# **Running Late: Testing Delayed Supermassive Black Hole Growth Models Against the Quasar Luminosity Function**

Megan Taylor Tillman,<sup>1,2,3</sup>\* Sarah Wellons,<sup>2</sup> Claude-André Faucher-Giguère<sup>2</sup>, Luke Zoltan Kelley,<sup>2</sup> and Daniel Anglés-Alcázar<sup>4,5</sup>

<sup>1</sup>Department of Physics and Astronomy, Texas A&M University, College Station, TX 77843, USA

<sup>2</sup>Center for Interdisciplinary Exploration and Research in Astrophysics (CIERA) & Department of Physics and Astronomy, Northwestern University, Evanston, IL 60201, USA

<sup>3</sup>Department of Physics & Astronomy, Rutgers, The State University of New Jersey, 136 Frelinghuysen Rd, Piscataway, NJ 08854, USA

<sup>4</sup>Department of Physics, University of Connecticut, 196 Auditorium Road, U-3046, Storrs, CT 06269-3046, USA

<sup>5</sup>Center for Computational Astrophysics, Flatiron Institute, New York, NY 10010, USA

Accepted XXX. Received YYY; in original form ZZZ

## ABSTRACT

Observations of massive galaxies at low redshift have revealed approximately linear scaling relations between the mass of a supermassive black hole (SMBH) and properties of its host galaxy. How these scaling relations evolve with redshift and whether they extend to lower-mass galaxies however remain open questions. Recent galaxy formation simulations predict a delayed, or "two-phase," growth of SMBHs: slow, highly intermittent BH growth due to repeated gas ejection by stellar feedback in low-mass galaxies, followed by more sustained gas accretion that eventually brings BHs onto the local scaling relations. The predicted two-phase growth implies a steep increase, or "kink," in BH-galaxy scaling relations at a stellar mass  $M_* \sim 5 \times 10^{10} M_{\odot}$ . We develop a parametric, semi-analytic model to compare different SMBH growth models against observations of the quasar luminosity function (QLF) at  $z \sim 0.5 - 4$ . We compare models in which the relation between SMBH mass and galaxy mass is purely linear versus two-phase models. The models are anchored to the observed galaxy stellar mass function, and the BH mass functions at different redshifts are consistently connected by the accretion rates contributing to the QLF. The best fits suggest that two-phase evolution is significantly preferred by the QLF data over a purely linear scaling relation. Moreover, when the model parameters are left free, the two-phase model fits imply a transition mass consistent with that predicted by simulations. Our analysis motivates further observational tests, including measurements of BH masses and AGN activity at the low-mass end, which could more directly test two-phase SMBH growth.

**Key words:** cosmology: theory – galaxies: evolution – galaxies: active – galaxies: luminosity function – quasars: supermassive black holes

# **1 INTRODUCTION**

The co-evolution of supermassive black holes (SMBHs) and their host galaxies has been an active area of research for over two decades, driven in large part by observations (for a comprehensive review, see Kormendy & Ho 2013). Indeed, studies made possible by the Hubble Space Telescope revealed that all (or nearly all) low-redshift, massive galaxies host a nuclear SMBH. Moreover, detailed observations have shown that masses of SMBHs, M<sub>BH</sub>, correlate surprisingly tightly with properties of their host galaxies, such as the stellar bulge mass  $M_{*bulge}$  (e.g., Magorrian et al. 1998; Häring & Rix 2004; Marconi & Hunt 2003), or the velocity dispersion  $\sigma$  of the stellar bulge (e.g., Ferrarese & Merritt 2000; Gebhardt et al. 2000; Tremaine et al. 2002).

The local  $M_{BH}$ - $M_{*bulge}$  scaling relation has generally been found to be nearly linear. However, this relation is largely constrained by measurements of relatively high mass galaxies in the local universe.

The empirical constraints on the  $M_{BH}$ - $M_{*bulge}$  relation remain comparatively poor at lower masses and high redshift, despite a number of important observational efforts. High-redshift measurements rely on indirect methods to estimate BH masses, such as emission line widths in active galactic nuclei (AGN; e.g., Treu et al. 2007; Merloni et al. 2010; Shen et al. 2015), while more direct measurements in low-redshift dwarf galaxies are limited by sample sizes (e.g., Läsker et al. 2016; Nguyen et al. 2019; Schutte et al. 2019). Recent studies have also noted that the scaling relation appears to depend on the sample selection, for example early- vs. late-type or blue vs. red galaxies (Graham & Scott 2013; Savorgnan 2016; Sahu et al. 2019). These observational results imply that extrapolating the usual, linear scaling relation down to lower-mass galaxies may not be correct in general.

Constraining the form of BH-galaxy scaling relations across the entire spectrum of galaxy masses and types, including in the early universe, is important not only for its own sake but also because it would allow a better understanding of how SMBHs grow with and affect their host galaxies via AGN feedback. AGN feedback is a

<sup>\*</sup> E-mail: mtt74@physics.rutgers.edu

key ingredient in modern galaxy formation theories, but how exactly SMBHs couple to their host galaxies and halos remains a major unknown (e.g., Somerville & Davé 2015; Naab & Ostriker 2017). In current models, AGN feedback is usually assumed to be critical for quenching star formation in massive galaxies, which is needed to explain the observed sequence of "red and dead" galaxies (e.g., Faber et al. 2007; Hopkins et al. 2008; Chen et al. 2020). However, there is increasing observational evidence of AGN-driven outflows in dwarf galaxies (e.g., Manzano-King et al. 2019; Liu et al. 2020). This suggests that AGN feedback could be important for lower-mass galaxies as well. AGN feedback can in principle affect BH-galaxy scaling relations either through its effect on star formation in the host galaxy or by regulating the growth of the nuclear BH.

This paper focuses on testing a prediction from a number of recent galaxy formation simulations concerning the growth of SMBHs. In the last several years, cosmological simulations of galaxy formation have advanced greatly both in resolution (especially in "zoom-in" simulations) as well as in how small-scale processes such as star formation, stellar feedback, and black hole physics are modeled (for reviews of recent progress, see Faucher-Giguère 2018; Vogelsberger et al. 2020). Although simulations using different codes differ in many important details, some predictions appear generic to relatively wide variations in simulation methodologies. A well-known example of this is the role of stellar feedback in shaping the low-mass end of the galaxy stellar mass function (e.g., Somerville & Davé 2015). Here, we are motivated by another prediction which appears generic to many different simulations, namely the delayed growth of SMBHs, with respect to stellar mass growth, due to gas ejection by stellar feedback.

It has been found in multiple simulations by different groups that SMBHs tend to grow in two different phases (e.g., Dubois et al. 2015; Bonoli et al. 2016; Bower et al. 2017; Habouzit et al. 2017; Prieto et al. 2017; Anglés-Alcázar et al. 2017b; McAlpine et al. 2018; Çatmabacak et al. 2020). In the early universe or in low-mass galaxies, feedback by stars (in particular supernovae) regularly ejects gas from galaxy centers. This results in extended periods of time during which there is little to no nearby gas for nuclear BHs to accrete. In this early phase, the BH mass lags behind while the host galaxy grows its stellar mass.<sup>1</sup> It is noteworthy that this result is robust to details of the BH accretion prescription used in the simulation, which are highly uncertain (see e.g., Hopkins et al. 2016; Anglés-Alcázar et al. 2021), as long as the accretion is tied to the gas reservoir in the immediate vicinity of the BH and the effects of stellar feedback are resolved (Anglés-Alcázar et al. 2017b). The simulations find that, eventually, the gas reservoir stabilizes in galactic nuclei. From that point on, nuclear BHs in star forming galaxies accrete at a much higher time-averaged rate and grow to masses comparable to those expected from local scaling relations.

