

The physical drivers of the atomic hydrogen–halo mass relation

Garima Chauhan^{1,2}, Claudia del P. Lagos^{1,2}, Adam R. H. Stevens^{1,2}, Danail Obreschkow^{1,2},
Chris Power^{1,2} and Martin Meyer^{1,2}

¹International Centre for Radio Astronomy Research, The University of Western Australia, 35 Stirling Highway, Crawley, WA 6009, Australia

²ARC Centre of Excellence for All Sky Astrophysics in 3 Dimensions (ASTRO 3D), Australia

Accepted 2020 July 27. Received 2020 July 26; in original form 2020 June 22

ABSTRACT

We use the state-of-the-art semi-analytic galaxy formation model, SHARK, to investigate the physical processes involved in dictating the shape, scatter, and evolution of the HI–halo mass (HIHM) relation at $0 \leq z \leq 2$. We compare SHARK with HI clustering and spectral stacking of the HIHM relation derived from observations finding excellent agreement with the former and a deficiency of HI in SHARK at $M_{\text{vir}} \approx 10^{12-13} M_{\odot}$ in the latter. In SHARK, we find that the HI mass increases with the halo mass up to a critical mass of $\approx 10^{11.8} M_{\odot}$; between $\approx 10^{11}$ and $10^{13} M_{\odot}$, the scatter in the relation increases by 0.7 dex and the HI mass decreases with the halo mass on average (till $M_{\text{vir}} \sim 10^{12.5} M_{\odot}$, after which it starts increasing); at $M_{\text{vir}} \gtrsim 10^{13} M_{\odot}$, the HI content continues to increase with increasing halo mass, as a result of the increasing HI contribution from satellite galaxies. We find that the critical halo mass of $\approx 10^{12} M_{\odot}$ is set by feedback from active galactic nuclei (AGNs) which affects both the shape and scatter of the HIHM relation, with other physical processes playing a less significant role. We also determine the main secondary parameters responsible for the scatter of the HIHM relation, namely the halo spin parameter at $M_{\text{vir}} < 10^{11.8} M_{\odot}$, and the fractional contribution from substructure to the total halo mass ($M_{\text{h}}^{\text{sat}}/M_{\text{vir}}$) for $M_{\text{vir}} > 10^{13} M_{\odot}$. The scatter at $10^{11.8} M_{\odot} < M_{\text{vir}} < 10^{13} M_{\odot}$ is best described by the black hole-to-stellar mass ratio of the central galaxy, reflecting the relevance of AGN feedback. We present a numerical model to populate dark matter-only simulations with HI at $0 \leq z \leq 2$ based solely on halo parameters that are measurable in such simulations.

Key words: galaxies: evolution – galaxies: formation – galaxies: haloes.

1 INTRODUCTION

Understanding the distribution and evolution of neutral atomic hydrogen (HI) in the Universe provides key insights into cosmology, galaxy formation, and the epoch of cosmic reionization (Blanton & Moustakas 2009; Pritchard & Loeb 2012; Somerville & Davé 2015; Rhee et al. 2018). A long-standing challenge in galaxy formation and evolution is addressing the relationship between stars, gas, and metals in galaxies, haloes, and the large-scale structure. HI is a primary ingredient for star formation (SF) and a key input to understand how various processes govern galaxy formation and evolution. The HI content of dark matter (DM) haloes forms an intermediate state in the baryon cycle that connects the largely ionized gas in the intergalactic medium (IGM), the shock heated gas at the virial radius and the star-forming cold gas in the interstellar medium (ISM) of galaxies (Krumholz & Dekel 2012; Putman, Peek & Joung 2012). Constraints on HI at all relevant scales (IGM, halo, and galaxy scales) are therefore key to reveal the role of gas dynamics, cooling, and regulatory processes such as stellar feedback, gas inflows, and outflows (Prochaska & Wolfe 2009; van de Voort et al. 2011), and the effect of environment in galaxy formation (Fabello et al. 2012; Zhang et al. 2013).

When studying galaxy formation and evolution, the exploration of scaling relations is particularly useful as a way of reducing the inherent complexity of the process and providing a quantitative means of examining physical properties of galaxies. The dependence of the abundance of baryons on the host halo mass is considered one of the most fundamental scaling relations (Wechsler & Tinker 2018). In particular, the stellar–halo mass relation has been studied in detail, and has been shown to have little scatter (≈ 0.2 dex, see Behroozi, Conroy & Wechsler 2010; Moster et al. 2010) and a shape that reflects the mismatch between the halo and stellar mass functions – the latter has a much shallower low-mass end slope and a more abrupt break at the high-mass end than the former (see review Wechsler & Tinker 2018). The scatter around these scaling relations is particularly useful because it helps to pinpoint how a halo’s assembly history affects its baryon content (Mitchell et al. 2016; Matthee et al. 2017; Kulier et al. 2019).

Stellar mass can be inferred observationally for large statistical samples, unlike the gas content of galaxies and haloes. However, given that stellar mass is only a small contribution to the baryon content of the Universe (Fukugita, Hogan & Peebles 1998; Driver et al. 2018), it is imperative to explore how the abundance of different gas phases correlate with halo mass. HI is particularly interesting because it is the intermediate state in the baryon cycle. The HI–halo mass scaling relation (HIHM) is likely to be much more complex than the stellar–halo mass relation because observations show that the correlation between HI mass and stellar mass is characterized by a large scatter (e.g. Catinella et al. 2010; Brown et al. 2015,

* E-mail: garima.chauhan@research.uwa.edu.au (GC); claudia.lagos@icrar.org (CDPL)

2017; Catinella et al. 2018). This is implied by the work of Chauhan et al. (2019), who used galaxy formation simulations to show that the correlation between H I mass and H I velocity width – a tracer of a galaxy’s dynamical mass – is complex, with variations of >2 dex in H I mass at fixed velocity width.

Several empirical studies have inferred limits on the form of the HIHM relation. Eckert et al. (2017) attempted to measure the ‘cold’ baryon mass (stars plus ISM mass) versus halo mass relation, for which they combined 21 cm-derived H I masses with empirical estimates of the gas mass in galaxies based on the correlation between the H I mass and optical colours in galaxies with detected H I. The difficulty with this approach is the unknown systematic effects in the application of the empirical estimation to a wider parameter space than probed by actual H I detections (see Eckert et al. 2015). Other approaches use H I-clustering measurements to infer an HIHM relation (Padmanabhan & Refregier 2017; Obuljen et al. 2019), as well as H I spectral stacking, which has been used to calculate the mean H I content of groups identified in optical redshift surveys (Guo et al. 2020). H I clustering provides an indirect way of measuring the HIHM relation because it relies on abundance matching to match the H I with the respective halo that will be expected to host galaxies of the observed H I mass. In contrast, H I stacking provides a direct measurement of the *mean* H I mass inside haloes of a given mass, typically using an estimate of the halo radius to choose the stacking area. However, it relies on group finders and halo mass estimates based on optical redshift surveys and so care must be taken because of the well-known issue that optically selected and H I-selected galaxies do not fully overlap, such that H I-selected surveys typically miss the most massive, gas-poor galaxies (e.g. de Blok, McGaugh & van der Hulst 1996; Schombert, McGaugh & Eder 2001). The HIHM relation is also expected to differ from the stellar–halo mass relation because, as previous work has shown, the distribution of H I-selected galaxies depend not only on halo mass but also on the halo’s formation history (Gao, Springel & White 2005; Guo et al. 2017) and on halo spin parameter (Maddox et al. 2015; Obreschkow et al. 2016; Lutz et al. 2018).

While these observational inferences provide highly valuable constraints on the *average* HIHM relation, they do not constrain the scatter. The HIHM relation has been investigated extensively using different theoretical models, including semi-analytic models of galaxy formation (Kim et al. 2017; Baugh et al. 2019; Spinelli et al. 2020) and hydrodynamical simulations (Villaescusa-Navarro et al. 2018), which have consistently shown that the HIHM relation is characterized by a large scatter (especially in the region $10^{12} M_{\odot} < M_{\text{vir}} < 10^{13} M_{\odot}$) – much larger than the stellar–halo mass relation, by >0.5 dex. However, the predicted scatter of the HIHM appears to be largely model-dependent and no observational constraints have been obtained yet. For instance, both Baugh et al. (2019) and Spinelli et al. (2020) attribute the scatter in the relation to feedback from active galactic nuclei (AGNs), which suppresses gas cooling in the halo, preventing further gas accretion on to the central galaxy. Spinelli et al. (2020) also find that the HIHM relation depends on the detailed assembly history of haloes, which agrees with inferences based on H I clustering studies in Guo et al. (2017). Villaescusa-Navarro et al. (2018), using the IllustrisTNG hydrodynamical simulations, also report a larger scatter in their HIHM relation at M_{vir} between 10^{12} and $10^{13} M_{\odot}$, compared to what is found for the stellar–halo mass relation in their simulation.

