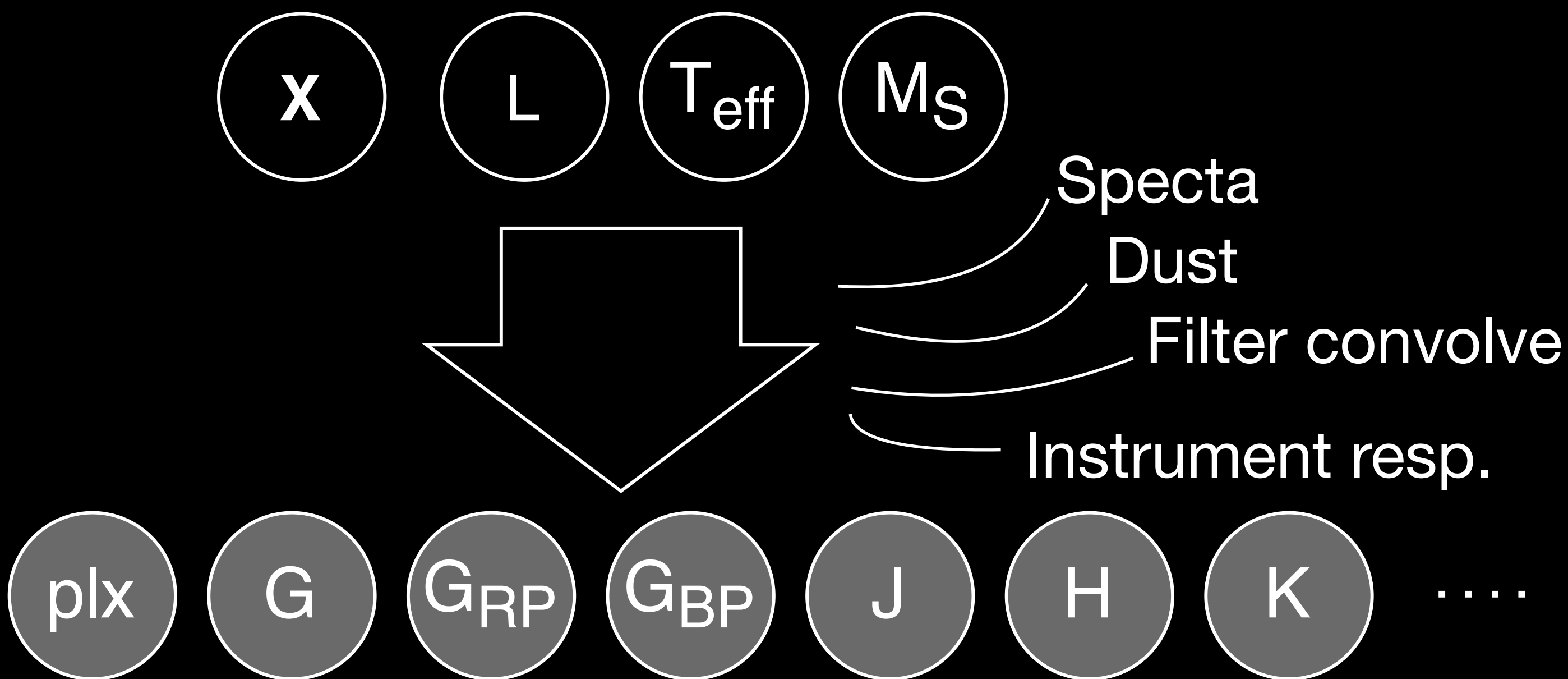




# Fwd model



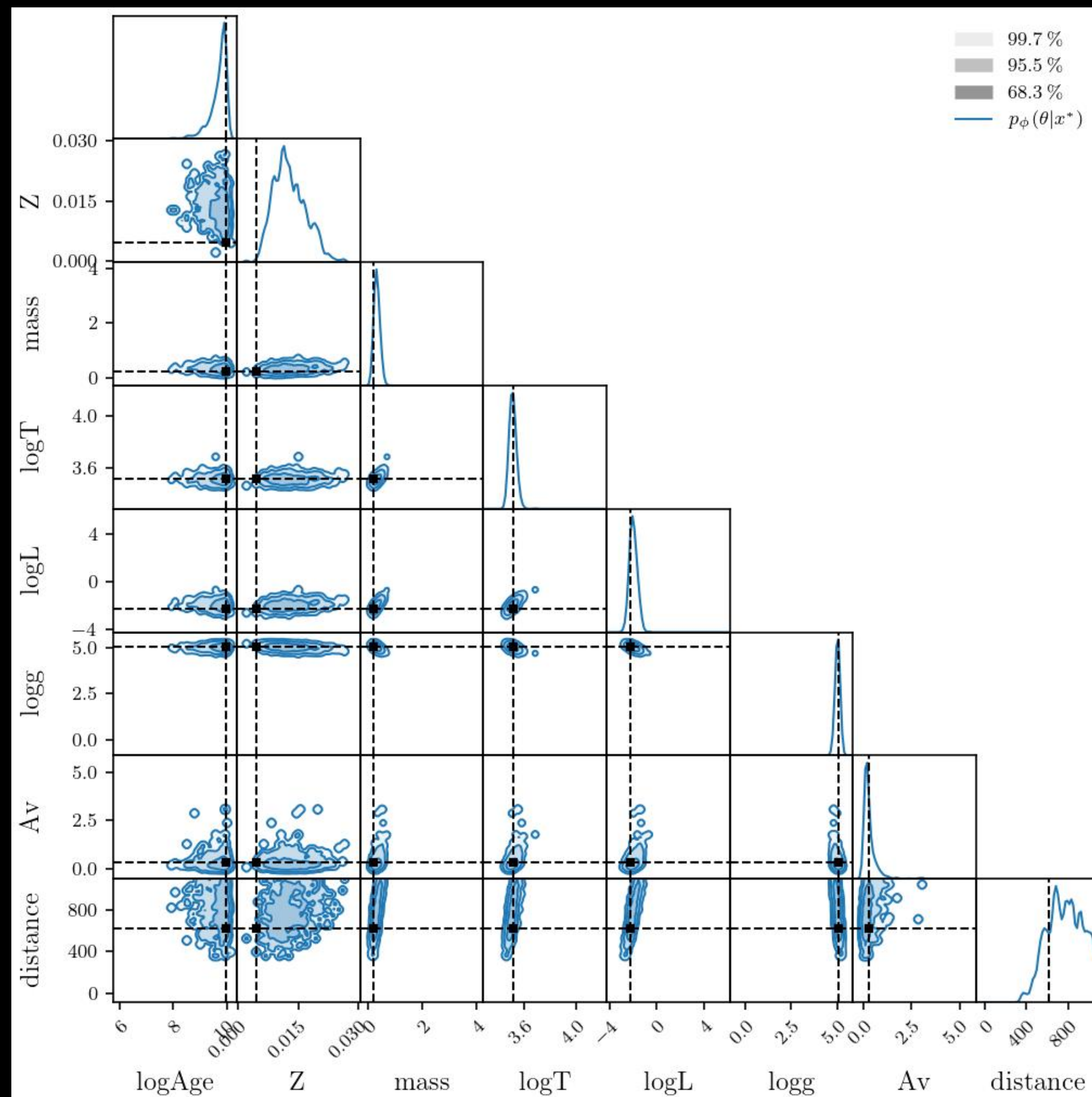
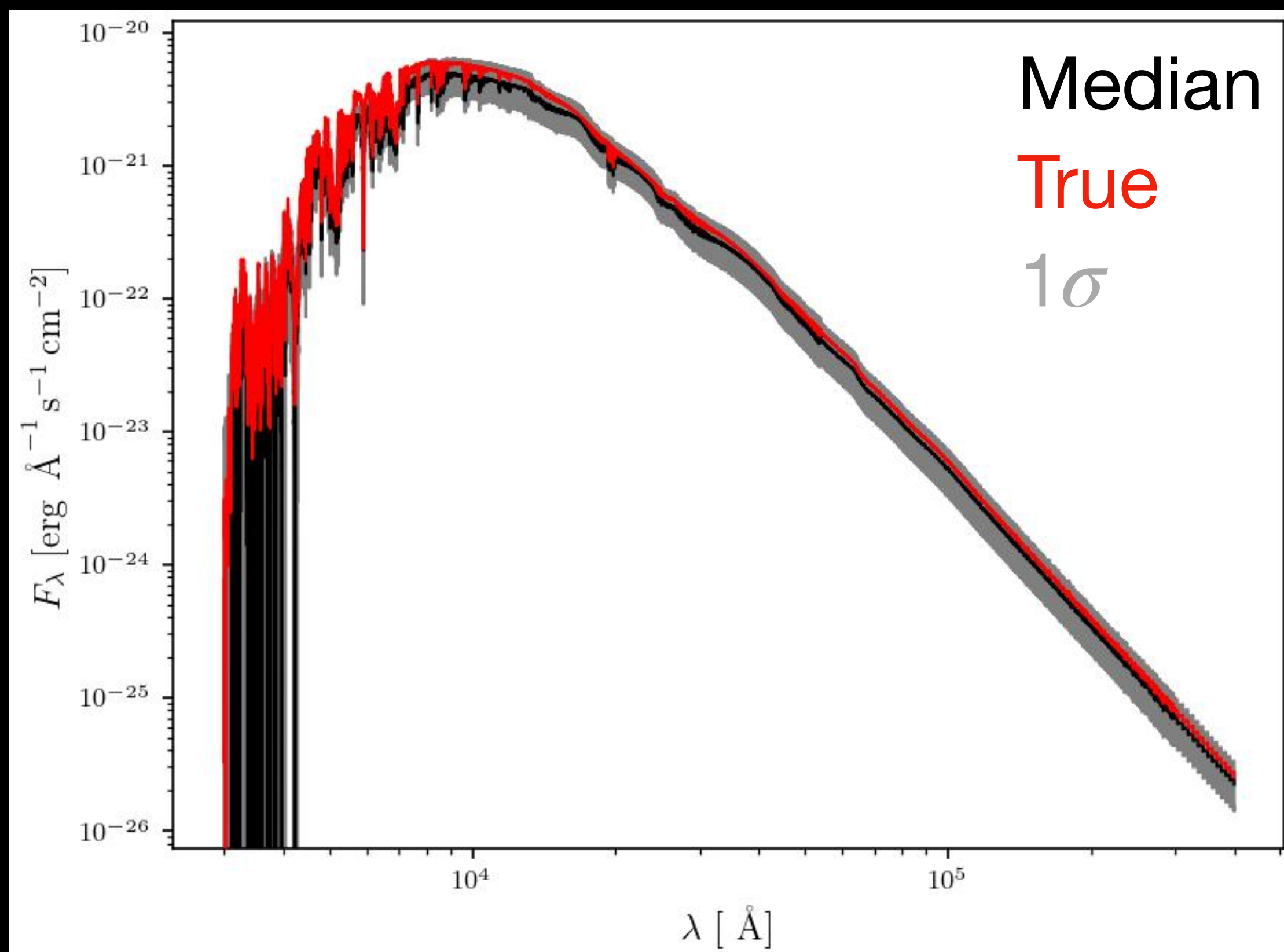
# Fwd model