Figure 1 shows results from galaxy formation simulations from the FIRE ("Feedback In Realistic Environments") project (Hopkins et al. 2014, 2018)<sup>2</sup> illustrating the "two-phase" SMBH growth. At the high-mass end, corresponding to later times for the galaxies tracked, the BHs end with masses roughly consistent with locally observed scaling relations (shown here in terms of  $M_{BH}$  vs. total galaxy stellar mass  $M_*$ ). However, in low-mass galaxies BHs can be under-massive relative to their host galaxies by more than an order of magnitude. As Figure 1 shows, in FIRE this produces a relation between BH and galaxy masses which has a prominent "kink" at a galaxy stellar mass  $M_{*}\sim5\times10^{10}~M_{\odot}.$  There is not yet agreement on the primary cause of the change in SMBH fueling regimes, but different possibilities have been discussed including an increase in the escape velocity in the galactic nucleus (Dubois et al. 2015; Anglés-Alcázar et al. 2017b; Lapiner et al. 2021), a change in the buoyancy of galactic winds due to formation of a hot gaseous halo (Bower et al. 2017), and a change in the stability of the gaseous galactic disk, possibly owing to confinement by a hot inner circumgalactic medium (CGM; Stern et al. 2021; Gurvich et al., in prep.; Byrne et al., in prep.). In other simulations, the transition between BH fueling regimes also does not necessarily occur at a fixed stellar mass. For example, in EAGLE the transition is better approximated by a threshold in virial temperature of the halo (McAlpine et al. 2018).

Since the relations between BH and galaxy masses are not well constrained at low masses and at high redshifts, it is not immediately clear whether a kink at  $M_* \sim 5 \times 10^{10} M_{\odot}$  is consistent with BH mass measurements. In this paper, our goal is to test two-phase growth by focusing on another set of observations: the AGN luminosity function. Since AGN are powered by accretion onto SMBHs, the AGN luminosity function is sensitive to the growth history of BHs (e.g., Soltan 1982; Small & Blandford 1992; Yu & Tremaine 2002), including the form of scaling relations as a function of redshift. In practice, this is complicated by the fact that individual observations (e.g., in the optical or X-ray) typically probe only a fraction of the accretion power and by the fact that a large fraction of this accretion power can be missed entirely due to obscuration (e.g., Hickox & Alexander 2018). To circumvent these difficulties, we employ previous studies which have modeled these effects to infer the bolometric luminosity function. We use the results from Shen et al. (2020) which updates the classic analysis of Hopkins et al. (2007) that combined a large set of AGN luminosity function measurements over the redshift interval z = 0 - 6. These authors then obtained a bolometric luminosity function which self-consistently reproduced the observations in different bands, taking into account the luminosity dependence of intrinsic AGN spectra as well as the luminosity dependence of their obscuring columns.

We want to account for the fact that the details of the BH growth histories can differ from simulation to simulation, depending on the specific physics prescriptions used. Thus, instead of comparing exact predictions from a specific set of simulations, we construct a general, parameterized semi-analytic model intended to capture a range of possible variations around the type of kinked scaling relation shown in Figure 1. Although our model includes AGN of a wide range of luminosities, i.e. not only the more luminous AGN commonly known as quasars, we will follow the common practice of referring to the AGN luminosity function also as the quasar luminosity function, or QLF for short. The model scaling relation is convolved with the redshiftdependent galaxy stellar mass function, as well as a distribution of accretion rates, to derive AGN luminosity function predictions. By comparing model AGN luminosity functions produced in this way with the observations, we can test whether a kinked scaling relation consistent with what is predicted by simulations is allowed or even preferred by the luminosity function data. Overall, we find that a two-phase BH growth model can successfully reproduce the QLF in the redshift range  $z \sim 0.5 - 4$ , where it is best constrained. Although this should not be interpreted as proof due to the modeling assumptions and the limited data compared to, we furthermore find evidence that a two-phase growth model is significantly favored over a simpler

<sup>&</sup>lt;sup>1</sup> It is not guaranteed that early BHs will be located at galaxy centers, e.g. if the timescale for "sinking" to the center is too long (e.g., Ma et al. 2021). The delayed BH growth included in our models could also, at least in part, be caused by dynamic effects such as this.

<sup>&</sup>lt;sup>2</sup> See the FIRE project web site: http://fire.northwestern.edu.

model in which the relationship between BH mass and galaxy mass is linear at all masses.

The plan of this paper will be as follows. §2 describes our modeling approach in more detail. Model luminosity functions are fit and compared to observations in §3. We discuss the results in §4, and summarize the main take-aways in §5.

# 2 QLF MODELING METHODOLOGY

When describing the two-phase (2P) model, we will refer to the early/low-mass and late/high-mass growth phases as the "low-mass" and "high-mass" regimes. In addition to parameter variations of the 2P model, we explore a simpler model which assumes a purely linear relation between BH mass and galaxy stellar mass (the L model). By comparing best-fit 2P and L models, we can assess whether the QLF data prefer one over the other.

## 2.1 Galaxy Stellar Masses and Star Formation Rates

All of our models are anchored to observations of the galaxy stellar mass function and designed such that the implied growth of the BH population is consistent with the observed growth of the host galaxy population. We utilize Universe Machine (UM, Behroozi et al. 2019) to model galaxy properties as a function of redshift. In particular, we use UM to model the redshift-dependent galaxy stellar mass function (SMF) and the mean specific star formation rates (sSFRs) of galaxies. In UM, these properties are self-consistently constrained based on a wide range of observations, including different luminosity and correlation functions. UM provides the SMF and sSFR data for a large number of stellar mass and redshift bins. In our modeling code, we interpolate smoothly between the values provided.

## 2.2 The Two-Phase and Linear Models

In this section, we describe in more detail how we implement the two-phase model. The linear model is simply a special case of the 2P model in which there is no transition between distinct regimes.

# Black Hole Masses

In the 2P model, BHs are assumed to follow, on average, a relation between BH mass and total galaxy stellar mass that has two distinct phases, corresponding to the low-mass and high-mass regimes:

$$M_{BH} = \begin{cases} M_{BHlo} & M_* < M_{*crit} \\ \\ M_{BHhi} & M_* > M_{*crit}, \end{cases}$$
(1)

where the two regimes are separated by a critical stellar mass  $M_{*crit}$  (the 'transition mass'). The BH mass scalings in the two regimes are parameterized by three dimensionless parameters: the power-law slope in the low-mass regime ( $\alpha$ ), the asymptotic  $M_{BH}/M_*$  ratio at high mass ( $\beta$ ), and the factor by which the BH is undermassive at  $M_{*crit}$  relative to the linear relation ( $\gamma$ ).

In the low-mass regime, when BHs are undermassive, we use the following parameterization to describe the scaling relation:

$$M_{BHlo} = M_{BHcrit} \left(\frac{M_*}{M_{*crit}}\right)^{\alpha},$$
(2)

where MBHcrit sets the normalization in the low-mass regime. For the

Parameter	Description
M <sub>*crit</sub> (transition mass)	The stellar mass at which the low-mass regime ends and the high-mass regime begins.
$\alpha$ (low-mass slope)	The power-law slope of the $M_{\ast}\mbox{-}M_{BH}$ relation in the low-mass regime.
β (high-mass normalization)	The asymptotic $M_{BH}/M_*$ ratio in the high-mass regime.
γ (break factor)	The factor by which the BH is undermassive at $M_{*crit}$ relative to the high-mass power law.
$\sigma_{\ln X \log}$ (low-mass $\sigma_{\ln X}$ )	The log-normal standard deviation of the normalized BH accretion rate distribution in the low-mass regime.
$\sigma_{\ln Xhi}$ (high-mass $\sigma_{\ln X}$ )	Same as above but for the high-mass regime.