The current paucity of observational constraints on the shape, scatter, and evolution of the HIHM is likely to change in the coming decade, ultimately with the Square Kilometer Array (SKA; see Abdalla & Rawlings 2005), but also with its pathfinders (e.g. MeerKAT, see Holwerda et al. 2011 and the Australian SKA Pathfinder,

ASKAP; see Duffy et al. 2012; Koribalski et al. 2020). With these transformational instruments on the horizon, it is imperative that we use current galaxy formation models and simulations to explore the physics shaping the HIHM relation to offer predictions and aid the interpretation of these upcoming observations. This is the main motivation of this paper.

Another important challenge is the fact that the SKA is expected to probe cosmological volumes much larger than those we currently use to study galaxy formation (Power et al. 2015), even in the case of semi-analytic models of galaxy formation – whose accessible volumes are already 2–3 orders of magnitude larger than what we can reliably do with hydrodynamical simulations. In the case of semi-analytic models, the typically used cosmological volumes are usually limited by the fact that we require enough resolution to accurately model the assembly and growth history of the haloes. The challenge is even greater if we focus on cosmological studies with the SKA, which require thousands of statistical realizations of the universe with trustworthy models describing how to populate haloes with H I mass. This demands a physically motivated way of populating DM-only simulations with H I without the need of running computationally expensive physical galaxy formation models on them. This is an important second motivation for our work.

These motivations require an in-depth exploration of the astrophysical processes that shape the HIHM relation and the development of an analytical model for how to populate dark matter haloes with H I. We aim to understand what physical parameters are responsible for how H I populates haloes, and what drives the shape and scatter of the relation. For this, it is necessary to assess how the baryon physics included in galaxy formation simulations and halo formation history affect the HIHM relation across cosmic time. We explore which (other) halo properties affect the HIHM relation (e.g. spin, substructure mass fraction etc.). We do this by the use of the SHARK semi-analytic model of galaxy formation (Lagos et al. 2018) and leverage its modularity and flexibility to test the effect of different physical models and parameters on the shape of the HIHM relation. We expect our numerical model showing how to populate DM haloes with H I to be beneficial for designing H I-stacking and H I-intensity mapping experiments.

The structure of this paper is as follows. Section 2 summarizes the relevant features of SHARK. Section 3 validates our semi-analytic model against the local Universe H I observations that capture the average HIHM relation. In Section 4, we delve into the properties responsible for the shape and scatter of the HIHM relation, and see how much impact these properties have. In Section 5, we present our physically motivated HIHM relation along with providing information on its evolution with redshift. We draw conclusions in Section 6. The Appendices show how the HIHM relation evolves with redshift and provide tabulated fits to populated haloes with H I mass.

2 MODELLING THE H I CONTENT OF GALAXIES AND HALOES

In this section, we describe the semi-analytical model that is used in the study, and which prescriptions are applied to calculate the H I content of galaxies and haloes. The results of using these models are discussed in Section 4.

2.1 The SHARK semi-analytical model of galaxy formation

We use the semi-analytical model of galaxy formation (SAM), SHARK (Lagos et al. 2018). SAMs use halo merger trees, which are produced from a cosmological DM only N -body only simulation, and follow

the formation and evolution of galaxies by solving a set of equations that describe all the physical processes that (we think) are relevant for the problem (see reviews by Baugh 2006; Somerville & Davé 2015).

SHARK¹ is an open-source, flexible and highly modular SAM that models the key physical processes of galaxy formation and evolution. These include

- (i) the collapse and merging of DM haloes;
- (ii) the accretion of gas on to haloes, which is governed by the DM accretion rate;
- (iii) the shock heating and radiative cooling of gas inside DM haloes, leading to the formation of galactic discs via conservation of specific angular momentum of the cooling gas;
- (iv) the formation of a multiphase interstellar medium and subsequent SF in galaxy discs;
- (v) the suppression of gas cooling due to photoionization;
- (vi) chemical enrichment of stars and gas;
- (vii) stellar feedback from evolving stellar populations;
- (viii) the growth of supermassive black holes (SMBH) via gas accretion and merging with other SMBHs;
- (ix) heating by AGNs;
- (x) galaxy mergers driven by dynamical friction within common DM haloes, which can trigger bursts of SF and the formation and/or growth of spheroids; and
- (xi) the collapse of globally unstable discs leading to bursts of SF and the creation and/or growth of bulges.

SHARK also includes several different prescriptions for gas cooling, AGN feedback, stellar, and photoionization feedback, and SF.

Using these models, SHARK computes the exchange of mass, metals, and angular momentum between the key baryonic reservoirs in haloes and galaxies, which include hot and cold halo gas, the galactic stellar and gas discs and bulges, central black holes, as well as the ejected gas component that tracks the baryons that have been expelled from haloes. In Section 2.3, we describe in detail the modelling of SF, AGN feedback, stellar feedback, reionization, and gas stripping in satellite galaxies, all of which are relevant for the discussions in Sections 4–6.

The models and parameters used in this study are the SHARK defaults, as described in Lagos et al. (2018) and used in Chauhan et al. (2019) to study the HI content of galaxies. These have been calibrated to reproduce the $z = 0, 1,$ and 2 stellar mass functions; the $z = 0$ black hole–bulge mass relation; and the disc and bulge mass–size relations. This model also successfully reproduces a range of observational results that are independent of those used in the calibration process. These include the total neutral, atomic, and molecular hydrogen–stellar mass scaling relations at $z = 0$; the cosmic star formation rate (SFR) density evolution up to $z \approx 4$; the cosmic density evolution of the atomic and molecular hydrogen at $z \leq 2$ or higher in the case of the latter; the mass–metallicity relations for gas and stellar content; the contribution to the stellar mass by bulges; and the SFR–stellar mass relation in the local Universe. Davies et al. (2018) show that SHARK reproduces the scatter around the main sequence of SF in the SFR–stellar mass plane; Chauhan et al. (2019) show that SHARK can reproduce the HI mass and velocity widths of galaxies observed in the ALFALFA survey; and Amaratidis et al. (2019) show that the predicted AGN luminosity functions (LFs) agree well with observations in X-rays and radio wavelengths.

¹<https://github.com/ICRAR/shark/>

Table 1. SURFS simulation parameters of the runs being used in this paper. We refer to L40N512 and L210N1536 as micro-SURFS and medi-SURFS, respectively.

Name	Box size $L_{\text{box}} [\text{cMpc } h^{-1}]$	Number of particles N_p	Particle mass $m_p [M_\odot h^{-1}]$	Softening length $\epsilon [\text{ckpc } h^{-1}]$
L40N512	40	512^3	4.13×10^7	2.6
L210N1536	210	1536^3	2.21×10^8	4.5

In addition, Lagos et al. (2019) has shown that SHARK can reproduce the panchromatic emission of galaxies throughout cosmic time; most notably, SHARK reproduces the number counts from GALEX UV to the JCMT 850 μm band, the redshift distribution of submillimetre galaxies, and the ALMA bands number counts (Lagos et al. 2020). Bravo et al. (2020) show that SHARK also reproduces reasonably well the optical colour distribution of galaxies across a wide range of stellar masses and redshift, as well as the fraction of passive galaxies as a function of stellar mass.

We use the SURFS suite of DM only N -body simulations for our study (Elahi et al. 2018), which consist of N -body simulations of differing volumes, from 40 to $210 h^{-1}$ cMpc on a side, and particle numbers, from ~ 130 million up to ~ 8.5 billion particles. The simulations adopt the Lambda cold dark matter (Λ CDM) Planck cosmology (Planck Collaboration XIII 2016), which assumes total matter, baryon, and dark energy densities of $\Omega_m = 0.3121$, $\Omega_b = 0.0491$, and $\Omega_\Lambda = 0.6751$, and a dimensionless Hubble parameter of $h = 0.6751$.