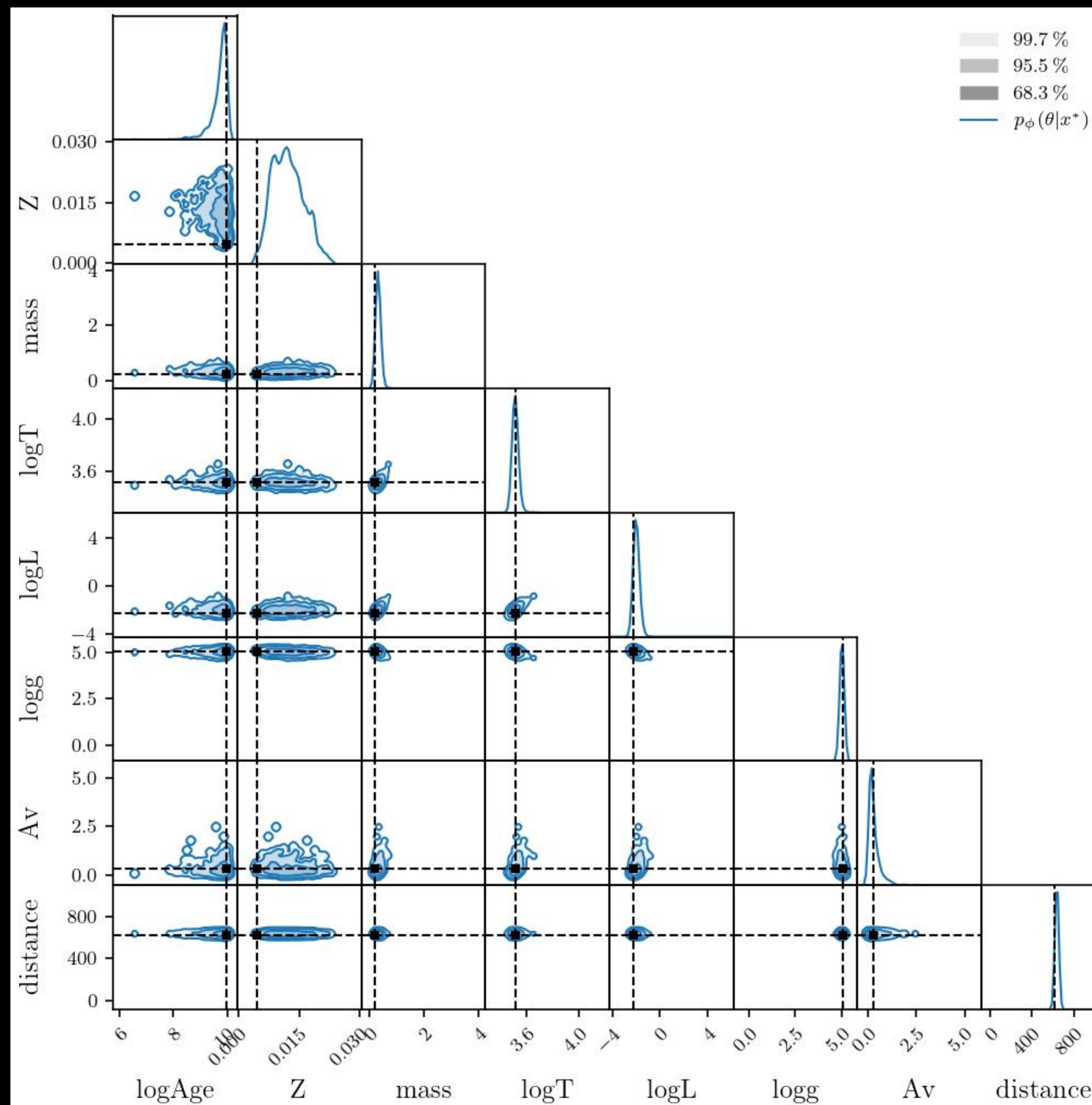
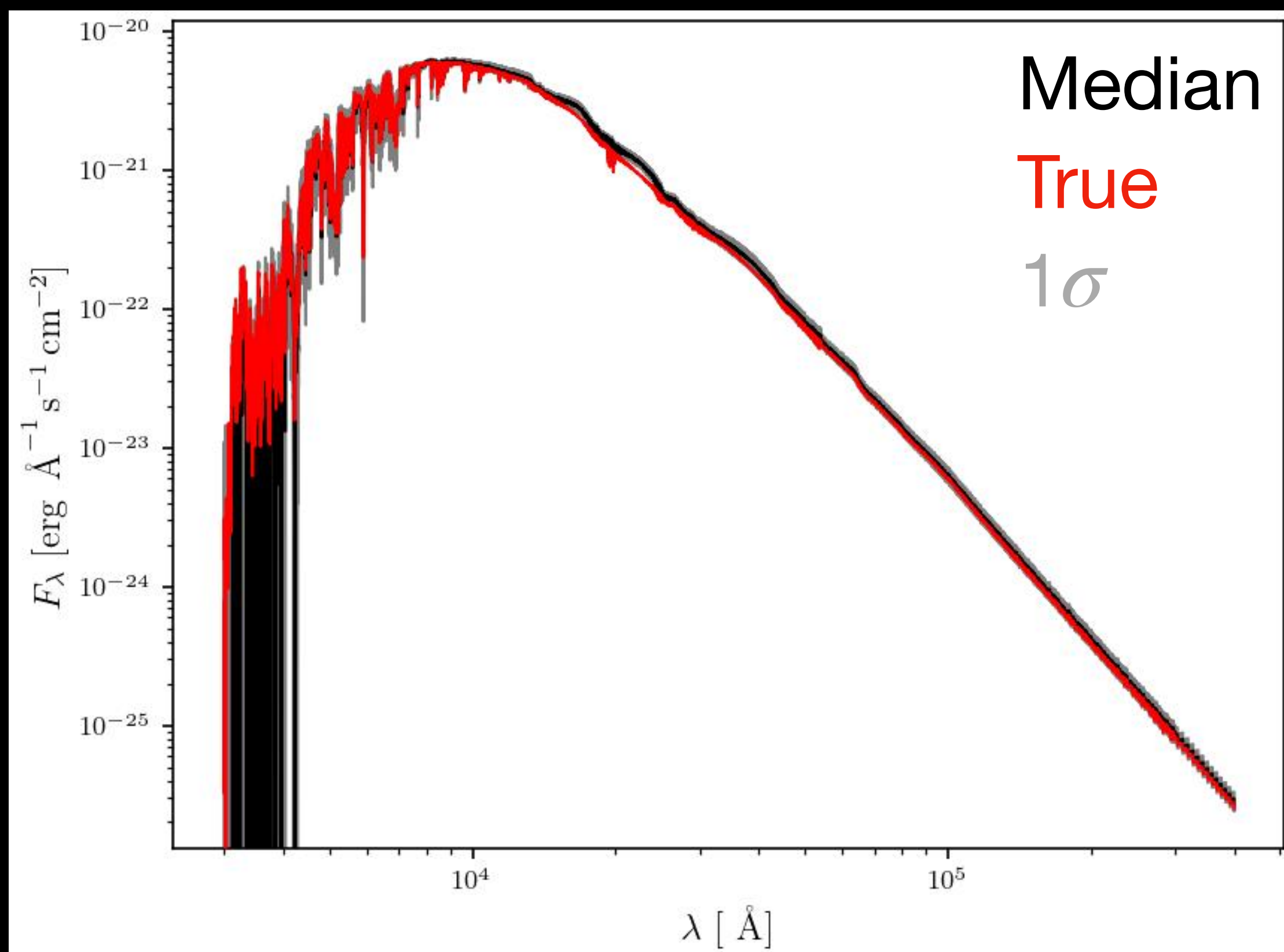
**Example: low mass star**