Table 1. Overview of the two-phase (2P) model parameters.

high-mass regime, we set the requirement that the relation approaches linearity,

$$M_{BHhi}(M_* >> M_{*crit}) \to \beta M_*, \tag{3}$$

but we do not enforce a strictly linear relation between  $M_{\rm BH}$  and  $M_*$ . The linear slope at high masses approximates the scaling relations often found in observations (e.g., Marconi & Hunt 2003; Häring & Rix 2004; Kormendy & Ho 2013). Rather, we use continuity considerations described in more detail below to connect the low- and high-mass regimes. We then introduce the dimensionless parameter  $\gamma$  to set the normalization of the low-mass scaling relation,

$$\gamma = \frac{\beta M_{*crit}}{M_{BHcrit}}.$$
(4)

To obtain a functional form for  $M_{BHhi}$  that connects to the lowmass regime, we make the ansatz that as soon as galaxies enter the high-mass regime, the BH mass and stellar mass of the galaxy start growing in proportion to each other:

$$\dot{M}_{\rm BH} = \beta \dot{M}_*,\tag{5}$$

implying that in the high-mass regime,

$$M_{BHhi} = M_{BHcrit} + \beta \left( M_* - M_{*crit} \right).$$
(6)

Figure 1 shows an example of the  $M_{BH}$ - $M_*$  relation from the 2P model, overplotted on FIRE simulation data.

We note that in observations  $M_{BH}$  has often been found to correlate more tightly with the bulge mass rather than the total stellar mass of the galaxy. We use total stellar mass in our model because the simulations predict the 2P growth behavior in the  $M_{BH}$ - $M_*$  relation (see Fig. 1 and Byrne et al., in prep.) and not only when  $M_{BH}$  is plotted as a function of a proxy for bulge mass (Anglés-Alcázar et al. 2017b). We did however explore model variations in which we use bulge mass and found that our main conclusions regarding the twophase model are not changed. This makes sense because, at least at the low redshift where detailed bulge/disk decompositions are available, the bulge mass on average dominates the total stellar mass above ~  $M_{*crit}$  (e.g., Benson et al. 2007). Thus, the effects of distinguishing between bulge mass and total stellar mass are partially degenerate with the parameters that describe the break in the  $M_* - M_{BH}$  relation.



**Figure 1. Left:** The relation between BH mass and galaxy stellar mass. The gray curves show the trajectories of the nuclear BHs in four different A-series simulations of massive galaxies from the FIRE project, analyzed in more depth in Anglés-Alcázar et al. (2017b). The black curves show analytic models for the  $M_*$ - $M_{BH}$  relation: a simple linear model (L; dashed) and a two-phase model (2P; solid) in which BH growth is suppressed at low masses. The 2P model parameters are chosen to approximate the average trajectories of simulated BHs. The red background region corresponds to the "low-mass" regime during which BH growth is slow, whereas the blue background region corresponds to the "high-mass" regime during which BH growth is more efficient. The transition mass between these regimes is indicated by the vertical black dotted line. **Right:** Similar but for the BH accretion rate vs. stellar mass. For clarity, we show the BH accretion rate for a single example simulation and only the 2P model. On each panel, the red labels indicate how each parameter of the 2P model affects the relations shown. These parameters are defined in more detail in §2.2 and in Table 1.

#### Black Hole Accretion Rates

Next we describe how we model the distribution of BH accretion rates. At any given redshift, the distribution of BH masses is obtained by convolving the galaxy stellar mass function with the  $M_{BH} - M_*$  model from the previous section. Since the galaxy stellar mass function evolves with redshift, this implies a specific redshift evolution for the BH mass function. To constrain the distribution of accretion rates, we enforce the requirement that the mean BH accretion rate for any given stellar mass and redshift is consistent with the time evolution of the BH mass function.

The BH accretion rate can be related to the stellar mass growth rate via the following identity:

$$\frac{d\mathbf{M}_{\rm BH}}{dt} = \frac{\mathbf{M}_{\rm BH}}{\mathbf{M}_*} \frac{d\ln\mathbf{M}_{\rm BH}}{d\ln\mathbf{M}_*} \frac{d\mathbf{M}_*}{dt}.$$
(7)

Averaging both sides of this equation (treating  $M_*$  as the independent variable and the slope of the scaling relation as a constant at fixed  $M_*$ ) and combining with equation (1) we obtain for the  $\dot{M}_{BH} - \dot{M}_*$  relation:

$$\left\langle \dot{\mathbf{M}}_{\mathrm{BH}} \right\rangle = \begin{cases} \alpha \mathbf{M}_{\mathrm{BHo}} \frac{\left\langle \mathbf{M}_{*} \right\rangle}{\mathbf{M}_{*}} & \mathbf{M}_{*} < \mathbf{M}_{*\mathrm{crit}} \\ \\ \beta \left\langle \dot{\mathbf{M}}_{*} \right\rangle & \mathbf{M}_{*} > \mathbf{M}_{*\mathrm{crit}} \end{cases}. \tag{8}$$

This form is convenient because the terms involving the galaxy stellar mass and its rate of growth can be self-consistently modeled using results from UM. In doing so, we assume that the factor  $\langle \dot{M}_* \rangle / M_*$  above equals the mean sSFR as a function of stellar mass and redshift which we calculate from UM.<sup>3</sup> We note that, using the relations in

the previous sections, the other terms can be fully parameterized by  $M_{*crit}$  and the dimensionless parameters  $\alpha$ ,  $\beta$ , and  $\gamma$ .

To account for the strong variability in BH accretion rates, we define a distribution of accretion rates around the mean. We use a log-normal distribution for simplicity and to capture a wide range of accretion rates. To minimize the number of free parameters, we make the ansatz that the distribution of BH accretion rates can be parameterized by the distribution of dimensionless fluctuations

$$X \equiv \frac{M_{BH}}{\left< \dot{M}_{BH} \right>}.$$
(9)

The X distribution is assumed to be independent of redshift and to depend only on whether the galaxy is in the low-mass regime or the high-mass regime. The log-normal distribution is then defined as:

$$p\left[\ln \mathbf{X}\right] = \frac{1}{\sqrt{2\pi}\sigma_{\ln \mathbf{X}}} \exp\left(\frac{-\left(\ln \mathbf{X} - \mu_{\ln \mathbf{X}}\right)^2}{2\sigma_{\ln \mathbf{X}}^2}\right). \tag{10}$$

Using the requirement that by definition  $\langle X \rangle \equiv 1$ , implying  $\mu_{\ln X} = -0.5\sigma_{\ln X}^2$ , the distribution is fully characterized by the single parameter  $\sigma_{\ln X}$ .

The low-mass and high-mass distributions parameters are labeled  $\sigma_{\ln Xlo}$  and  $\sigma_{\ln Xhi}$ , respectively. Since accretion is more sporadic in the low-mass regime, we expect that  $\sigma_{\ln Xlo} \ge \sigma_{\ln Xhi}$  and impose this requirement when fitting the model to observational data (see

Kravtsov 2011). We neglect this difference because it is degenerate with the normalization of the scaling relation, which is a free parameter. We also neglect the fact that the stellar mass of a galaxy can grow via mergers. This is a fair approximation for our purposes because this primarily affects the most massive galaxies (e.g., Anglés-Alcázar et al. 2017a) and we do not expect this to significantly change the possible signatures of a change in BH accretion properties at intermediate stellar masses.