For this analysis, we use the L40N512 and L210N1536 runs, referred to as micro-SURFS and medi-SURFS, respectively, and whose properties are described in Table 1. By using two different resolution runs of different volumes, we can probe over six orders of magnitude in DM halo mass, thus giving us an optimal dynamic range for exploring the HIHM scaling relation. We show the results of SHARK using micro-SURFS at halo masses below $10^{11.2} M_\odot$, while medi-SURFS is used for higher halo masses. This transition mass is chosen as according to Elahi et al. (2018) at this mass haloes in medi-SURFS comprise ≥ 200 particles, making them reliable for our calculation (because their merger trees will be sufficiently well resolved). Merger trees and halo catalogues were constructed using the phase-space finder VELOCIRAPTOR (Cañas et al. 2019; Elahi et al. 2019a) and the halo merger tree code TREFROG (Poultou et al. 2018; Elahi et al. 2019b).

We define three types of galaxies in our analysis: *centrals*, *satellites*, and *orphans*. SHARK uses the merger trees and subhalo catalogues as a skeleton, that is required to evolve our galaxies, and so we use this information to describe our galaxy types as well. In SHARK, the central subhalo of every halo in the catalogue is defined as the most massive subhalo of every existing halo at $z = 0$, and then subsequently making the main progenitor of those centrals as the centrals of their respective halo. Every subhalo/halo is connected to its progenitor(s) and descendant subhalo/halo, which is connected to the merger tree they belong to. Haloes point to their central and satellite subhaloes, with the subsequent subhaloes pointing to the list of galaxies they may contain. Following the subhalo and merger tree information, we define *centrals* or $\text{type} = 0$ to be the central galaxy of the central subhalo. We only allow the central subhaloes to host the central galaxy, which in turn becomes the central galaxy of the hosthalo. The *satellite* or $\text{type} = 1$ galaxies are the central galaxies of the other existing subhaloes for that hosthalo (satellite subhaloes). The galaxies belonging to a subhalo that merges on to another one and is not the main progenitor become the *orphan* or $\text{type} = 2$ galaxies. A central subhalo in SHARK can have only

one central galaxy and any number of orphan galaxies, whereas the satellite subhalo can only have one `type = 1` galaxy. When a subhalo becomes a satellite subhalo, any orphan galaxies in that subhalo are transferred to the central subhalo.

2.2 Halo properties as calculated in SHARK

SHARK assumes the masses of DM haloes (M_{halo}) to be those calculated by VELOCIRAPTOR. The virial mass is defined as $M_{\text{halo}} \equiv M_{200} = 4\pi R_{200}^3 \Delta\rho_{\text{crit}}/3$, with ρ_{crit} being the critical density of the universe, with M_{200} and R_{200} being the mass and radius of the halo, respectively, when the density within the halo becomes 200 times of the critical density of the universe. It is assumed that the mass profile of the halo follows an NFW profile (Navarro, Frenk & White 1997). The halo concentration is estimated using the Duffy et al. (2008) relation between concentration, the halo’s virial mass and redshift. The spin parameter of the haloes is drawn from a lognormal distribution of mean 0.03 and width 0.5. These parameters correspond to those measured in SURFS with the well-resolved haloes (Elahi et al. 2018).

2.3 Modelling of key physical processes in SHARK

As stated in Section 2.1, SHARK is a modular SAM, and so the user can adopt a range of models for different physical processes. Although we use the default SHARK model for the derivation of the HIHM scaling relation, we also want to understand what drives the shape of the HIHM relation, and so varying the models and parameters adopted in SHARK is necessary. Here, we describe a subsample of the models and physical processes that are relevant for the HIHM relation.

We compare the H I in haloes based on two different ISM gas-phase models, different AGNs and stellar feedback efficiencies, and different ram pressure stripping considerations, as well as altering the photoionization of H I in haloes.

2.3.1 Gas phases in the interstellar medium and star formation

In the default SHARK model, hereafter referred to as SHARK-ref, we use the prescription described in Blitz & Rosolowsky (2006), hereafter referred to as BR06, to compute the amount of atomic and molecular hydrogen (H I and H₂, respectively) in the gas disc and bulge of the galaxy. The gas, once it cools, is assumed to settle in an exponential disc of half-mass radius, $r_{\text{gas, disc}}$. In BR06, the ratio of the molecular to atomic hydrogen gas surface density in galaxies is a function of the local hydrostatic pressure in the mid-plane of the disc, with a power-law index close to 1,

$$R_{\text{mol}} \equiv \frac{\Sigma_{\text{H}_2}}{\Sigma_{\text{H}_1}} = \left(\frac{P}{P_0}\right)^{\alpha_P}, \quad (1)$$

where P_0 and α_P are parameters measured in observations and have values $P_0/\kappa_B = 1500\text{--}40\,000\text{ cm}^{-3}\text{ K}$ and $\alpha_P \approx 0.7\text{--}1$ (Blitz & Rosolowsky 2006; Leroy et al. 2008). The hydrostatic pressure from the surface densities of gas and stars is calculated following (Elmegreen 1989)

$$P = \frac{\pi}{2} G \Sigma_{\text{gas}} \left(\Sigma_{\text{gas}} + \frac{\sigma_{\text{gas}}}{\sigma_*} \Sigma_* \right), \quad (2)$$

where Σ_{gas} and Σ_* are the total gas (atomic, molecular, and ionized) and stellar surface densities, respectively, and σ_{gas} and σ_* are the gas and stellar velocity dispersions. The stellar surface density is assumed to follow an exponential profile with a half-mass stellar radius of

$r_{*, \text{disc}}$. We adopt $\sigma_{\text{gas}} = 10\text{ km s}^{-1}$ and calculate $\sigma_* = \sqrt{\pi G h_* \Sigma_*}$, where $h_* = r_*/7.3$ (Kregel, Van Der Kruit & Grijps 2002), with r_* being the half-stellar mass radius. The H I surface densities cannot extend to infinitely small surface densities because the UV background will ionize very low density gas; thus a minimum threshold of $\Sigma_{\text{thresh}} = 0.1\text{ M}_\odot\text{ pc}^{-2}$ is applied, following the results of the hydrodynamical simulations of Gnedin (2012). All the gas at lower densities is considered to be ionized.

In order to understand how the default SHARK ISM prescription works against another available ISM model in SHARK, we carry out another run using an alternative prescription – in this case, Gnedin & Draine (2014), hereafter referred to as GD14. The GD14 model uses the dust-to-gas ratio, D_{MW} , and the local radiation field, U_{MW} , with respect to that of the solar neighbourhood, to estimate the ratio of H I to H₂ in the gas disc. These two parameters are estimated as $D_{\text{MW}} = Z_{\text{gas}}/Z_\odot$ and $U_{\text{MW}} = \Sigma_{\text{SFR}}/\Sigma_{\text{MW}}$, where Z_{gas} is the metallicity of the ISM. The values $Z_\odot = 0.134$ (Asplund et al. 2009) and $\Sigma_{\text{MW}} = 2.5\text{ M}_\odot\text{ yr}^{-1}$ (Bonatto & Bica 2011) are estimates from the solar neighbourhood. Hence, D_{MW} and U_{MW} are quantities that vary with galaxy properties. Using the argument presented in Wolfire et al. (2003), where it is stated that the pressure balance between the warm and the cold neutral media can only be achieved if the density is larger than a minimum density, we can approximate the minimum density to be proportional to U_{MW} . Hence, assuming that the pressure equilibrium between warm/cold media is a necessity for the formation of ISM, then U_{MW} will be proportional to ρ_{gas} , with ρ_{gas} being the gas density. As galaxies show an almost constant σ_{gas} , it can be assumed that the gas scale height is also close to constant, which allows us to replace ρ_{gas} by Σ_{gas} above. Based on D_{MW} and U_{MW} we calculate R_{mol} following (Gnedin & Draine 2014)

$$R_{\text{mol}} = \left(\frac{\Sigma_{\text{gas}}}{\Sigma_{R=1}}\right)^{\alpha_{\text{GD}}}, \quad (3)$$

where

$$\alpha_{\text{GD}} = 0.5 + \frac{1}{1 + \sqrt{U_{\text{MW}} D_{\text{MW}}^2/600}}, \quad (4)$$

$$\Sigma_{R=1} = \frac{50\text{ M}_\odot\text{ pc}^{-2}}{g} \frac{\sqrt{0.01 + U_{\text{MW}}}}{1 + 0.69\sqrt{0.01 + U_{\text{MW}}}}, \quad (5)$$

and

$$g = \sqrt{D_{\text{MW}}^2 + D_*^2}. \quad (6)$$

Here, $D_* \approx 0.17$ for scales $> 500\text{ pc}$.