# Gaia + 2MASS + low SNR parallax

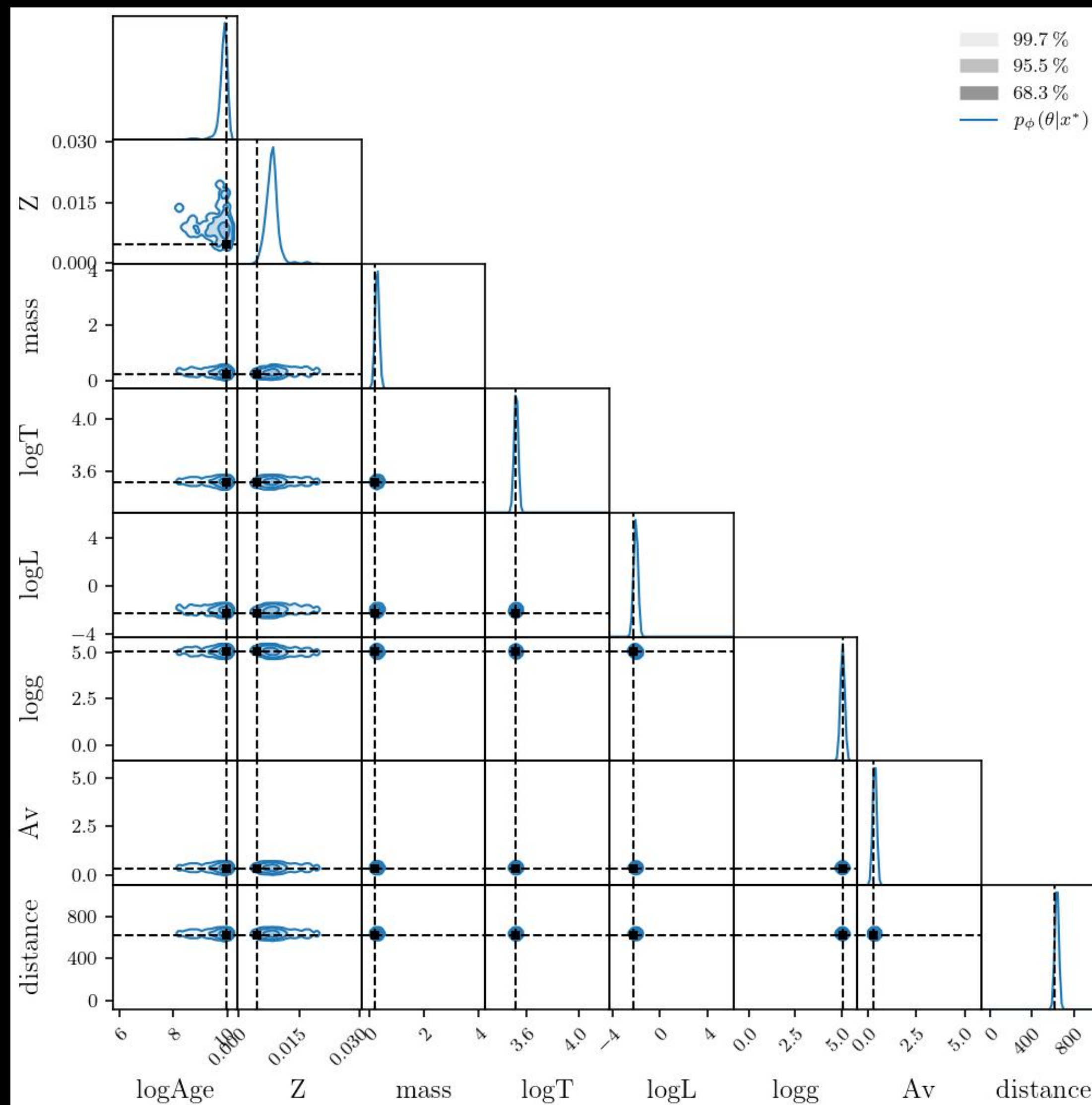
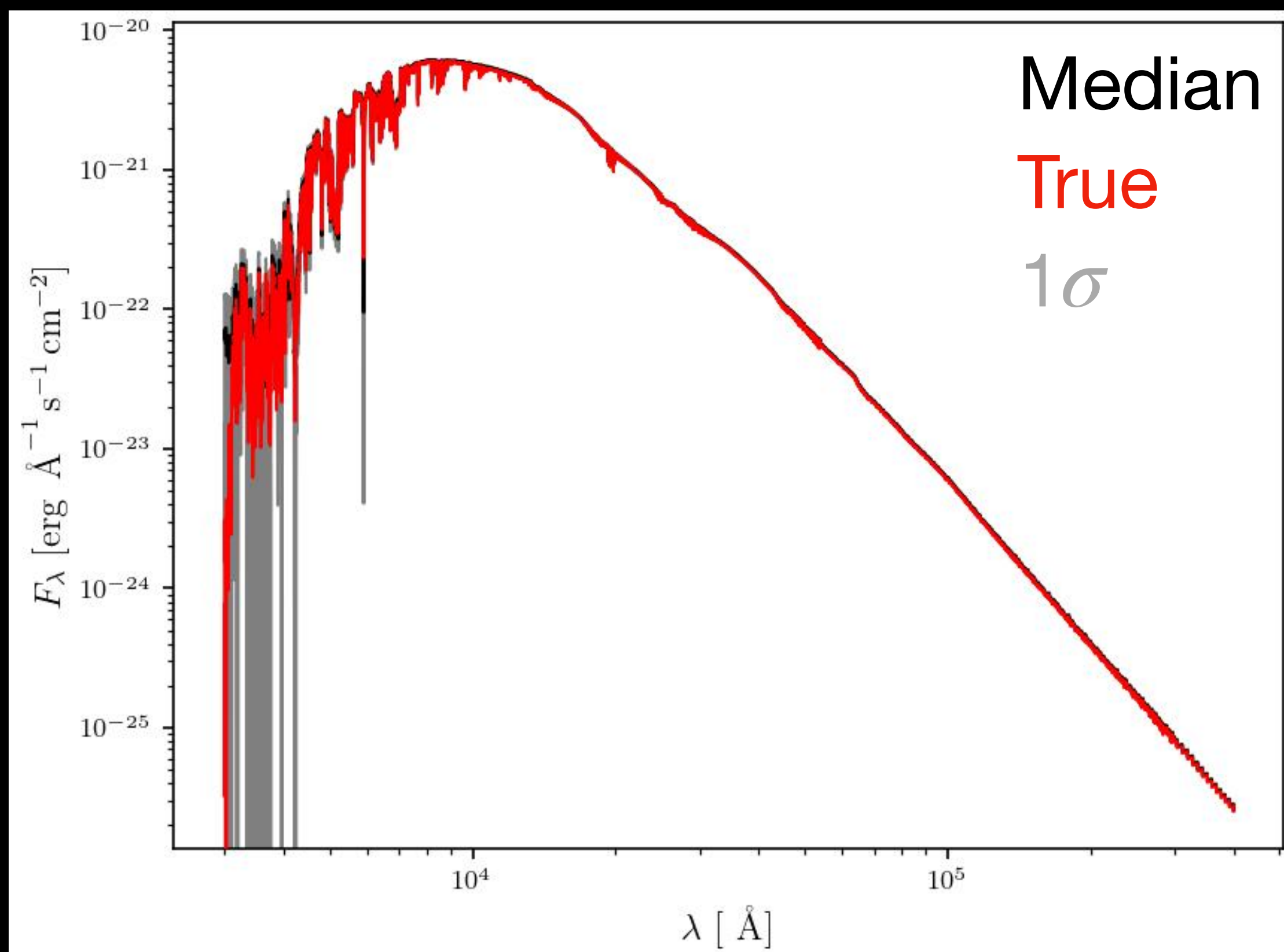