 $<sup>^3</sup>$  This identification neglects a  $\sim 20\%$  difference between the SFR and the net stellar mass growth rate  $\dot{M}_*$  owing to stellar mass loss (e.g., Leitner &



**Figure 2.** Illustration of different ingredients used in producing model QLFs, at z = 1. **Left column:** The models are anchored to the observed galaxy stellar mass function and the self-consistently derived mean specific star formation rates of galaxies, which we model using Universe Machine. The top panel shows the SMF (units of Mpc<sup>-3</sup> log<sub>10</sub>  $[M_*]^{-1}$ ) and the bottom panel shows the mean sSFR vs. stellar mass. **Middle column:** The top branch shows the M<sub>BH</sub> – M<sub>\*</sub> relation for a two-phase model with a break at M<sub>\*crit</sub> and the corresponding mean specific BH accretion rate vs. BH mass. The bottom branch shows the same quantities but for a purely linear model (no break in M<sub>BH</sub> – M<sub>\*</sub>). **Right column:** The resulting bolometric QLFs (in units of Mpc<sup>-3</sup> log<sub>10</sub>  $[L_{bol}]^{-1}$ ) for the two-phase and linear models. The solid black curves show the best-fit model QLFs at z = 1 when the model parameters are assumed to be independent of redshift (fixed) and the models are simultaneously fit to observations at z = 0.5, 1, 2, 3, and 4 (see §3.2). The solid red and blue curves in the top panel correspond to the low-mass and high-mass contributions to the two-phase QLF, respectively. The orange shaded regions correspond to different stellar mass bins which can be seen in the left column of this figure. The data points are from the Shen et al. (2020) bolometric QLF and correspond to inferences for observations in different bands.

§3.1). Furthermore, to avoid a discontinuity in  $\sigma_{\ln X}$  at  $M_{*crit}$ , the lognormal dispersion is implemented such that it changes continuously from  $\sigma_{\ln X \log}$  at  $M_* \leq M_{*crit}$  to  $\sigma_{\ln X hi}$  at a stellar mass 0.5 dex higher than  $M_{*crit}$ . The interpolation is done linearly in  $\sigma_{\ln X}$  vs. log  $M_*$ space.

To further characterize BH growth in the models, we define the mean specific BH accretion rate sBHAR  $\equiv \langle \dot{M}_{BH} \rangle / M_{BH}$ , which is a function of redshift and BH mass. Note that this quantity is related to the Eddington ratio, but it is expressed in different units and averaged over the BH population, for a given  $M_{BH}$ .

#### The Linear Limit

To assess whether two-phase SMBH growth is preferred by the QLF data, we also explore a model in which the relation between BH mass and stellar mass is purely linear and with no transition in the normalized accretion rate distribution, i.e. with constant  $\sigma_{\ln X}$ . This linear model is fully specified by the  $\beta = M_{BH}/M_*$  ratio and a single log-normal accretion rate dispersion  $\sigma_{\ln X}$ . We note that, as for the 2P model, our L model is self-consistently anchored to the redshift-dependent SMF and sSFR data from UM. This is significant because some previous QLF models also assumed a log-normal accretion rate distribution at any given mass and redshift, but did not

include the requirement that the distribution of accretion rates must self-consistently connect the implied BH mass functions at different redshifts (e.g., Conroy & White 2013; though see Veale et al. 2014).

## 2.3 Producing the QLF

The steps described so far yield a distribution of BH accretion rates for any given stellar mass and redshift. We convert this straightforwardly to a bolometric luminosity distribution by assuming a constant radiative efficiency:

$$L_{bol} = \epsilon \dot{M}_{BH} c^2.$$
(11)

For simplicity we set  $\epsilon = 0.1$  to represent radiatively efficient accretion disks (e.g., Abramowicz & Fragile 2013).

The observed QLF includes accretion in galaxies of all stellar masses, so to obtain the luminosity function at any redshift, we integrate over the galaxy SMF:

$$\frac{d\mathbf{N}}{d\ln \mathbf{L}_{bol}} = \int p \left[ \ln \mathbf{L}_{bol} | \ln \mathbf{M}_* \right] \frac{d\mathbf{N}}{d\ln \mathbf{M}_*} d\ln \mathbf{M}_*, \tag{12}$$

where

$$p \left[ \ln L_{bol} | \ln M_* \right] = \frac{1}{\sqrt{2\pi\sigma_{\ln\dot{M}_{BH}}^2}} \exp\left(\frac{-\left(\ln \dot{M}_{BH} - \mu_{\ln\dot{M}_{BH}}\right)^2}{2\sigma_{\ln\dot{M}_{BH}}^2}\right), \quad (13)$$



Figure 3. Illustration of how varying each of the six free parameters of the two-phase model affects the results, for the representative redshift z = 1 (see Table 1 for parameter definitions). Top row: Effects on the mean  $M_* - \dot{M}_{BH}$  (black) and  $M_* - M_{BH}$  (purple) relationships. The gray region corresponds to the  $\pm \sigma_{\ln X}$  range that characterizes the log-normal distribution of accretion rates. Bottom row: Effects on the resulting QLF. The black curves correspond to the total QLF, while the red and blue curves show contributions from the low-mass and high-mass regimes, respectively. The values of the free parameters that are varied are, from left to right:  $M_{*crit} = (10^{9.67}, 10^{10.67}, 10^{11.67}) M_{\odot}$ ,  $\alpha = (0.5, 1.0, 1.5)$ ,  $\beta = (10^{-2.53}, 10^{-2.33}, 10^{-2.13})$ ,  $\gamma = (10^{1.34}, 10^{1.84}, 10^{2.34})$ ,  $\sigma_{\ln Xho} = (1.85, 2.35, 2.85)$ , and  $\sigma_{\ln Xhi} = (0.83, 1.33, 1.83)$ . The results for the different parameter values are denoted by different line styles respectively as (dotted, dashed, solid).

 $\mu_{\ln \dot{M}_{\rm BH}} = \mu_{\ln \rm X} + \ln \left< \dot{M}_{\rm BH} \right>, \label{eq:multiplicative}$ 

and

 $\sigma_{\ln \dot{M}_{BH}} = \sigma_{\ln X}.$ 

We do not explicitly model radiatively inefficient accretion, based on the assumption that doing so would only significantly affect the predicted QLF at luminosities too low to affect our conclusions regarding the effects of a change in BH fueling at ~  $M_{*crit}$ . We also do not explicitly model scatter in the  $M_{BH} - M_{*}$  relation. Instead, we assume that the net effect of such scatter can be modeled implicitly as a contribution to the scatter in the accretion rate distribution.

Figure 2 illustrates our process to produce model QLFs, using z = 1 data. The leftmost column displays the SMF and sSFR data. The middle column shows example M<sub>\*</sub>-M<sub>BH</sub> and sBHAR relations for the 2P model (top) and the linear model (bottom). The rightmost column shows the implied QLF for each model. The panel for the 2P model additionally shows the QLF contributions from the low-mass and high-mass regimes. For this figure, the model parameters are best fits to the Shen et al. (2020) QLF data when parameters are assumed to be independent of redshift (fixed) and the models are simultaneously fit to observations at z = 0.5, 1, 2, 3, and 4 (we describe the different fits we explore in the next section).

Figure 3 shows the effects the different free parameters have on the predicted QLF, as well as the M<sub>\*</sub>-M<sub>BH</sub> relation and the M<sub>\*</sub>-M<sub>BH</sub> relations. The effects on the predicted QLF are mostly as may be expected. For example, increasing  $\beta$  increases BH masses and BH accretion rates linearly in the high-mass regime, so this increases the QLF at the high end. Increasing the break factor  $\gamma$  decreases the same quantities in the low-mass regime, so this tends to decrease the QLF at the low end, and can help with imprinting a knee shape (in phenomenological QLF fits, the "knee" is where there is a break in power-law slope between the low-luminosity and high-luminosity regimes). Interestingly, we see that the scatter in the accretion rate distribution (via the  $\sigma_{\ln X}$  parameters) is important in determining the number of high-luminosity quasars (as has also been found in previous studies, e.g. Veale et al. 2014). We note that varying the  $\sigma_{\ln X}$  parameters changes the mean accretion rates because the accretion rate distribution is normal in the logarithm, so the linear mean shifts with  $\sigma_{\ln X}$ .