Independent of how the H I/H I/H₂ is computed, our default SF model assumes the SFR surface density to be proportional to the H₂ surface density. The SFR surface density is then calculated by assuming a constant depletion time for the molecular gas, following:

$$\Sigma_{\text{SFR}} = \nu_{\text{SF}} f_{\text{mol}} \Sigma_{\text{gas}}. \quad (7)$$

Here, ν_{SF} is the inverse of the H₂ depletion time-scale with $f_{\text{mol}} = \Sigma_{\text{mol}}/\Sigma_{\text{gas}}$, where Σ_{mol} is the molecular gas surface density and Σ_{gas} is the total gas surface density; Σ_{SFR} is integrated over a radii range of $0\text{--}10\text{ }r_{\text{gas, disc}}$. Equation (7) applies to both the BR06 and GD14 models. Note that two different values of ν_{SF} are adopted in SHARK. For SF in discs, $\nu_{\text{SF}} = 1\text{ Gyr}^{-1}$, while for starbursts triggered by galaxy mergers and disc instabilities, $\nu_{\text{SF}} = 0.1\text{ Gyr}^{-1}$ (these values are based on Sargent et al. 2014). This is motivated by the bimodality observed in the $\Sigma_{\text{SFR}} - \Sigma_{\text{mol}}$ plane for normal star-forming galaxies and starbursts (Genzel et al. 2010).

2.3.2 AGN feedback

AGN feedback influences the amount of gas that cools and hence replenishes the ISM content of galaxies. The default AGN feedback model used in SHARK is that of Croton et al. (2016), hereafter referred to as **Croton16**. **Croton16** assumes a Bondi–Hoyle (Bondi 1952) like accretion mode

$$\dot{M}_{\text{BH, hh}} = 2.5\pi G^2 \frac{m_{\text{BH}}^2 \rho_0}{c_s^3}, \quad (8)$$

where c_s and ρ_0 are the sound speed and average density of the hot gas in the halo that accretes on to the SMBH, respectively, where $c_s \approx V_{\text{vir}}$ and V_{vir} is the halo’s virial velocity. $\dot{M}_{\text{BH, hh}}$ is the accretion rate calculated for the hot-halo mode, as described below. ρ_0 is calculated by equating the sound traveltime across a shell of diameter twice the Bondi radius to the local cooling time. This is also termed the ‘maximal cooling flow’ by Nulsen & Fabian (2000), which leads to

$$\dot{M}_{\text{BH, hh}} = \kappa_{\text{agn}} \frac{15}{16} \pi G \mu m_p \frac{\kappa_{\text{B}} T_{\text{vir}}}{\Lambda} m_{\text{BH}}. \quad (9)$$

κ_{agn} is a free parameter that was introduced in Croton et al. (2006) to counteract the approximations used to derive the accretion rate. κ_{B} and Λ are the Boltzmann constant and the cooling function that depends on T_{vir} and the hot gas metallicity. From equation (9), we can estimate the BH luminosity (L_{BH}) in this accretion mode, which in turn is used to calculate the heating provided by the BH for the halo as shown

$$\dot{M}_{\text{heat}} = \frac{L_{\text{BH}}}{0.5V_{\text{vir}}^2}, \quad (10)$$

where $L_{\text{BH}} = \eta \dot{M}_{\text{BH, hh}} c^2$, with η and c being the luminosity efficiency (based on Lagos, Padilla & Cora 2009) and speed of light, respectively.

To understand the effect of the AGN feedback, we vary the value of the free parameter κ_{agn} between 0 (no AGN feedback) to 1. Note that the default value in SHARK is 0.002.

2.3.3 Stellar feedback

The stellar feedback in SHARK is separated into two main components: the outflow rate of the gas that escapes from the galaxy, \dot{m}_{outflow} , and the ejection rate of the gas that escapes from the halo, \dot{m}_{ejected} . Lagos et al. (2018) describe $\dot{m}_{\text{outflow}} = \psi f(z, V_{\text{circ}})$, where ψ is the instantaneous SFR, z is the redshift and V_{circ} is the maximum circular velocity of the galaxy, where the ejection rate is > 0 only when the total injected energy of the outflow is greater than the binding energy of the halo. The terminal wind velocity, V_w , is based on the FIRE simulation suite (Muratov et al. 2015)

$$\frac{V_w}{\text{kms}^{-1}} = 1.9 \left(\frac{V_{\text{circ}}}{\text{kms}^{-1}} \right)^{1.1}. \quad (11)$$

The terminal wind velocity is required to compute the excess energy that will be used to eject the gas out of the halo:

$$E_{\text{excess}} = \epsilon_{\text{halo}} \frac{V_w^2}{2} f(z, V_{\text{circ}}), \quad (12)$$

where ϵ_{halo} is a free parameter. The net ejection rate can then be calculated as

$$\dot{m}_{\text{ejected}} = \frac{E_{\text{excess}}}{V_{\text{circ}}^2/2} - \dot{m}_{\text{outflow}}. \quad (13)$$

If $\dot{m}_{\text{ejected}} < 0$ no ejection from the halo takes place and we limit $\dot{m}_{\text{outflow}} = E_{\text{excess}}/(V_{\text{circ}}^2/2)$.

In SHARK-ref, we use the modelling presented in Lagos, Lacey & Baugh (2013), referred to as Lagos13, where they follow the evolution of the expansion of SNe driven bubbles from an early epoch of adiabatic expansion to the momentum-driven phase of expansion. They used this model to estimate \dot{m}_{outflow} and find

$$f = \epsilon_{\text{disc}} \left(\frac{V_{\text{circ}}}{v'_{\text{hot}}} \right)^\beta, \quad (14)$$

$$v'_{\text{hot}} = v_{\text{hot}}(1+z)^{z_p}. \quad (15)$$

SHARK-ref uses the default values as described in Lagos et al. (2018) with $\epsilon_{\text{disc}} = 1$ and $z_p = 0.12$. We vary the value of β from 0.5 to 5 in increments of 1, with the default value in SHARK-ref being 4.5, to understand how stellar feedback influences the amount of HI in haloes. For the no-stellar-feedback run, we set $\epsilon_{\text{disc}} = 0$.

2.3.4 Photoionization feedback

Photoionization feedback refers to the feedback arising from the ionizing radiation background produced by the first generation of stars, galaxies, and quasars during the epoch of reionization. The large ionizing radiation density affects small haloes, keeping the baryon temperature higher than the virial temperature, thus suppressing radiative cooling.

SHARK-ref follows the results of the one-dimensional collapse simulations of Sobacchi & Mesinger (2013), which suggest that the effects of reionization can be captured by allowing only those haloes that satisfy a redshift-dependent threshold velocity to be occupied. SHARK-ref use the *Sobacchi & Mesinger parametric* form, as adapted by Kim et al. (2015), which depends on the halo’s V_{circ} based on the spherical collapse model of Cole & Lacey (1996) instead. This predicts $M_{\text{halo}} \propto V_{\text{circ}}^3$. Thus, haloes with circular velocities below $v_{\text{thresh}}(z)$ are not allowed to cool their halo gas, where

$$v_{\text{thresh}}(z) = v_{\text{cut}}(1+z)^{\alpha_v} \left[1 - \left(\frac{1+z}{1+z_{\text{cut}}} \right)^2 \right]^{2.5/3}. \quad (16)$$

Here, v_{cut} , z_{cut} , and α_v are free parameters that are constrained by the Sobacchi & Mesinger (2013) simulation. We use different v_{cut} values, ranging from 20 to 50 km s⁻¹, to study the effect on the HI content of the haloes. The default value in SHARK-ref is 35 km s⁻¹. We keep the other two parameters fixed to $z_{\text{cut}} = 10$ and $\alpha_v = -0.2$, which are default in SHARK-ref.