# Gaia + 2MASS + high SNR parallax



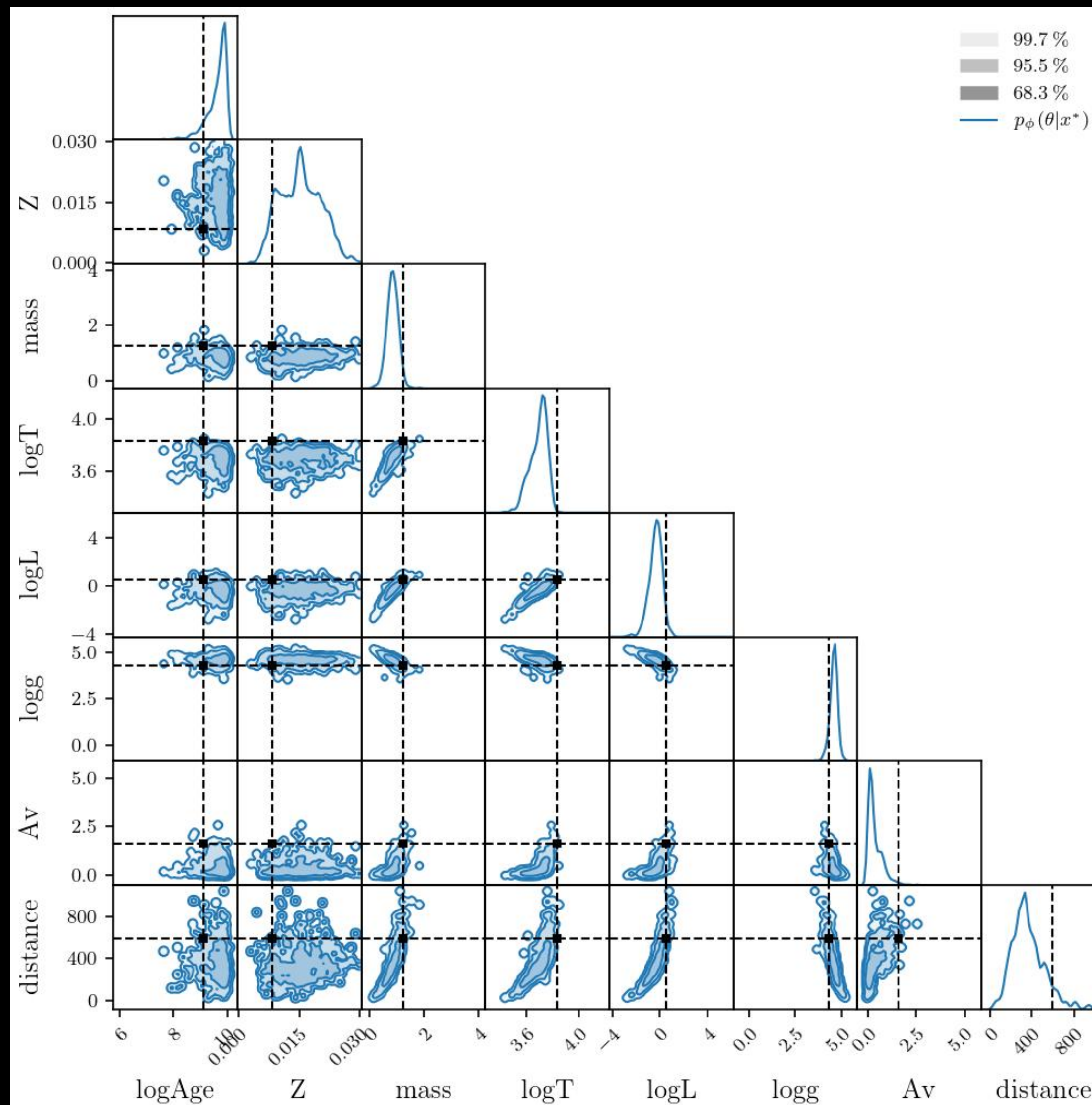
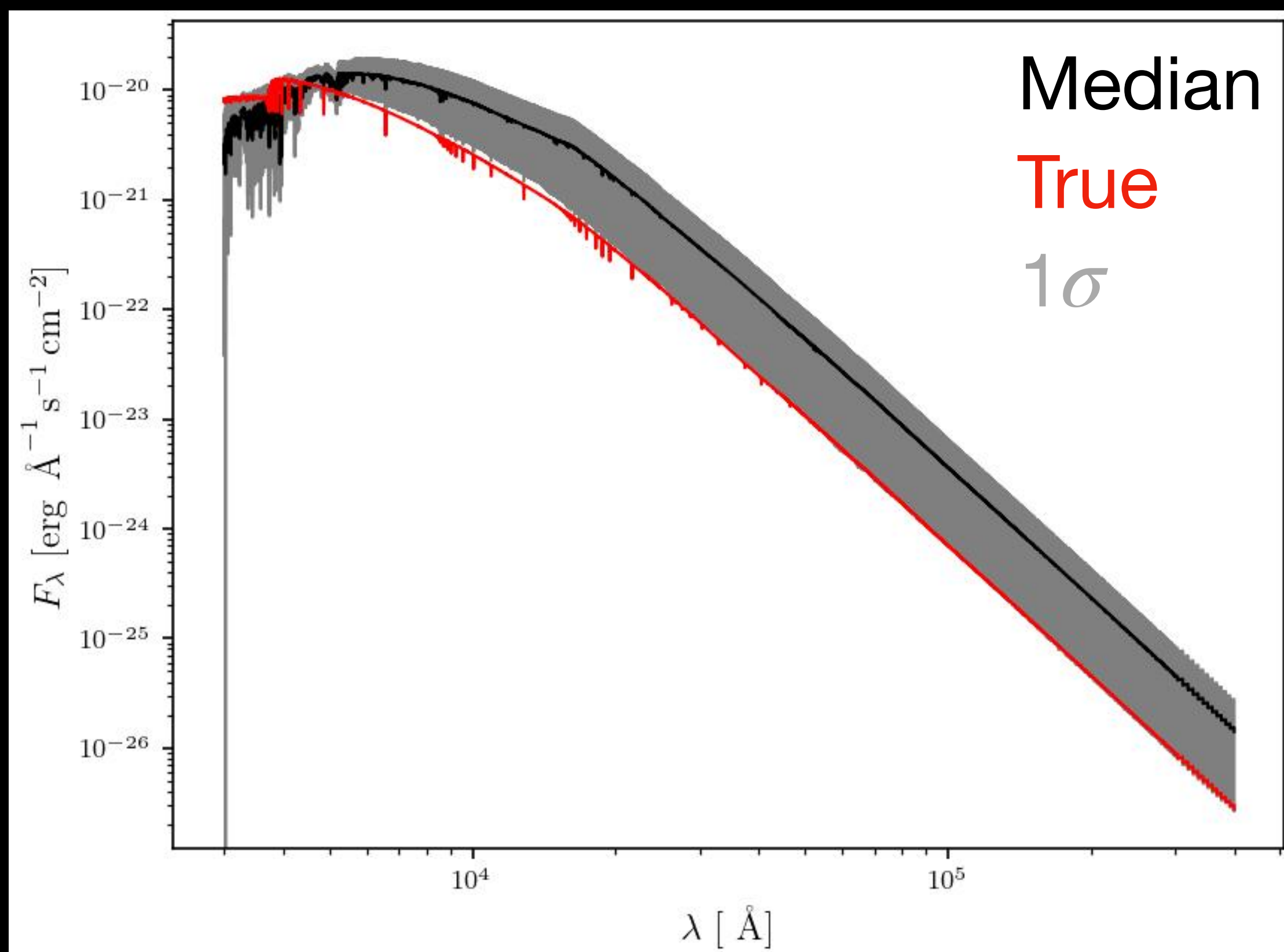
# Gaia + 2MASS + high SNR parallax + LAMOST