## **3 RESULTS**

## 3.1 Fitting to Observational Data

We conduct fits of both our 2P and L model QLFs to observational data on the bolometric QLF compiled in Shen et al. (2020). Their study compiles multi-wavelength observational data, including updates from the last decade, using a quasar SED model and bolometric/extinction corrections to update constraints on the observed bolometric QLF from the earlier Hopkins et al. (2007) study. Shen et al. (2020) provides constraints on the QLF from z = 0 to z = 7, but in our study we focus our fits on z = 0.5, 1, 2, 3 and 4 since the observational data is most complete within this range.

To find best-fit parameters, we employ a least-squares method over a finite parameter space. The goodness-of-fit is evaluated at each point in parameter space using the reduced  $\chi^2$  statistic computed in the standard way,  $\chi^2_R = \chi^2/(N-p)$ , where N is the number of data points in the fit and p is the number of free parameters. We evaluate  $\chi^2_R$  over a multidimensional Cartesian grid and find the best fit by minimizing  $\chi^2_R$ .

We carried out several versions of the fit. For each of the 2P and L models, we tried both fitting the model parameters to each redshift independently (the "evolving" fits) and requiring a fixed set of parameter values to simultaneously fit the data at all redshifts (the "fixed" fits). For the evolving fits,  $\chi_R^2$  is evaluated for each redshift, but for the fixed fits, a single  $\chi_R^2$  is evaluated which includes data at all redshifts.

For the 2P model, we conducted one fit with all six model parameters (summarized in Table 1) free, and another with only five free parameters in which  $\alpha$  was fixed to a value of 1. The motivation for the latter fit is that  $\alpha$  parameterizes the low-mass slope of the  $M_{\rm BH} - M_*$  relation and may not be well constrained because



Figure 4. Best-fit QLFs for the two-phase (2P) and linear (L) models, in solid black. 1st column: 2P model with parameters allowed to vary with redshift (evolving). 2nd column: 2P model but with fixed parameters fit simultaneously to all redshifts shown. 3rd column: L model with evolving parameters. 4th column: L model with fixed parameters. Each row displays results for the redshift indicated in the 1st column. The models are fit to bolometric QLF values inferred by Shen et al. (2020) based on observations in different bands, plus extrapolated points at the low-luminosity end (in red) to penalize models that predict a change in slope where observations are not available. The best-fit double power law model from Shen et al. (2020) is shown in gray for reference (their "Global A" fit). The  $\chi_R^2$  quantity shown in panels is the reduced  $\chi^2$  value for the plotted fit. For the fixed fit, there is a single  $\chi_R^2$  corresponding to the global fit. As discussed in the text, the 2P fits shown here assume a fixed  $\alpha = 1$  value because it is found that freeing this parameter does not significantly improve the quality of the fits.

of luminosity limits on the observations. We found that the reduced  $\chi^2$  values were nearly identical for the fits with  $\alpha$  free in the range  $0 \le \alpha \le 2$  (which resulted in a best fit value of  $\alpha = 0.9$ ) versus fixed at  $\alpha = 1$ . For this reason, and to minimize degeneracies associated with multiple parameters, we focus on 2P results for  $\alpha = 1$  fixed for the rest of the paper. Thus, the 2P fits that follow have five free parameters while our linear fits have two free parameters.

Our approach is to let the fits explore wide ranges of parameter values (i.e., to avoid prescribing constraining priors), so that we can determine the parameter values favored by the QLF data alone. The range considered for  $\beta$  was broadly motivated by BH mass measurements at the high end, but nevertheless set wide enough to allow the QLF data to determine a favored value. We confirmed that the likelihood for each free parameter peaks well within the explored range. The parameter ranges used for our final fits are (in the notation for the 2P model):  $8 \le \log_{10} M_{*crit} \le 12, -3.1 \le \log_{10} \beta \le -1.8, 0 \le \log_{10} \gamma \le 4, 0 \le \sigma_{\ln Xho} \le 5$ , and  $0 \le \sigma_{\ln Xhi} \le 5$ . When fitting the 2P model, we implement a physical prior requiring  $\sigma_{\ln Xhi} \le \sigma_{\ln Xho}$ . This is because in hydrodynamical simulations, the accretion rate variability is predicted to be higher in the low-mass regime.

In addition to fitting the observation data points from Shen et al. (2020), we include in the fits additional "extrapolated points" at the low-luminosity end at each redshift. The extrapolated data points are intended to penalize models that imply a low-luminosity QLF shape that diverges strongly from the usual power-law form, which is found where low-luminosity data are available (e.g., down to  $L_{bol} < 10^{42}$  erg s<sup>-1</sup> at z = 0.5). The extrapolated points are most important at high redshift, where the observations do not directly constrain the low-luminosity end (e.g., only down to  $L_{bol} \sim 10^{44}$  erg s<sup>-1</sup> at z = 4). For each redshift, the extrapolated points are introduced starting at the low luminosity limit of the observational data (for the given z), and continue every 0.5 dex down to  $L_{bol} = 10^{41}$  erg s<sup>-1</sup>. We assume these points lie on the Shen et al. (2020) "Global A" best fit and the uncertainty is calculated based on the uncertainty of the Global A fit's parameters.

#### 3.2 QLF Fit Results

Figure 4 shows the best-fit QLFs for both the 2P and L models, and for the evolving vs. fixed fits. The  $\chi_R^2$  values for the fits are indicated on the figure panels (for the evolving fit, there is a  $\chi_R^2$  for each redshift, but for the fixed fits, there is a single value for the global fit to all redshifts simultaneously, indicated in the top row).

We must first acknowledge that none of the fits are ideal from a statistical point of view because the  $\chi^2_R$  values are all significantly above unity. There are a couple likely reasons for this. One is that a close examination of the observational data points suggests that not all the data points are consistent with each other. For example, there are mid-IR and soft X-ray data points with small error bars that appear systematically above other observations and above the Shen et al. (2020) best fit (e.g., around the knee of the QLF at z = 0.5). This suggests that there are some systematic effects not accounted for in the error bars, and that even a perfect model would produce a fit with a  $\chi^2_{R}$  exceeding unity. In fact, our evolving 2P fit appears to describe very well the z = 0.5 observed QLF over the entire luminosity range plotted, even though it has a  $\chi^2_R = 3.038$ . We also note that our models are relatively simple and that there is no a priori guarantee that they can capture all the complexities of the observed data. In particular, our 2P model includes the minimum number of parameters necessary to describe a two-phase scenario in which the parameter values are not fixed. Both the 2P and L models furthermore

assume a simple log-normal distribution of normalized accretion rates (eq. 10), which may not be a fully accurate characterization of AGN activity. Nevertheless, the fits are useful to assess the degree to which a two-phase model is favored over a linear model.