2.3.5 Gas stripping in satellite galaxies

Following the model of ‘instantaneous ram-pressure stripping’ described in Lagos et al. (2014), SHARK assumes that as soon as galaxies become satellites, their halo gas is instantaneously stripped and transferred to the hot gas of the central galaxy, a process that is commonly referred to as ‘strangulation’. Thus, gas can only accrete on to the central galaxy in the halo and not on to satellite galaxies. Cold gas in the discs of galaxies is not stripped. SHARK also allows us to switch off this process, in turn assuming that satellite galaxies can retain their hot halo gas, and hence their ISM can continue to be replenished for some time, until their halo gas reservoir is exhausted. We note that the quenching of satellites also happens in this case as satellite subhaloes, where satellite galaxies reside, are cut-off from cosmological accretion, and hence their halo gas reservoir is not replenished. We test the effect of turning on and off the ‘instantaneous ram pressure stripping’ on the overall HI mass contained in haloes, with stripping ‘on’ being used in SHARK-ref.

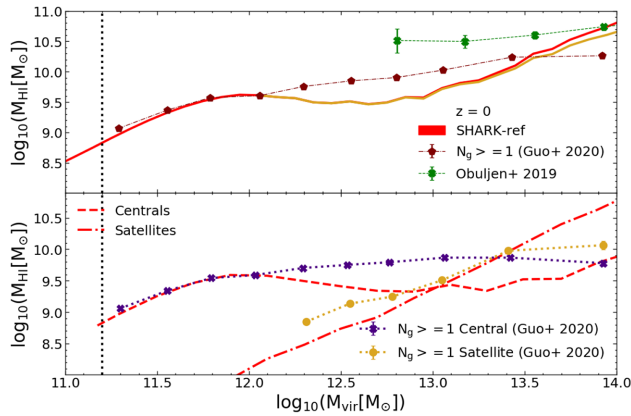


Figure 1. The mean of the total H I content in haloes as a function of halo mass at $z = 0$. In the upper panel, the red line shows predictions from SHARK-ref, with the vertical dashed line showing the convergence point between micro-SURFS and medi-SURFS. The yellow line shows the H I contained in subhaloes that are associated with the host-halo and are within one virial radius of the host halo. The symbols with error bars show the observed values of H I shown in Guo et al. (2020) and Obuljen et al. (2019), as labelled. Note that SHARK-ref predicts the H I content of $N_g \geq 1$ reasonably well until $M_{\text{vir}} \approx 10^{12} M_{\odot}$, with all the points agreeing with SHARK-ref, at which point SHARK-ref starts to deviate from the Guo et al. (2020) and Obuljen et al. (2019) points, either overpredicting or underpredicting the content at various points. The lower panel shows the central and satellite H I contribution from Guo et al. (2020) compared with SHARK-ref. We see the centrals agreeing with SHARK-ref until $M_{\text{vir}} \approx 10^{12} M_{\odot}$, but the satellite population agrees reasonably well with SHARK-ref over the entire range.

Regardless of whether the stripping is ‘on’ or ‘off’, the gas that is ejected from satellite galaxies due to stellar feedback is transferred to the ejected gas reservoir of the central galaxies, and hence that gas cannot be reincorporated into the hot halo gas of the satellites.

3 VALIDATION OF THE SHARK MODEL AGAINST LOCAL UNIVERSE H I OBSERVATIONS AND PREVIOUS MODELS

In this section, we describe how the total H I in the haloes compares with available observations, with the aim of validating the model before we analyse in detail what drives the shape and scatter of the HIHM relation. In particular, we compare with the observed HIHM relation (Section 3.1) and H I clustering (Section 3.2). We remind the reader that previous papers have shown that SHARK-ref reproduces well the H I mass function, H I–stellar mass scaling relation (Lagos et al. 2018), H I mass and velocity width distributions and the H I mass–velocity width relation observed in ALFALFA (Chauhan et al. 2019).

3.1 The local Universe HIHM relation

In Fig. 1, we compare the results from SHARK-ref with observations. We use the results shown in Guo et al. (2020), where they calculate the H I content of groups from the Sloan Digital Sky Survey DR7 Main Galaxy (SDSS; Lim et al. 2017) sample by stacking the H I spectra obtained from ALFALFA survey. SDSS is a major multispectral and spectroscopic redshift survey that covers over 35 per cent of the sky. We use data from the main SDSS galaxy survey, which is sensitive to 17.77 r -band magnitude. The ALFALFA (Arecibo Legacy Fast ALFA) survey, on the other hand, is a blind H I survey

covering 6900 deg² in the Northern hemisphere, with $\sim 31\,000$ direct H I detections (Giovannelli et al. 2005; Haynes et al. 2018) and going out to redshift $z = 0.06$.

Guo et al. (2020) use the SDSS DR7 group catalogue to identify galaxies with available spectroscopic redshifts, which is about 98 per cent complete. The halo masses of these groups were calculated using the proxy of galaxy stellar mass, with the halo radius, r_{200} , estimated from the definition that the mean mass density within r_{200} is 200 times the mean density of the universe at a given redshift. For stacking the H I for these groups and galaxies, they use ALFALFA IDL (see Fabello et al. 2011), which integrates over a square aperture and returns the H I spectrum. They have used $2 r_{200}$ as the aperture for groups, with 200 kpc being the apertures for centrals. They were able to extract 25 906 group spectra and 25 868 central spectra for their analysis. We present their final sample (with an occupancy number $N_g \geq 1$), which includes all the haloes with 1 or more galaxies in it, and compare with SHARK-ref.

We also use the data from Obuljen et al. (2019), who estimate the H I masses in dark matter haloes by directly integrating the H I mass functions over the available range of H I masses. Obuljen et al. (2019) model the abundance and clustering of neutral hydrogen through a halo-model based approach, where they parametrize the HIHM relation as a power law with an exponential mass cut-off (see equation 6 in Obuljen et al. 2019). In contrast to Guo et al. (2020) and Obuljen et al. (2019) do not directly measure the H I content of haloes, but instead use empirical relations to derive it. There is clearly some tension between these two approaches because they appear to be more than 2σ away from each other at $M_{\text{halo}} > 10^{13} M_{\odot}$. Some of this may be due to the SDSS group catalogue not sampling the high halo mass end with enough statistics, as well as the Obuljen et al. (2019) model not correctly capturing the H I mass in the massive haloes (where the H I content of galaxies is generally undetected by ALFALFA).

In the upper panel of Fig. 1 compares SHARK-ref with observations. We calculate the error on the mean H I content of SHARK-ref haloes via bootstrapping. The error is too small to be noticeable in the plot shown here. The observational data plotted are taken from Guo et al. (2020) and Obuljen et al. (2019). It can be seen that SHARK-ref is consistent with the H I mass content of groups until $M_{\text{vir}} < 10^{12} M_{\odot}$. For the H I-stacking points with $M_{\text{vir}} > 10^{12} M_{\odot}$, SHARK-ref consistently underpredicts H I in haloes, while it overpredicts it for $M_{\text{vir}} > 10^{13.2} M_{\odot}$. The inferred relation of Obuljen et al. (2019) seems to be flatter than our predictions, which results in the model under-(over-) predicting the H I content of haloes at $M_{\text{halo}} < (>) 10^{13.8} M_{\odot}$.

In the lower-panel of Fig. 1, we compare the H I contribution from the satellite and central populations to the H I content of haloes at $z = 0$. We also show the H I-stacking results for the contribution of H I from centrals and satellites as presented in (Guo et al. 2020). The errorbars for centrals (from the observational data) are the values presented in Guo et al. (2020). We estimate the errors, Δ , for the satellites from those reported for the total H I and central galaxy contributions as $\Delta_{\text{sat}} = \sqrt{\Delta_{\text{total}}^2 + \Delta_{\text{central}}^2}$, with Δ_{total} and Δ_{central} being the errors calculated for the total H I content of the halo and centrals, respectively. We find that the observed centrals data are consistent with the SHARK-ref predictions until $M_{\text{vir}} < 10^{12} M_{\odot}$, and thereafter SHARK-ref underpredicts the H I contained in the centrals. The satellites data, in contrast, are in better agreement with SHARK-ref predictions.