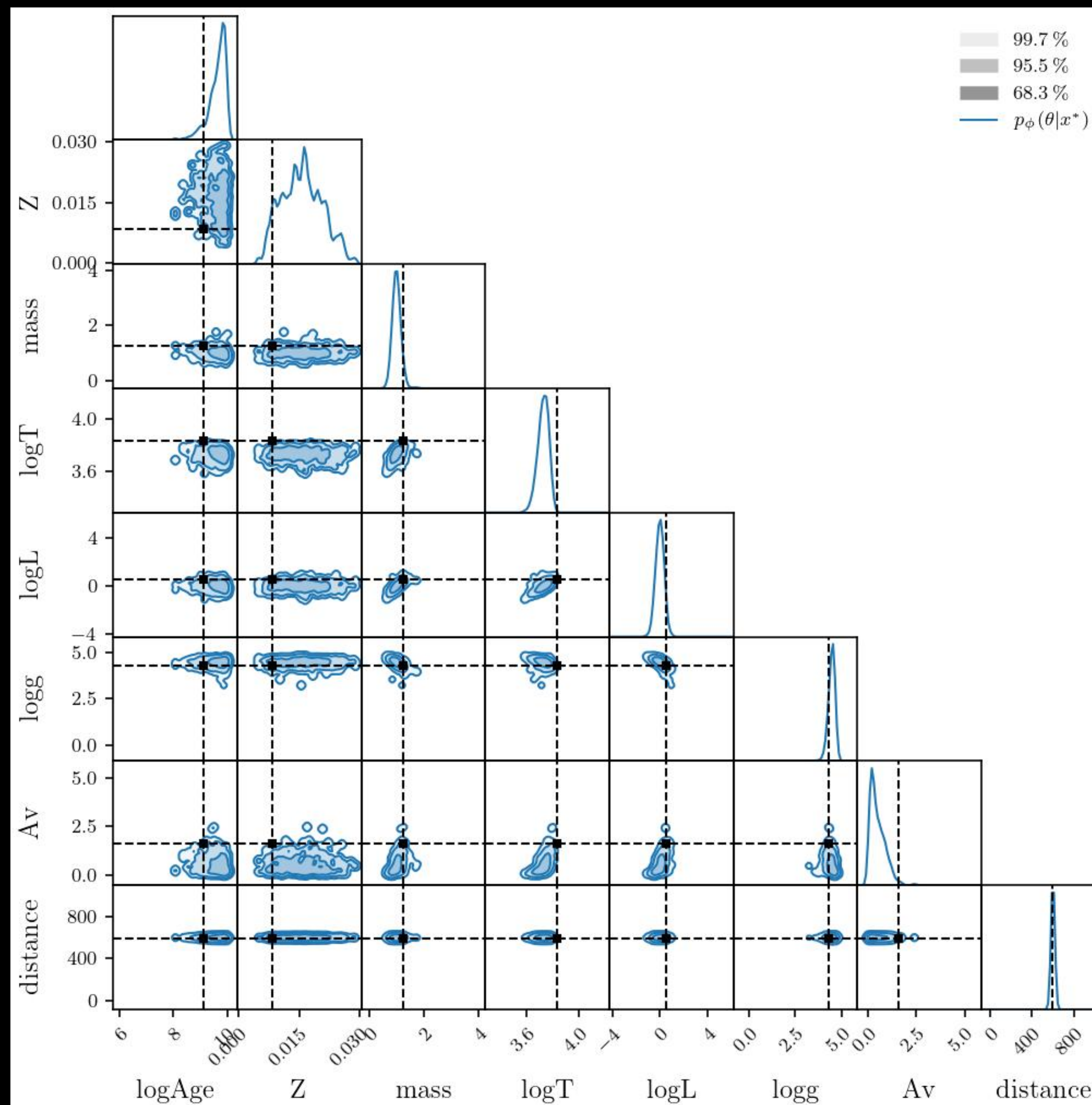
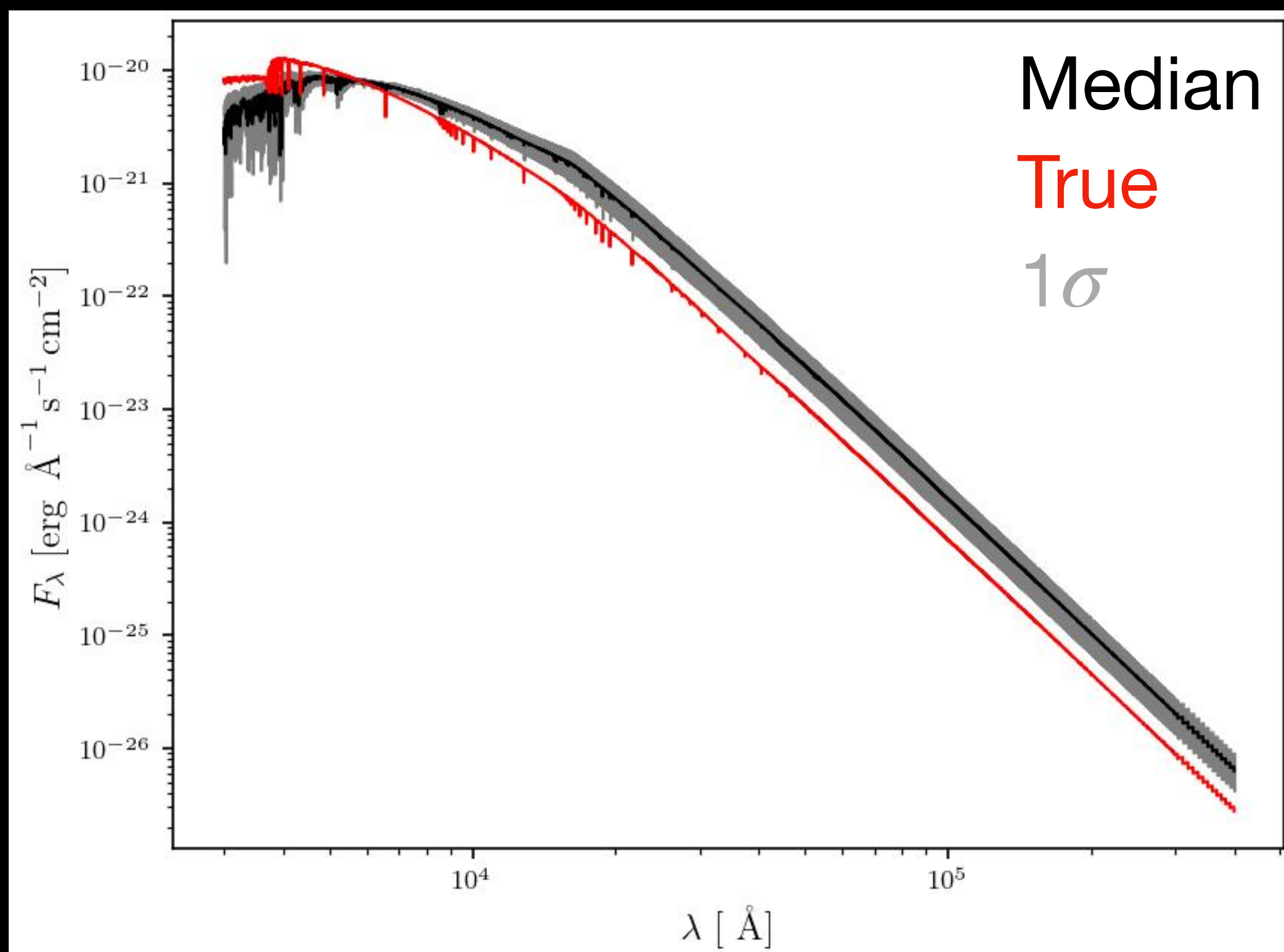
**Example: high mass star**