Figure 4 shows that 2P fits are systematically better than the L fits. This is the case for the fixed fits, as well as for the evolving fits for each redshift. For the fixed fits, the difference in reduced  $\chi^2$  is  $\Delta \chi_{R}^{2} = 6.910 - 4.048 = 2.862$  in favor of the 2P model. As mentioned above, the 2P evolving fit appears to describe the z = 0.5 observations very well, and this is the redshift for which the observations cover the largest luminosity range. At higher redshifts  $z \ge 2$ , the 2P fits imply a "bumpy" QLF shape and a low-luminosity end with a steeper slope than the empirical double power-law fit from Shen et al. (2020). We note, however, that these effects are seen where low-luminosity data become sparse or non-existent (other than through the extrapolated points), so this may be largely due to the poor constraints on the fits. We therefore do not consider the "bumpy" low-luminosity QLF shape favored by some 2P fits in higher-redshift bins to be a robust prediction of two-phase growth. On the other hand, the L fits fail in a generic way to describe the observed QLF because they are unable to reproduce a clear knee. This is especially evident at intermediate redshifts z = 2 - 3, where the observations show a distinct knee in the QLF which the best-fit L models entirely fail to capture.

Figure 5 summarizes the best-fit parameters for the 2P model, and shows how the evolving parameters compare to the fixed parameters. This provides a useful check on the 2P model, because physically we expect that if the model is a good description of the AGN population, the model parameters should evolve smoothly with redshift. This is supported by the results in the figure, which show that the best-fit evolving parameters are either stable or only modestly evolving with redshift. This result also explains why the fixed fits are only slightly worse than the evolving fits. Interestingly, the evolving fits suggest some evolution in  $\beta$  (the high-mass M<sub>BH</sub>-M<sub>\*</sub> normalization), such that M<sub>BH</sub>/M<sub>\*</sub> increases by a factor ~ 2 from z ~ 3-4 to z = 0.5. An increase in the high-mass M<sub>BH</sub>/M<sub>\*</sub> ratio by a comparable amount over this redshift interval is also suggested by the recent Trinity empirical model (Zhang et al. 2021).

For reference, the best-fit values for the fixed 2P model are  $\log_{10} M_{*crit} = 10.7^{+0.3}_{-0.5} M_{\odot}$ ,  $\log_{10} \beta = -2.3^{+0.3}_{-0.4}$ ,  $\log_{10} \gamma = 1.8^{+0.6}_{-0.8}$ ,  $\sigma_{\ln Xlo} = 2.4^{+0.4}_{-0.4}$ , and  $\sigma_{\ln Xhi} = 1.3^{+0.9}_{-1}$ . The best-fit values for the fixed L model are  $\log_{10} \beta = -3.00^{+0.08}_{-0.08}$  and  $\sigma_{\ln X} = 2.6^{+0.2}_{-0.2}$ . The errors are  $1\sigma$  and determined by where the 1D profile likelihood for each parameter drops to  $e^{-1/2}$  of its peak value (see the 1D likelihood panels in Fig. 6).

#### **4 DISCUSSION**

In the previous section, we showed that a 2P model fits the QLF data better than an L model. In this section, we assess in more detail the evidence that the QLF data supports two-phase SMBH growth.

One potential concern is that the 2P model might fit the QLF observations better simply because this model has more free parameters, i.e. that the better fit does not necessarily imply that SMBHs grow in two phases. We can assess this by examining the parameter likelihoods for the 2P fits. Figure 6 shows a corner plot for the fixed 2P fit (with  $\alpha = 1$ ). The 2D panels quantify the degeneracies between pairs of parameters, while the 1D panels show the likelihoods for each of the five free parameters. The likelihoods shown are profile likelihoods, which correspond to maximum values of the 5D likelihood function as a function of the one or two parameters considered in each panel (i.e., the likelihood of the best-fit model with the one



**Figure 5.** Comparison of the best-fit parameters for the two-phase model when the parameter values are allowed to evolve with redshift (black dots) vs. are fixed (dashed lines). Overall, the best-fit parameter are relatively stable with redshift even when allowed to vary, though some evolution in  $\beta$  (the high-mass M<sub>BH</sub>-M<sub>\*</sub> normalization) is suggested by the fits. As in Fig. 4, the fits shown here assume a constant  $\alpha = 1$ .

or two parameter values fixed and all others free). First, we note that each of the free parameters is well constrained by the data, in that the profile likelihood has a well-defined peak within the range explored.

We can then ask whether the parameter values preferred by the 2P fit correspond to a "physical" two-phase model, in that the implied M<sub>BH</sub>-M<sub>\*</sub> relation is consistent with the predictions of galaxy formation simulations that motivated our investigation (see Fig. 1). This is a useful question because we constructed our 2P model to be sufficiently flexible that it can represent a wide range of M<sub>BH</sub>-M<sub>\*</sub> relations, including a purely linear limit with no break, or a break of any magnitude at any stellar mass. Using the results from the FIRE simulations shown in Figure 1 as our primary reference point, the main characteristics of a physical two-phase model are a transition mass  $M_{*crit} \sim 5 \times 10^{10} M_{\odot}$  and a break factor  $\gamma \gtrsim 10$ . The transition mass is roughly constant in the FIRE simulations (see also Byrne et al., in prep.). The magnitude of the break factor is not robustly predicted by existing simulations because it is sensitive to the assumed "seed" mass for SMBHs, which determines the normalization of the MBH-M\* relation at the low-mass end (e.g., Anglés-Alcázar et al. 2017b). We consider, however, that a two-phase scenario implies a break of at least one order of magnitude. As Figure 6 shows, it is remarkable that the 2P fit to the QLF observations favors parameter values consistent with these simulation results for a two-phase scenario, with a best-fit transition mass  $M_{*crit} \approx 10^{10.7} M_{\odot} \approx 5 \times 10^{10} M_{\odot}$  and a best-fit break factor  $\gamma \approx 10^{1.8} \approx 60$ . The best-fit high-mass normalization  $\beta = M_{\rm BH}/M_* \sim 0.5\%$  is furthermore similar to the value inferred more directly from BH mass measurements in massive galaxies (e.g., McConnell & Ma 2013; Kormendy & Ho 2013), although we do not compare the exact values in detail because in our QLF analysis the  $\beta$ normalization is degenerate with the assumed radiative efficiency.

There is also some direct observational evidence for a break in the  $M_{BH}$ - $M_*$  relation. In a study that included AGN in dwarf galaxies (with BH masses inferred from broad lines), Reines & Volonteri (2015) found that the low-mass galaxies in their sample had lower average  $M_{BH}/M_*$  than the high-mass galaxies with dynamical BH mass measurements (their Fig. 8). Interestingly, the stellar mass and magnitude of the break in the  $M_{BH}$ - $M_*$  relation from Reines & Volonteri (2015) appear similar to the values favored by our 2P fits to the QLF (quoted above), albeit with large scatter. We note, however, that in a study which measured stellar velocity dispersions  $\sigma_*$ in eight active dwarf galaxies, Baldassare et al. (2020) found that the M<sub>BH</sub>- $\sigma_*$  relation for dwarfs is consistent with a power-law extrapolation from higher masses. These results suggest that when testing models for scaling relations, it is important to consistently compare the models and observations for relations with respect to the same galaxy property (stellar mass, bulge mass, velocity dispersion, etc.). As mentioned in the introduction, other observational studies have also reported evidence for changes in BH-galaxy scaling relations at low masses, for different cuts of the observational samples and different measures of host galaxy properties (Graham & Scott 2013; Savorgnan 2016; Sahu et al. 2019).

We noted in the introduction that multiple different simulations predict two-phase SMBH growth qualitatively similar to what is shown for FIRE simulations in Figure 1. However, the different simulations do not agree in all quantitative details, so it is also interesting to consider how our best fits derived from the QLF data compare with the SMBH growth transitions found in other simulations. We focus here on comparing our results with the EAGLE simulations, for which there have been detailed studies of the SMBH growth transition (Bower et al. 2017; McAlpine et al. 2018). Rather than stellar mass, McAlpine et al. (2018) characterize the transition in terms of the properties of the dark matter halo, finding that the transition occurs in halos of a roughly constant virial temperature  $T_{\rm vir} \approx 10^{5.6}$  K, corresponding to a critical halo mass that decreases with increasing redshift. To compare with McAlpine et al. (2018), we can infer the halo masses and virial temperatures corresponding to the transition stellar masses M\*crit favored by our fits. To do so, we use our "evolving" 2P fits in which the model parameters are allowed to evolve freely with redshift, so that we can infer how the preferred transition halo mass and virial temperature change with redshift. To infer halo mass from stellar mass, we use the median, redshift-dependent stellar mass-halo mass relation from UniverseMachine (Behroozi et al. 2019). The virial temperature is then evaluated using standard relations for virialized halos (Bryan & Norman 1998; Barkana & Loeb 2001).