Note that the relation derived in Fig. 1 has not taken into account limitations that are inherent in observational surveys. Bravo et al. (2020), using a SHARK-derived lightcone to produce an analogue of the Galaxy and Mass Assembly (GAMA) survey (e.g.

Robotham et al. 2011), showed that assigning galaxies to groups and classifying them as centrals and satellites in the same way as is done in observations has an important impact on how we understand satellite/central galaxy quenching (also see Stevens & Brown 2017). This is because ~ 15 percent of satellites/centrals are wrongly classified as such (according to the intrinsic definition provided by the halo/subhalo catalogue). In this work, we compare directly VELOCIRAPTOR groups to the stacking results of Guo et al. (2020) without considering the effects shown in Bravo et al. (2020). Because the SDSS group catalogue used by Guo et al. (2020) is expected to have an even higher contamination than the GAMA groups analysed by Bravo et al. (2020) (see Robotham et al. 2011 for details), we expect this to play an even greater role in our comparison. In future work, we will make a detailed comparison with observations by mimicking the HI stacking procedure, with the aim of quantifying the systematic effects above. As is shown in Chauhan et al. (2019), accounting for observational limitations and producing mock-catalogues for comparison is essential when comparing simulations with observational data.

After comparing with the observations, we compare SHARK against other SAMs, such as GALFORM (Cole et al. 2000) and GAEA (Galaxy Evolution and Assembly; Xie et al. 2017). Baugh et al. (2019) analysed the HIHM relation in a recalibrated GALFORM variant, using the Planck Millennium N -body simulation, which is the latest addition to the ‘Millennium’ series of simulations of structure formation. For reference, Planck Millennium has a DM particle mass of $2.12 \times 10^9 M_\odot$ and a box of length $542.6 h^{-1}$ cMpc (Baugh et al. 2019).

GAEA on the other hand was run on the Millennium I (Springel et al. 2005) and Millennium II simulations (Boylan-Kolchin et al. 2009), whose DM particle masses are 1.7×10^{10} and $1.4 \times 10^8 M_\odot$, respectively, in boxes of length of 500 and $100 h^{-1}$ cMpc, respectively. We also compare to the HIHM relation derived from the hydrodynamical simulation Illustris-TNG100 (Nelson et al. 2018; Pillepich et al. 2018), which is publicly available (Nelson et al. 2019). This simulation has a box size of $75 h^{-1}$ cMpc and a DM particle mass of $7.5 \times 10^6 M_\odot$. The HI content of Illustris-TNG100 galaxies was calculated in post-processing, following the ‘inherent’ method outlined in Stevens et al. (2019), using the Gnedin & Draine (2014) prescription. We exclusively sum the HI masses of Illustris-TNG100 galaxies within R_{vir} to calculate a halo’s total HI mass. In other words, we intentionally exclude any HI contribution from the CGM. This makes the results from Illustris-TNG directly comparable to SAMs, which do not include HI in the CGM by design. We also compare with the semi-empirical HIHM relation described in Padmanabhan & Kulkarni (2017), which was derived at $z \sim 0$ by abundance matching dark matter haloes with HI-selected galaxies. They use the HI-mass function from HIPASS (Meyer et al. 2004) and ALFALFA (Martin et al. 2012) along with the Sheth & Tormen (2002) dark matter halo-mass function to match the HI-selected galaxies to dark matter haloes. They assume that each dark matter halo hosts one HI galaxy with its HI mass is proportional to the host dark matter halo mass. By construction, this means that the most massive HI galaxies inhabit the most massive haloes.

In Fig. 2, we plot the median of the total HI content as a function of halo mass for the SAMs, SHARK, GALFORM (Baugh et al. 2019), and GAEA (Spinelli et al. 2020); the hydrodynamical simulation Illustris-TNG100; and the empirical relation by Padmanabhan & Kulkarni (2017).

Both GALFORM and SHARK predict qualitatively similar curves, which display a prominent dip in the median HI mass of haloes at intermediate masses. The exact mass at which the dip happens

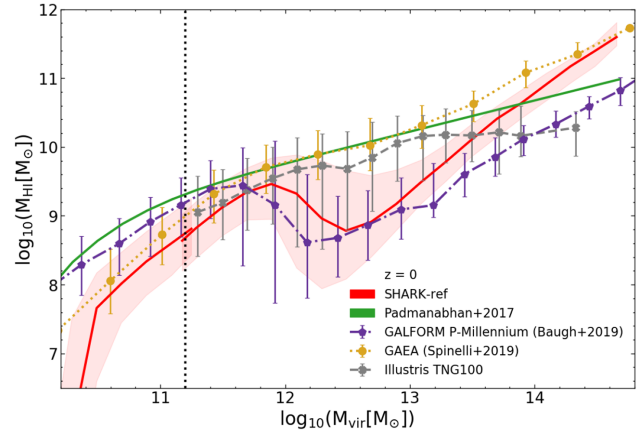


Figure 2. The median and the 16th–84th percentile range of HI content of haloes as a function of halo mass by GALFORM (Baugh et al. 2019), GAEA (Spinelli et al. 2020), TNG100 (Stevens et al. 2019), and SHARK (Lagos et al. 2018). The dip in the median HI mass occurs at lower halo masses for GALFORM than for SHARK-ref, though for GAEA and TNG100 we do not see a prominent dip at all. This is an effect of different AGN feedback and SF models implemented by the different SAMs presented here, as the strength of the AGN feedback affects the position and shape of the drop. As for TNG100, the dip and HI value depends on how it has been calculated, as for the current comparison, the CGM contributions to the HI in the haloes has been removed from the TNG100 to make it more comparable with the SAMs presented. The purple-dotted line with the errorbars are the HI values as shown in Baugh et al. (2019), with the errorbars showing the 10th–90th percentile range of the distribution, whereas the yellow-dotted line represents the values obtained from GAEA with the errorbars showing the 16th–84th percentile range of the distribution. The grey dashed line represents TNG100, with errorbars showing the 16th–84th percentile range of the distribution. The solid green line represents the HI–halo scaling relation developed by Padmanabhan & Kulkarni (2017). The red solid line is the prediction from SHARK-ref with the shaded region representing the 16th–84th percentile range of the distribution.

differs between the models, with GALFORM predicting this to take place at $M_{\text{halo}} \approx 10^{12} M_\odot$, while for SHARK this happens at $M_{\text{halo}} \approx 10^{12.5} M_\odot$. At lower (higher) halo masses, GALFORM predicts a higher (lower) median HI mass than SHARK. GAEA on the other hand, displays a very weak dip in the median HI mass with halo mass.

The Padmanabhan & Kulkarni (2017) semi-empirical relation by construction shows a monotonically increasing HI mass versus halo mass. This behaviour is qualitatively very different to the SAMs shown here, particularly SHARK and GALFORM. We show in Section 4.1 that the non-monotonic relation between the HI and halo mass is due to the modelling of AGN feedback. The difference in the sharpness of the dips seen in SHARK and GALFORM is due to the AGN feedback modelling used in the SAMs. As mentioned in Section 2.3.2, SHARK uses the Croton et al. (2016) model for AGN feedback, where the BH heating is estimated based on the luminosity of the BH, which is then used to adjust the cooling rate to respond to the heating. The heating radius is then estimated based on the radius within which the energy injected by the AGN equals that of the halo gas internal to that radius that would be lost if the gas were to cool. Whereas when looking at the AGN feedback in GALFORM, which is based on the Bower et al. (2006) model, AGNs are assumed to quench gas cooling only if the available AGN power is comparable to the cooling luminosity. The latter makes the AGN heating a binary mode, resulting in a sharper transition in GALFORM. GAEA produces massive galaxies that are less quenched than observations suggest at stellar masses $> 10^{10} M_\odot$ (see e.g. fig. 3 in Xie et al. 2020), which

may be an indication that their AGN feedback is not efficient enough. Illustris-TNG100, on the other hand, displays a mild dip at around $M_{\text{halo}} \approx 10^{12.5} M_{\odot}$, but much weaker than that displayed in SHARK and GALFORM. This dip goes away when we include the CGM HI contribution in the total HI mass of the haloes (not shown here), strongly suggesting that the CGM makes up a non-negligible amount of the HI in groups. Unlike SAMs, Illustris-TNG100 predicts a flat median HI mass at $M_{\text{halo}} \gtrsim 10^{13} M_{\odot}$. We caution that the definition of M_{halo} is not the same in all these simulations, but differences in definitions are much smaller ($\lesssim 0.2$ dex) than the differences seen here in the position of the HI mass dip. The abrupt drop in the HI abundance of haloes at $M_{\text{vir}} \lesssim 10^{10.4} M_{\odot}$ is caused by the strength of the UV background being sufficient to keep the gas in those low-mass haloes ionized (see Section A1 for more details).