# Gaia + 2MASS + low SNR parallax

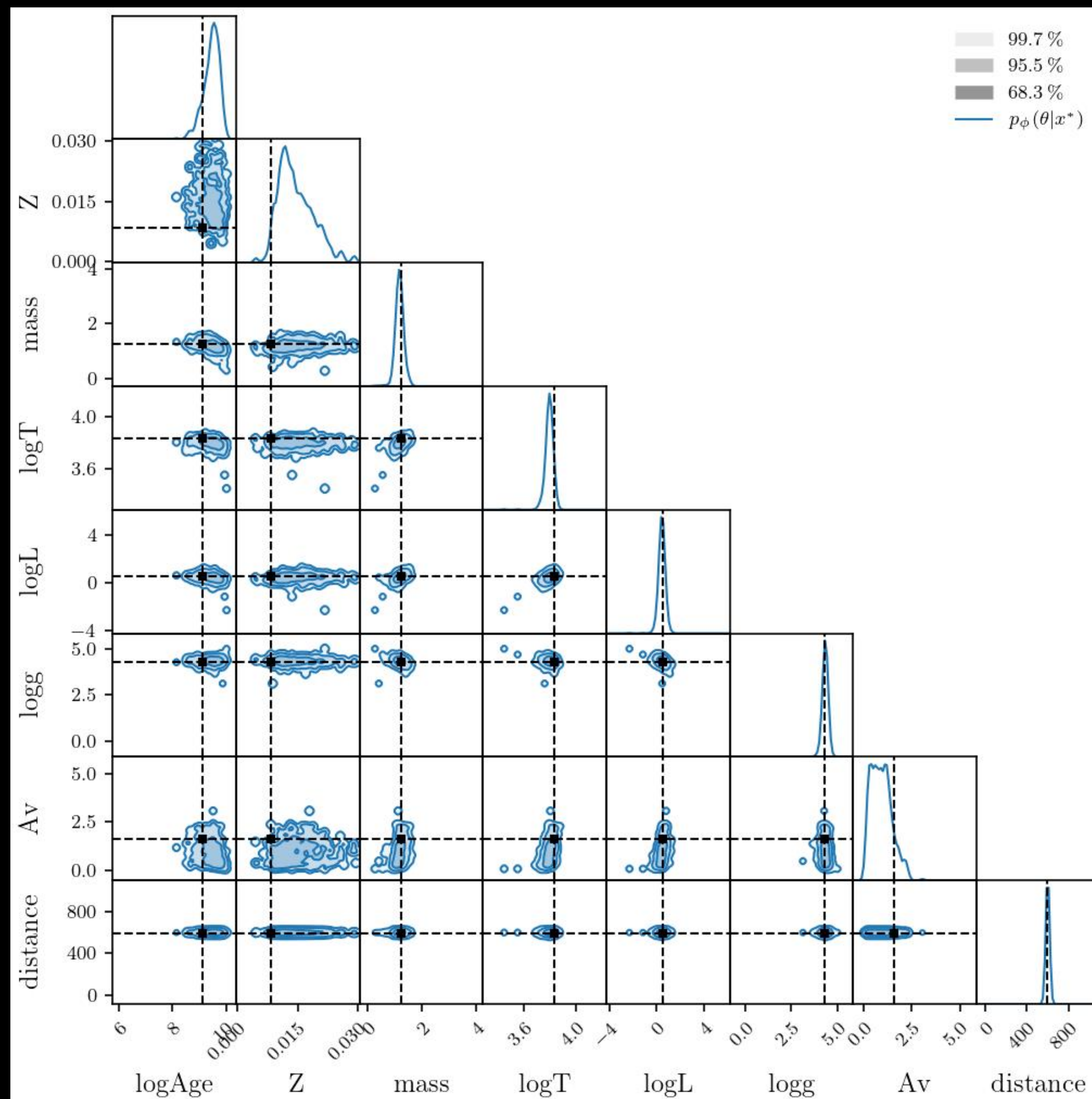
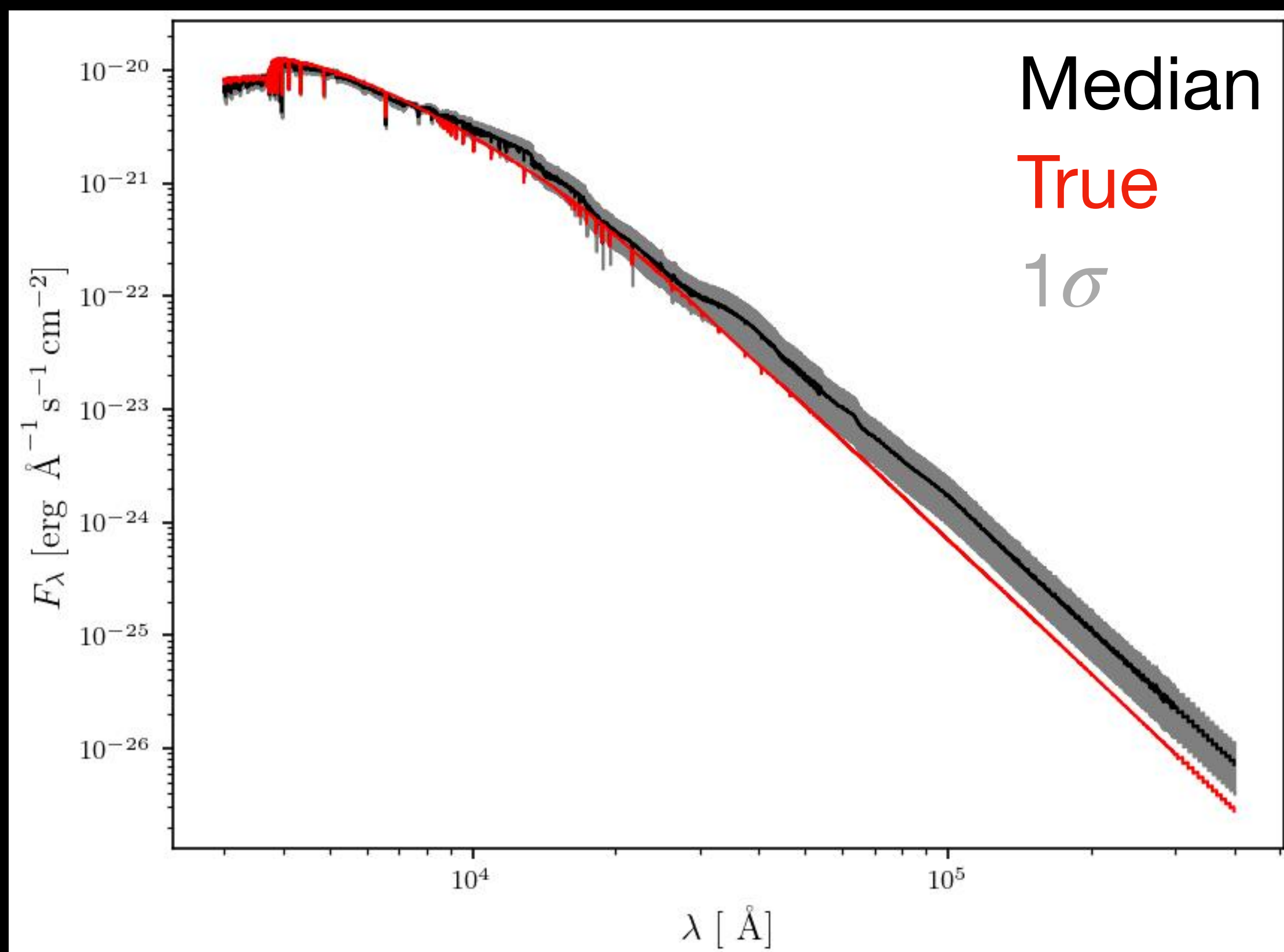