The results for the implied transition halo mass and virial temperature vs. redshift are shown in Figure 7. Interestingly, the 2P fits to the QLF data imply a transition halo virial temperature  $T_{\rm vir} \sim 10^{6.6}$  K that is roughly constant over the redshift interval z = 0.5 - 4 probed by our analysis, but a factor  $\sim 10 \times$  higher than in the EAGLE simulations analyzed by McAlpine et al. (2018). This is a large systematic offset, corresponding to larger halo masses at the SMBH growth transition. We note, however, that the stellar mass-halo mass relation is relatively flat around the best-fit "fixed"  $M_{*crit}\approx 10^{10.7}M_{\odot},$  so that a small change in stellar mass corresponds to large change in halo mass, according to the median relation. Moreover, in reality, there is scatter in the stellar mass-halo mass relation, as well as uncertainties in empirical determinations of this relation, such as the version in UniverseMachine that we have used here. To more robustly assess whether the QLF data may be consistent with a transition at a lower virial temperature similar to EAGLE, it would be necessary to carefully model scatter and uncertainties in the stellar mass-halo mass relation, which is beyond the scope of the present work.



Figure 6. Corner plot displaying the bi-variate normalized likelihood contours for the two-phase model. This is for the fit with parameters fixed with redshift, and with fixed  $\alpha = 1$  as in Figs. 4 and 5. The 1D panels show profile likelihoods (which indicate the maximum value of the likelihood function for the parameter, with all other parameters free) for the five free parameters. The fit to the observational data favors a transition mass  $M_{*crit} \approx 10^{10.7} M_{\odot} \approx 5 \times 10^{10} M_{\odot}$  and a break factor  $\gamma \approx 10^{1.8} \approx 60$  broadly consistent with physical expectations for a two-phase model from simulations.

# **5 CONCLUSIONS**

## 5.1 Summary of Main Results

We use observations of the quasar luminosity function at  $z \sim 0.5 - 4$  to test a two-phase scenario for SMBH growth motivated by a number of recent galaxy formation simulations. In this picture, SMBHs are undermassive relative to their host galaxies at low masses, or early times, leading to a M<sub>BH</sub>-M<sub>\*</sub> scaling relation with a break and/or increased scatter below a certain mass. In the FIRE simulations, the break occurs at stellar mass  $M_{*crit} \sim 5 \times 10^{10} M_{\odot}$  (e.g., Anglés-Alcázar et al. 2017b; Byrne et al., in prep.), though the exact threshold appears to depend on the simulation (e.g., McAlpine et al. 2018). We developed a flexible semi-analytic framework to predict the QLF implied for different two-phase model parameters, including a purely linear limit in which M<sub>BH</sub>  $\propto$  M<sub>\*</sub> at all masses. Our model incorporates a number of important physical constraints. All models are anchored to the observed galaxy stellar mass function and star

formation rates, as embodied by UniverseMachine (Behroozi et al. 2019). The models are also self-consistent in that the BH growth rates are required to consistently connect the implied BH mass functions at all redshifts. Obscured BH growth is taken into account by comparing to the bolometric quasar luminosity function, which synthesizes observations from the IR to the X-rays (Shen et al. 2020).

The model fits statistically favor a two-phase scenario over a purely linear model with a difference in reduced  $\chi^2$  values of  $\Delta \chi^2_R = 2.862$ . The two-phase model provides an excellent description of the observed QLF over the full luminosity range at z = 0.5, where the observations are most complete and probe the largest luminosity range. The linear model, on the other hand, fails to reproduce a well-defined QLF knee, a problem which is especially severe at intermediate redshifts z = 2 - 3, where a pronounced QLF knee is clearly implied by the data. There is also some evidence that the data favor a twophase scenario for physical reasons, rather than simply because the two-phase fits have more free parameters. Namely, when the model parameters are free within wide ranges, the fits favor values for the



Figure 7. Halo mass and virial temperature corresponding to the redshiftdependent, best-fit transition stellar mass  $M_{*crit}$  for the "evolving" 2P fits (see the black dots in Fig. 5). Halo masses are inferred assuming the median stellar mass-halo mass relation from UniverseMachine (Behroozi et al. 2019).

transition mass and break factor in the  $M_{BH}$ - $M_*$  relation that are very similar to the values predicted by simulations.

#### 5.2 Directions for Future Work

The main caveat to our conclusions is that, while our results indicate that a two-phase model is consistent with QLF observations, we cannot definitively rule out other scenarios. For example, the models explored in this work assume that accretion rates follow a simple log-normal distribution. It is possible that other forms of the accretion rate distribution would allow a linear model to better fit the observed QLF. Other extensions of the model which could modify the results would be to include scatter in the scaling relations, allow for variable radiative efficiencies, or explicitly distinguish between total stellar mass and bulge mass. It would also be interesting to explicitly include the effects of galaxy mergers. For example, McAlpine et al. (2018) find that in EAGLE, rapid BH growth is often initiated by galaxy interactions, especially at low redshift, even though in their simulations a characteristic halo virial temperature remains a requirement for rapid growth.

The model predictions should also be compared to other observations. In §4, we mentioned measurements of BH-galaxy scaling relations extending into the dwarf regime (see also Reines et al. 2013; Mezcua 2017). This kind of direct quantification of scaling relations across a wide range of galaxy masses provides stringent constraints on the models. It would be valuable to more rigorously compare with such observations, accounting for selection effects which can affect scaling relations (e.g., galaxies selected purely based on stellar mass vs. actively accreting). Focusing on active galaxies, it would be valuable to compare the models not only against the QLF (which integrates over the galaxy population at a given redshift) but also against more detailed measurements of accretion rates as a function of stellar mass and redshift (including summary statistics, such as active fractions, as well as full accretion rate distributions). Such observations are already available (e.g., Aird et al. 2018) and could help in breaking degeneracies. As emphasized by Hickox et al. (2014) and Veale et al. (2014), comparing predictions for the distributions of host galaxy properties as a function of AGN luminosity (as opposed to AGN luminosity vs. host properties) can also distinguish between models that otherwise make similar predictions. Finally, if massive BHs in low-mass galaxies depart substantially from standard scaling relations, this would have important implications for expected BH merger rates in dwarf galaxies, and in the early universe in particular. The effects of this may be detectable by future spaced-based gravitational wave detectors (Bailes et al. 2021).

# ACKNOWLEDGEMENTS

We thank the referee, Stuart McAlpine, for a very constructive review. MTT thanks Aaron Geller for hosting the 2018 Northwestern CIERA REU students and for welcoming her back in the summer of 2019. MTT also thanks Kim-Vy Tran and Louis E. Strigari for mentoring her Texas A&M University undergraduate thesis project based on this work. We thank Alex Gurvich for help with code, Alex Richings for data on bulge fractions, Jacob Shen and Phil Hopkins for collaboration on the quasar luminosity function, Philip Arevalo for advice on effective statistical presentations, and Tjitske Starkenburg, Michael Grudic, Jonathan Stern, Lindsey Byrne, and Zachary Hafen for advice during this project. This material is based upon work supported by NSF Grant No. AST-1757792. SW is supported by an NSF Astronomy and Astrophysics Postdoctoral Fellowship under award AST2001905. CAFG was supported by NSF through grants AST-1715216, AST-2108230, and CAREER award AST-1652522; by NASA through grant 17-ATP17-0067; by STScI through grant HST-AR-16124.001-A; and by the Research Corporation for Science Advancement through a Cottrell Scholar Award and a Scialog Award. This work was performed in part at Aspen Center for Physics, which is supported by National Science Foundation grant PHY-1607611. DAA was supported in part by NSF grants AST-2009687 and AST-2108944, and by the Flatiron Institute, which is supported by the Simons Foundation.