When looking at the scatter around the median HIHM relation for all simulations, we find that all galaxy formation simulations shown here (SHARK, GALFORM, GAEA, and Illustris-TNG100) agree in that the scatter is maximal at $M_{\text{halo}} \approx 10^{12-13} M_{\odot}$, although the exact mass at which this occurs, and the magnitude of the scatter, varies from simulation to simulation. SHARK, GALFORM, and Illustris-TNG100 produce a similarly large scatter ($\approx 1-1.5$ dex) at around the position where the dip in HI mass takes place, while GAEA predicts a much smaller scatter of ≈ 0.3 dex. This shows that observational constraints on the scatter of the HIHM relation are essential if we are to judge the success of the models.

3.2 The H I correlation function

The correlation function is defined as the excess clustering of a target distribution of galaxies over a random distribution, and thus is a measure of the spatial distribution of galaxies. It encodes information about both the underlying cosmology and the physics of galaxy formation, and its form is subject to how galaxies are selected (e.g. optically selected or HI selected).

We use the $z = 0$ medi- and micro-SURFS boxes to measure the projected two-point correlation function (2PCF) of galaxies with HI masses $> 10^8 M_{\odot}$ for medi-SURFS and HI masses $> 10^7 M_{\odot}$ for micro-SURFS. We employ the CorrFunc² (Sinha & Garrison 2020) PYTHON routine developed to compute correlation functions and other clustering statistics for simulated and observed galaxies, as follows:

$$\frac{w_p(r_p)}{r_p} = \frac{2}{r_p} \int_0^{\pi_{\text{max}}} \xi(r_p, \pi) d\pi. \quad (17)$$

Here, we have measured the correlation function as a two-dimensional histogram, $\xi(r_p, \pi)$, with the count of galaxy pairs as a function of both projected separation (r_p) and line-of-sight separation (π). By integrating $\xi(r_p, \pi)$ over π , we can account for the effect of peculiar velocities. The π_{max} values adopted for our micro- and medi-SURFS boxes are 10 and 30 h^{-1} cMpc, respectively. Different π_{max} values are used to incorporate the different box sizes of micro- and medi-SURFS. These values reproduce the observational measurements of Papastergis et al. (2013) and Meyer et al. (2007), with medi-SURFS using the same π_{max} values as were used in the observations. As for micro-SURFS, we opted for a lower π_{max} value because of the relatively small volume of the simulation box, which impacts the strength of clustering (Power & Knebe 2006).

In Fig. 3, we reproduce the clustering measurements using the criteria used by Papastergis et al. (2013) and Meyer et al. (2007),

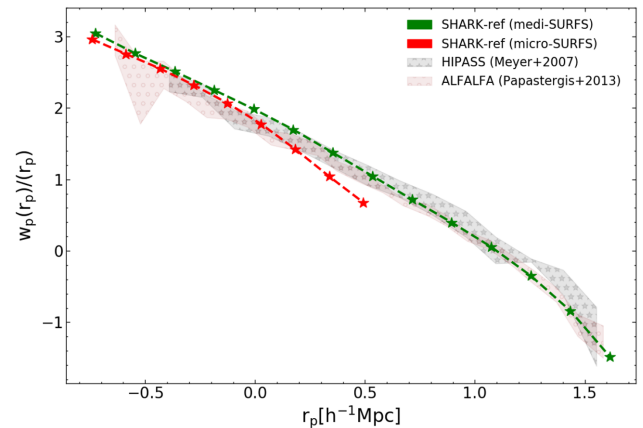


Figure 3. The projected two-point correlation function of the SHARK-ref model for micro-SURFS (red-dashed line) and medi-SURFS (green-dashed line) compared with the observations of Meyer et al. (2007) (grey-shaded region with stars) and Papastergis et al. (2013) (brown-shaded region with circles), for the HIPASS and ALFALFA 40 per cent surveys, respectively. There is good agreement between the predictions and observations within the errorbars. For micro-SURFS, the predictions deviate at $r_p \gtrsim 1 h^{-1}$ Mpc due to the small size of the simulated box.

and show the predicted 2PCF of HI selected galaxies in SHARK in both simulated SURFS boxes, micro-SURFS and medi-SURFS. We also show the observational measurements of Meyer et al. (2007) using HIPASS and Papastergis et al. (2013) using ALFALFA. Both these observational measurements apply a volume correction and hence are comparable to the 2PCF obtained from the simulated box, which by construction is volume-limited. Meyer et al. (2007) adopted a higher mass threshold of $M_{\text{HI}} \approx 10^9 M_{\odot}$ for their analysis, whereas Papastergis et al. (2013) utilize the entire 40 per cent ALFALFA data sample, with the HI masses limiting to $M_{\text{HI}} > 10^{7.5} M_{\odot}$. Despite also using different M_{HI} limits for our different resolution boxes, we find agreement between them, although the micro-SURFS predictions start deviating at about $r_p \gtrsim 1 h^{-1}$ Mpc as a result of the small volume of micro-SURFS.

HIPASS and ALFALFA have different volumes and depth, and hence they are expected to trace different HI mass distributions. This can, in principle, lead to different clustering signals if the 2PCF is HI-mass dependent. Papastergis et al. (2013) and Meyer et al. (2007) tested this dependence and found that the clustering amplitude was largely insensitive to the HI mass (see however Guo et al. 2017 for a different conclusion). Crain et al. (2017) also tested the HI-clustering dependency on the HI mass of the galaxies in EAGLE hydrodynamical simulation (Crain et al. 2015; Schaye et al. 2015) by looking at the clustering measurements of galaxies belonging to the same stellar bin, and found that HI-poor galaxies seem to be more clustered. We tested this in our simulated boxes and found that the clustering amplitude was independent of the HI mass selection (not shown here). This is also the reason why micro- and medi-SURFS agree well in Fig. 3 despite having different HI mass lower limits.

4 THE PHYSICAL DRIVERS OF THE HIHM RELATION

In this section, we explore the physical processes that drive the shape and the scatter of the HIHM relation. In what follows, we compute a halo’s HI mass by summing over the HI masses of all galaxies embedded in that halo. Note that SHARK does not model the atomic content of the intrahalo gas and hence our measurement only reflects

²<https://github.com/manodeep/Corrfunc>

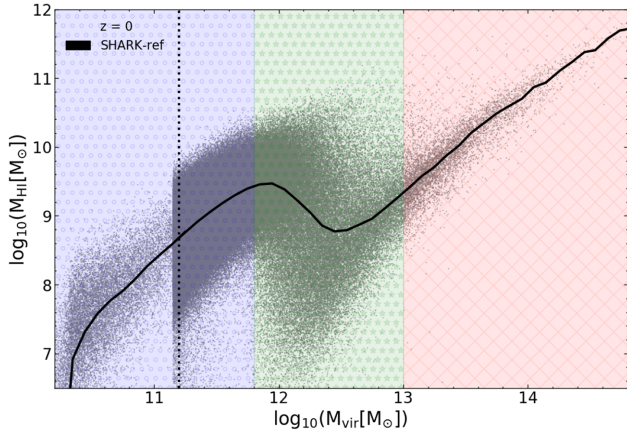


Figure 4. The HIHM relation in SHARK at $z = 0$. Each point is an individual halo, while the line shows the median of the relation. The three regions used to study the HIHM relation are shown with different shaded styles. The vertical dotted line represents the converging point of the two resolution boxes we are using – micro-SURFS and medi-SURFS.

the total HI content in the ISM of galaxies that belong to the same group.

In order to better understand the physical drivers of the HIHM relation, we divide the relation into three regions, as shown in Fig. 4:

(i) **Low-mass region:** includes haloes with $M_{\text{vir}} < 10^{11.8} M_{\odot}$. In this region the HI mass monotonically increases with halo mass. We show in Section 4.1 that here the majority of the HI content is in the central galaxy, with satellites contributing little to nothing, as many of these centrals are isolated (i.e. have no satellites).