# Gaia + 2MASS + high SNR parallax



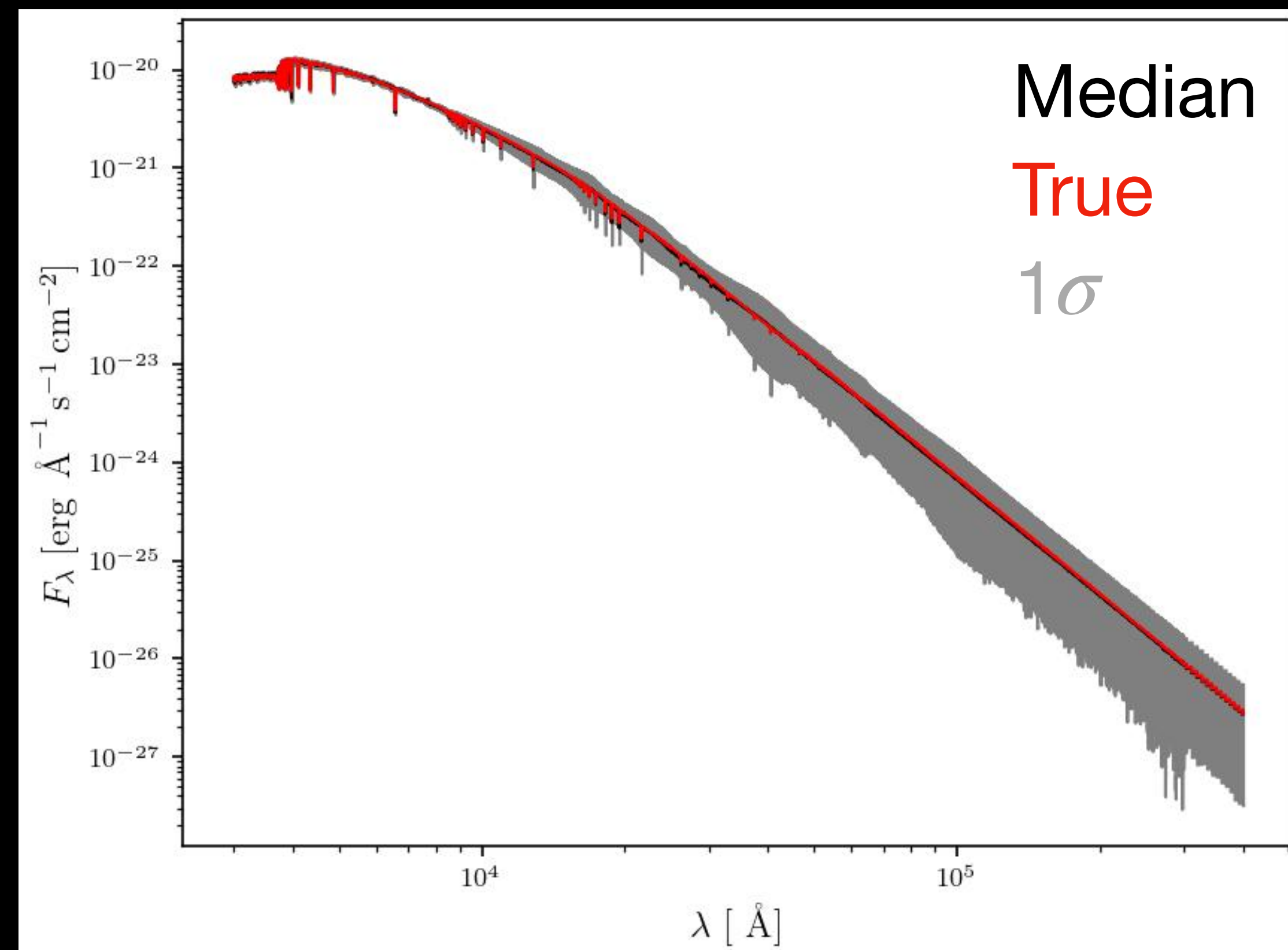
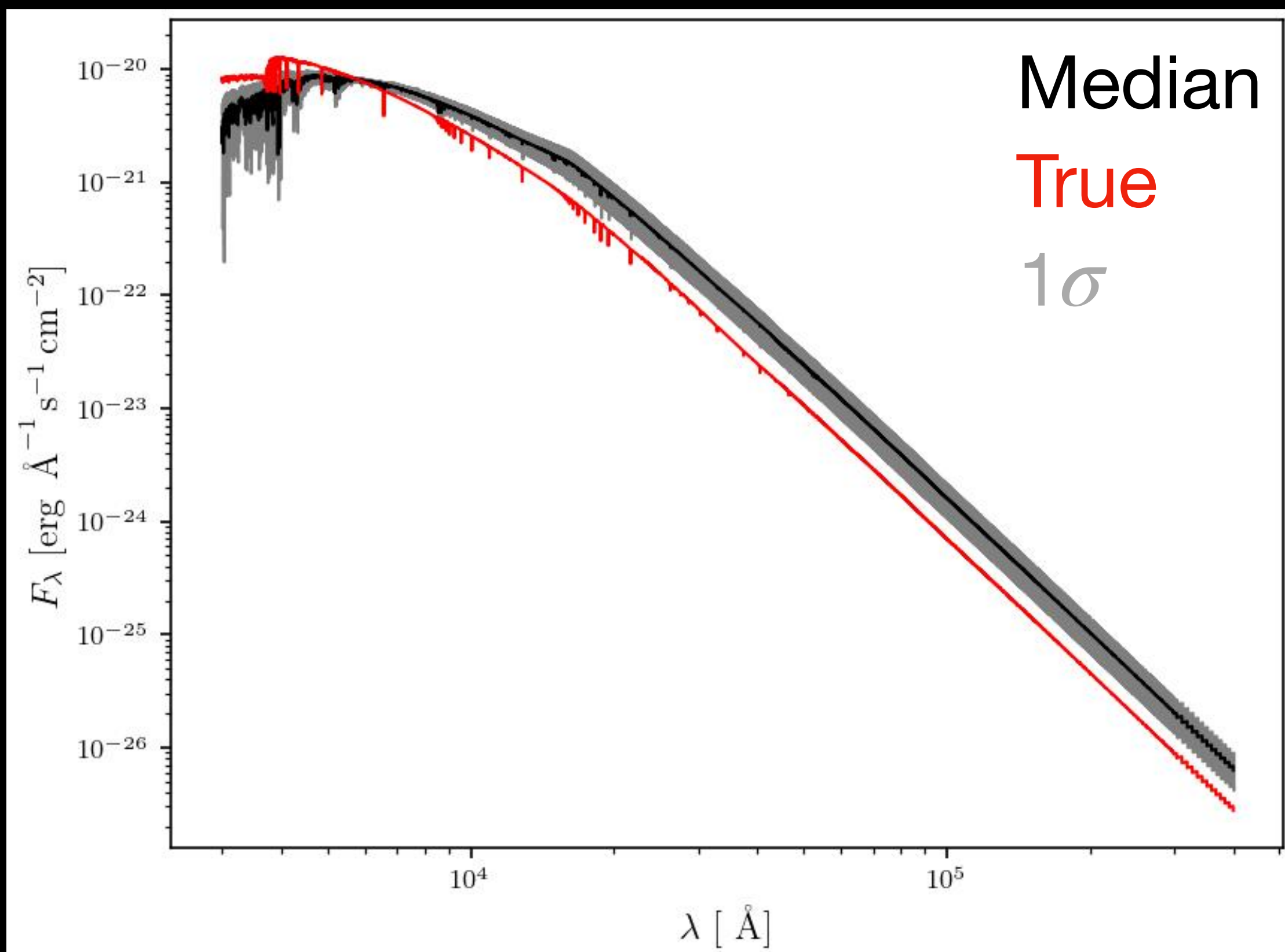
# Gaia + 2MASS + high SNR parallax + LAMOST