# DATA AVAILABILITY

The methods described in the paper should be sufficient to reproduce most results. Additional data can be obtained from the corresponding author on reasonable request.

#### REFERENCES

- Abramowicz M. A., Fragile P. C., 2013, Living Reviews in Relativity, 16, 1
- Aird J., Coil A. L., Georgakakis A., 2018, MNRAS, 474, 1225
- Anglés-Alcázar D., Faucher-Giguère C.-A., Kereš D., Hopkins P. F., Quataert E., Murray N., 2017a, MNRAS, 470, 4698
- Anglés-Alcázar D., Faucher-Giguère C.-A., Quataert E., Hopkins P. F., Feldmann R., Torrey P., Wetzel A., Kereš D., 2017b, MNRAS, 472, L109
- Anglés-Alcázar D., et al., 2021, ApJ, 917, 53
- Bailes M., et al., 2021, Nature Reviews Physics, 3, 344
- Baldassare V. F., Dickey C., Geha M., Reines A. E., 2020, ApJ, 898, L3
- Barkana R., Loeb A., 2001, Phys. Rep., 349, 125
- Behroozi P., Wechsler R. H., Hearin A. P., Conroy C., 2019, MNRAS, 488, 3143
- Benson A. J., Danovic D., Frenk C. S., Sharples R., 2007, Monthly Notices of the Royal Astronomical Society, 379, 841
- Bonoli S., Mayer L., Kazantzidis S., Madau P., Bellovary J., Governato F., 2016, Monthly Notices of the Royal Astronomical Society, 459, 2603
- Bower R. G., Schaye J., Frenk C. S., Theuns T., Schaller M., Crain R. A., McAlpine S., 2017, MNRAS, 465, 32
- Bryan G. L., Norman M. L., 1998, ApJ, 495, 80
- Chen Z., et al., 2020, ApJ, 897, 102
- Conroy C., White M., 2013, ApJ, 762, 70
- Dubois Y., Volonteri M., Silk J., Devriendt J., Slyz A., Teyssier R., 2015, MNRAS, 452, 1502
- Faber S. M., et al., 2007, ApJ, 665, 265
- Faucher-Giguère C.-A., 2018, Nature Astronomy, 2, 368
- Ferrarese L., Merritt D., 2000, ApJ, 539, L9
- Gebhardt K., et al., 2000, ApJ, 539, L13
- Graham A. W., Scott N., 2013, ApJ, 764, 151
- Habouzit M., Volonteri M., Dubois Y., 2017, MNRAS, 468, 3935

# 12 Tillman et al.

- Häring N., Rix H.-W., 2004, ApJ, 604, L89
- Hickox R. C., Alexander D. M., 2018, ARA&A, 56, 625
- Hickox R. C., Mullaney J. R., Alexander D. M., Chen C.-T. J., Civano F. M., Goulding A. D., Hainline K. N., 2014, ApJ, 782, 9
- Hopkins P. F., Richards G. T., Hernquist L., 2007, ApJ, 654, 731
- Hopkins P. F., Cox T. J., Kereš D., Hernquist L., 2008, ApJS, 175, 390
- Hopkins P. F., Kereš D., Oñorbe J., Faucher-Giguère C.-A., Quataert E., Murray N., Bullock J. S., 2014, MNRAS, 445, 581
- Hopkins P. F., Torrey P., Faucher-Giguère C.-A., Quataert E., Murray N., 2016, MNRAS, 458, 816
- Hopkins P. F., et al., 2018, MNRAS, 480, 800
- Kormendy J., Ho L. C., 2013, ARA&A, 51, 511
- Lapiner S., Dekel A., Dubois Y., 2021, MNRAS, 505, 172
- Läsker R., Greene J. E., Seth A., van de Ven G., Braatz J. A., Henkel C., Lo K. Y., 2016, ApJ, 825, 3
- Leitner S. N., Kravtsov A. V., 2011, ApJ, 734, 48
- Liu W., Veilleux S., Canalizo G., Rupke D. S. N., Manzano-King C. M., Bohn T., U V., 2020, ApJ, 905, 166
- Ma L., Hopkins P. F., Ma X., Anglés-Alcázar D., Faucher-Giguère C.-A., Kelley L. Z., 2021, arXiv e-prints, p. arXiv:2101.02727
- Magorrian J., et al., 1998, AJ, 115, 2285
- Manzano-King C. M., Canalizo G., Sales L. V., 2019, ApJ, 884, 54
- Marconi A., Hunt L. K., 2003, ApJ, 589, L21
- McAlpine S., Bower R. G., Rosario D. J., Crain R. A., Schaye J., Theuns T., 2018, MNRAS, 481, 3118
- McConnell N. J., Ma C.-P., 2013, ApJ, 764, 184
- Merloni A., et al., 2010, ApJ, 708, 137
- Mezcua M., 2017, International Journal of Modern Physics D, 26, 1730021
- Naab T., Ostriker J. P., 2017, ARA&A, 55, 59
- Nguyen D. D., et al., 2019, ApJ, 872, 104
- Prieto J., Escala A., Volonteri M., Dubois Y., 2017, ApJ, 836, 216
- Reines A. E., Volonteri M., 2015, ApJ, 813, 82
- Reines A. E., Greene J. E., Geha M., 2013, ApJ, 775, 116
- Sahu N., Graham A. W., Davis B. L., 2019, The Astrophysical Journal, 876, 155
- Savorgnan G. A. D., 2016, ApJ, 821, 88
- Schutte Z., Reines A. E., Greene J. E., 2019, ApJ, 887, 245
- Shen Y., et al., 2015, ApJ, 805, 96
- Shen X., Hopkins P. F., Faucher-Giguère C.-A., Alexander D. M., Richards G. T., Ross N. P., Hickox R. C., 2020, MNRAS, 495, 3252
- Small T. A., Blandford R. D., 1992, MNRAS, 259, 725
- Soltan A., 1982, MNRAS, 200, 115
- Somerville R. S., Davé R., 2015, ARA&A, 53, 51
- Stern J., et al., 2020, arXiv e-prints, p. arXiv:2006.13976
- Tremaine S., et al., 2002, ApJ, 574, 740
- Treu T., Woo J.-H., Malkan M. A., Blandford R. D., 2007, ApJ, 667, 117
- Veale M., White M., Conroy C., 2014, Monthly Notices of the Royal Astronomical Society, 445, 1144
- Vogelsberger M., Marinacci F., Torrey P., Puchwein E., 2020, Nature Reviews Physics, 2, 42
- Yu Q., Tremaine S., 2002, MNRAS, 335, 965
- Zhang H., Behroozi P., Volonteri M., Silk J., Fan X., Hopkins P. F., Yang J., Aird J., 2021, arXiv e-prints, p. arXiv:2105.10474
- Çatmabacak O., Feldmann R., Anglés-Alcázar D., Faucher-Giguère C.-A., Hopkins P. F., Kereš D., 2020, arXiv e-prints, p. arXiv:2007.12185
- This paper has been typeset from a TEX/LATEX file prepared by the author.