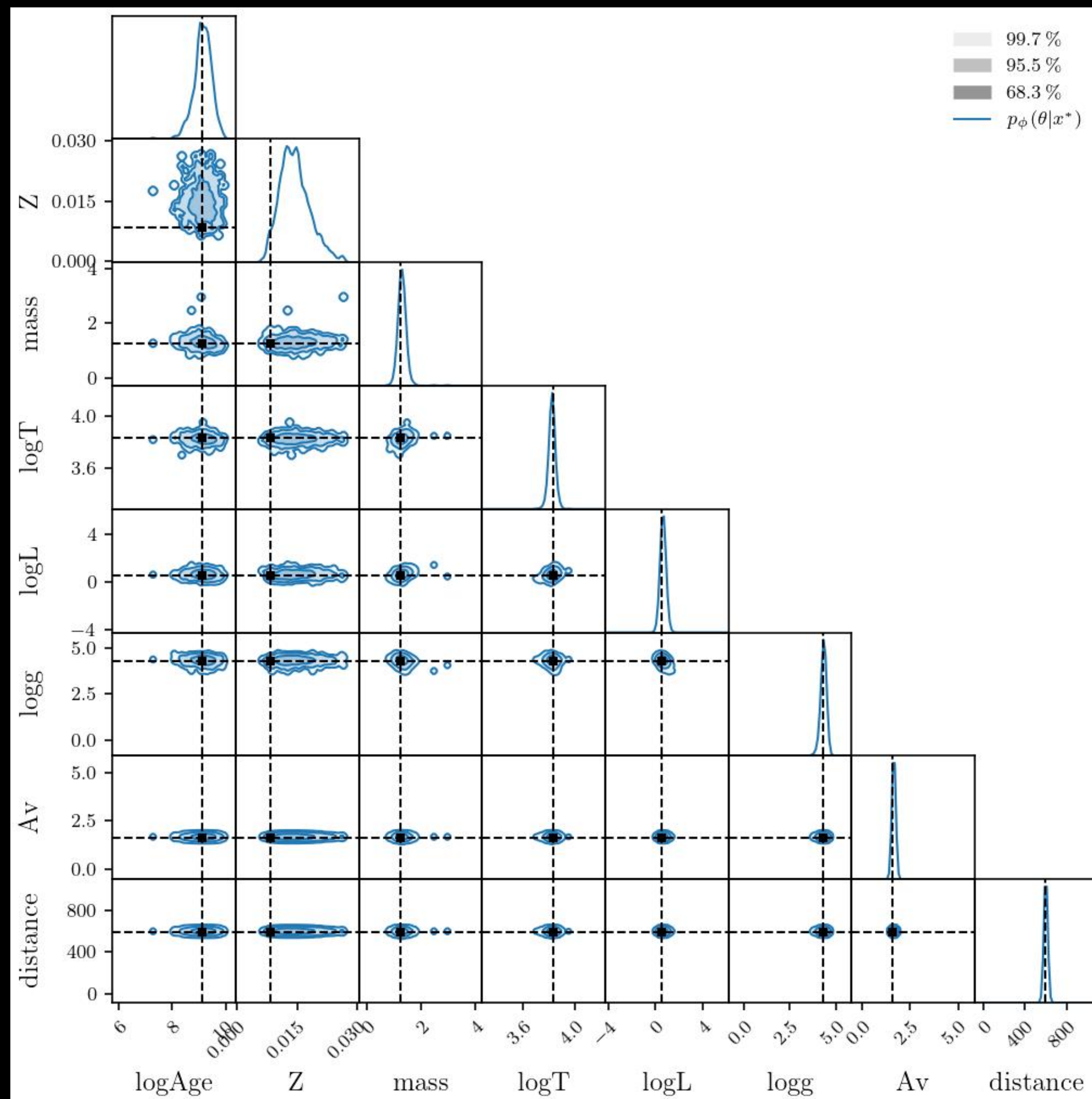
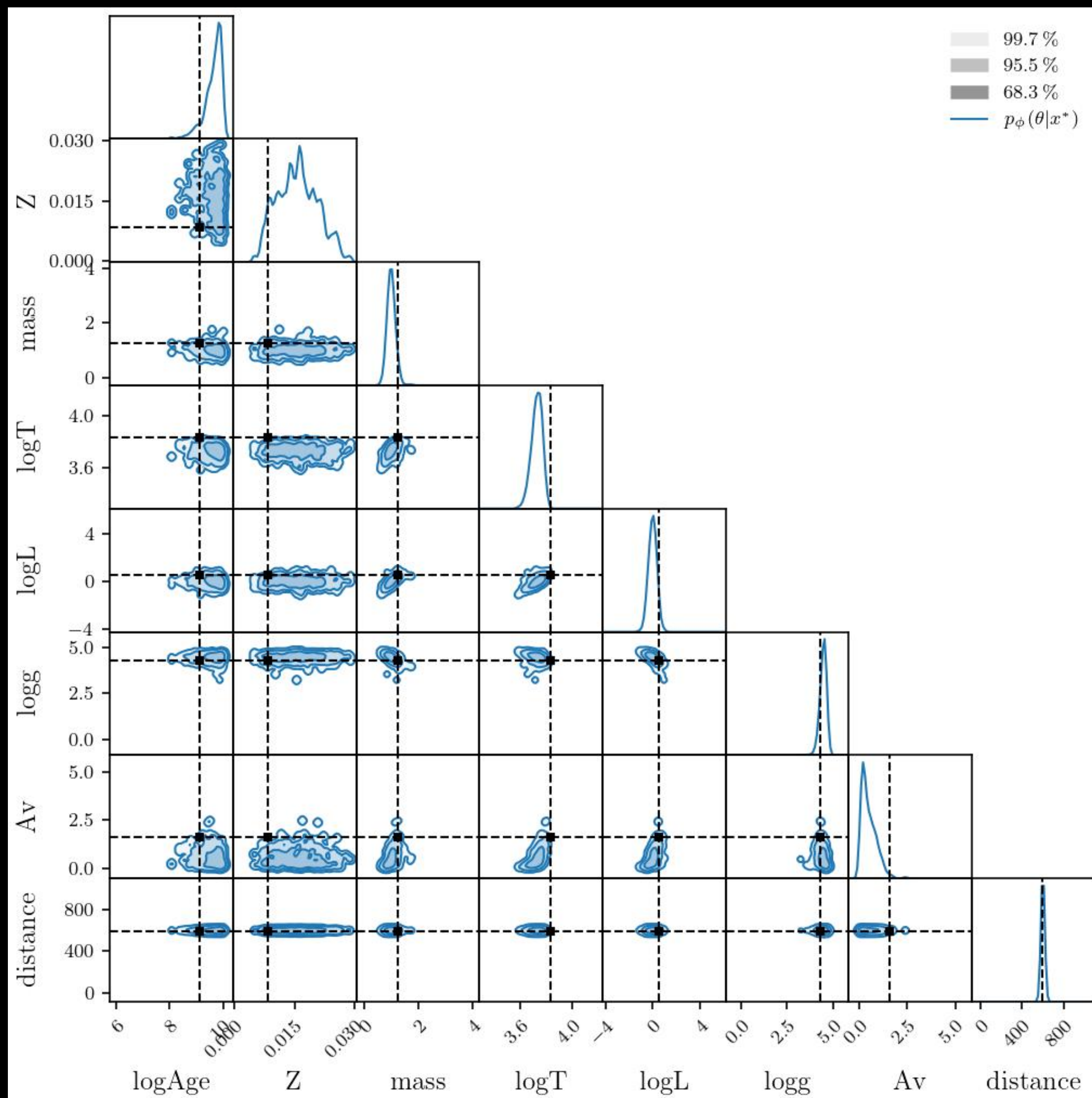
**Adding the guidance term**

# With guidance





# With guidance

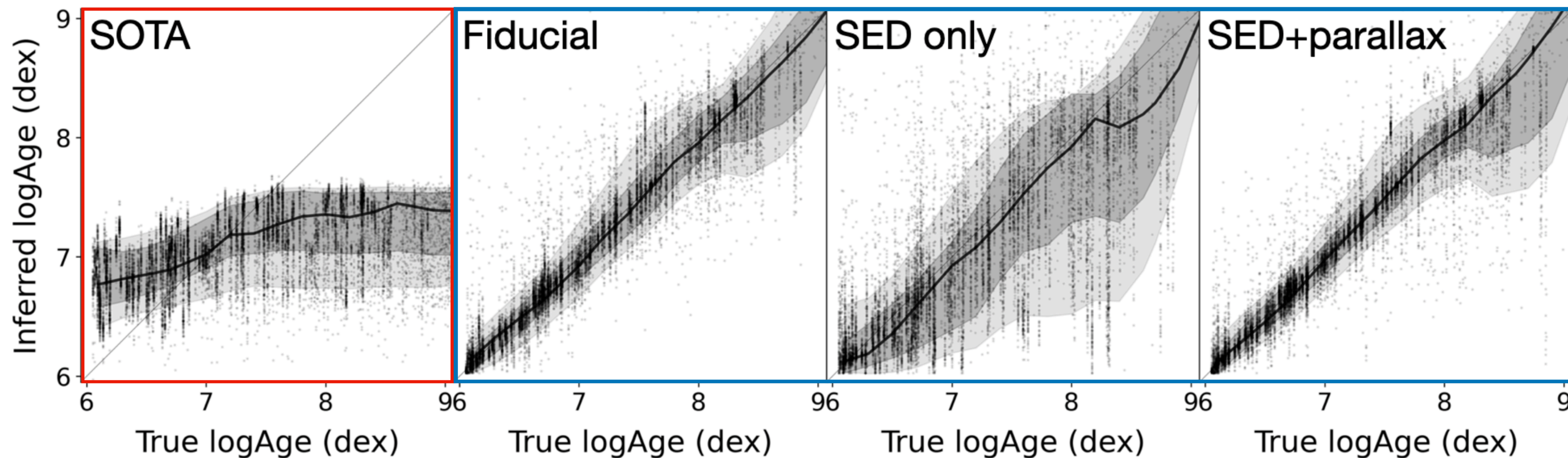




# Pilot study (Gaia + 2MASS + $A_v$ )

Existing

Proposed





# Summary

- SBI model using **flow matching + transformer model** to learn arbitrary **conditionals**
- Guidance via measurement operator constraints
- Obtain promising results on simulations

**Thank you**