(ii) **Transition region:** includes haloes with $10^{11.8} M_{\odot} \leq M_{\text{vir}} < 10^{13} M_{\odot}$. Here, the HI content of haloes displays a non-monotonic dependence on halo mass. In this region some haloes have most of their HI content in the central galaxy, while others are dominated by their satellites. As a result, this is the region of largest scatter.

(iii) **High-mass region:** includes haloes with $M_{\text{vir}} > 10^{13} M_{\odot}$. In this region, the HI mass returns to a monotonically increasing relation with the halo mass. Here, the majority of HI is contained in the satellite population.

4.1 Understanding the shape of the HIHM relation

In order to unveil the physical drivers behind the shape of the HIHM relation, we leverage on the flexibility and modularity of SHARK to explore different models and parameters for any one physical process. In this section, we show how the HIHM relation is affected by these variations and break down the analysis into the effect of different physical processes.

4.1.1 AGN feedback effect

As previously stated in Section 2.3.2, we vary the free parameter κ_{agn} (equation 9) that controls the strength of AGN feedback. In Fig. 5, we show how this efficiency affects the overall median HI content of the halo at $z = 0$. The different colours represent different values of κ_{agn} , with the shaded region representing the 16th–84th percentile range of the SHARK-ref model. We remind the reader that the vertical line demarcates the transition from the high resolution, small volume micro-SURFS box used at low-halo masses, to the moderate resolution, large volume, medi-SURFS box, used at high

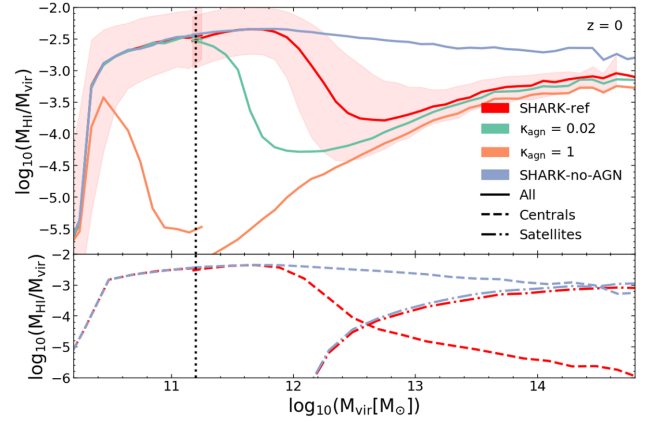


Figure 5. Median $M_{\text{HI}}/M_{\text{vir}}$ ratio as a function of M_{vir} . The shaded region represents the 1σ scatter on the median $M_{\text{HI}}/M_{\text{vir}}$ relation for our default (SHARK-ref) model, and other lines representing different strengths of the feedback (as labelled). κ_{agn} is the free parameter that regulates the AGN feedback efficiency (see equation 9); the higher the value, the stronger the feedback. It should be noted that as AGN feedback becomes more efficient, the knee of the relation shifts towards smaller virial masses, making AGN feedback efficiency a major contributor to the shape of the HIHM scaling relation. The vertical dotted line represents the shift from micro-SURFS (dashed-dotted lines) to medi-SURFS (solid lines). *Lower panel:* The median HI contribution from central and satellite galaxies to the total HI of the halo. For clarity we show this for the SHARK-ref and SHARK-no-AGN runs only. The centrals, which are major contributors to the HIHM relation at the transition region, are significantly affected by changes in the AGN feedback efficiency.

halo masses. This demarcation style is used throughout the figures in this paper, and has the purpose of increasing the dynamical range explored. Lagos et al. (2018) analysed the convergence between these two boxes and found that the stellar mass function was very well converged down to $10^8 M_{\odot}$ in medi-SURFS, while the HI mass function was converged at $10^{8.5} M_{\odot}$. We therefore adopt a transition between the boxes that roughly corresponds to these masses.

We find the $M_{\text{HI}}/M_{\text{vir}}$ ratio increases as M_{vir} increases, reaching a peak value and then rapidly dropping to a minimum (except for the SHARK-no-AGN run) to then gradually rise again. This drop corresponds to our *transition region* (for SHARK-ref), and is mostly influenced by the strength of the AGN feedback. As we move from $\kappa_{\text{agn}} = 0.002$ (the default in SHARK-ref), 0.02 and 1, the drop shifts from $M_{\text{vir}} = 10^{12} M_{\odot}$ to $10^{11.2} M_{\odot}$ to $10^{10.6} M_{\odot}$, respectively. As for the case of SHARK-no-AGN feedback, $\kappa_{\text{agn}} = 0$, we see that the ratio reaches a peak and then gradually decreases with the halo mass and this peak corresponds to the peak achieved by SHARK-ref model. This is because shock heating of the accreted gas on to haloes plays a role in slowing down the cooling in the more massive haloes and hence the replenishment of the ISM of central galaxies, producing the mild decrease in HI-to-halo mass ratio. It should be noted that despite the drop becoming steeper and taking place at lower halo masses with increasing AGN feedback efficiency, the HI contained in the haloes gradually rises up to similar values at the cluster regime ($M_{\text{vir}} > 10^{14.3} M_{\odot}$), which is a consequence of satellites dominating this regime. As for the smallest haloes, there is not much difference in their HI content as AGN feedback does not play a role here.

More efficient AGN feedback has the consequence of steepening the drop in the HI-to-halo mass ratio in the *transition region*, as this shifts to lower halo masses. This is driven by the fact that as AGN feedback becomes more efficient, gas cooling becomes

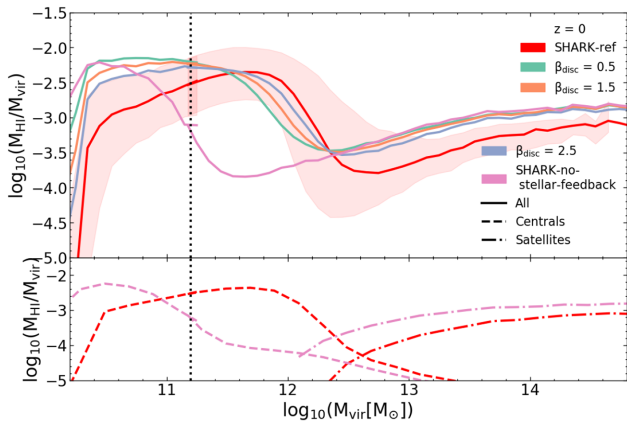


Figure 6. As in Fig. 5, but for different β_{disc} , which represents the power-law exponent in the circular velocity dependence of the mass loading due to stellar feedback (see equation 14). Although we see an effect of this parameter over the whole mass range, it is more prominent at low masses, with a weaker stellar feedback being associated with a higher H I-to-halo mass ratio. The impact of stellar feedback is anyway weaker than that of AGN. *Lower panel:* The median H I contribution from centrals and satellites to the total H I of the halo. For clarity we show the SHARK-ref and SHARK-no-stellar-feedback runs only. In the latter, centrals see a decrease of their H I content at very low halo masses compared to SHARK-ref.

extremely inefficient, hampering the replenishment of the ISM of central galaxies.

A related consequence is that satellite galaxies become more prominent H I reservoirs of the halo at lower halo masses as the AGN feedback efficiency increases, which can be seen in the lower panel of Fig. 5. The lower panel shows the central and satellite H I contributions for SHARK-ref and SHARK-no-AGN runs. We find that for the SHARK-no-AGN run, the centrals remain the primary H I reservoir of haloes as massive as $M_{\text{vir}} \approx 10^{14.6} M_{\odot}$, thereafter satellites become dominant. On the other hand, in the SHARK-ref run we find satellites start to become major H I contributors at much lower halo masses, $M_{\text{vir}} \approx 10^{12.5} M_{\odot}$.

4.1.2 Stellar feedback effect

Stellar feedback in SHARK is a two-step process: gas is first expelled from the galaxy, and then from the halo depending on the excess energy of the outflow compared to the bounding energy (discussed in detail in Section 2.3.3).

In the first step, the outflow rate from the galaxy depends on the maximum circular velocity of the galaxy to the power $-\beta_{\text{disc}}$. For reference, an energy conserved outflow should have a $\beta_{\text{disc}} = 2$, while a momentum-conserved outflow has $\beta_{\text{disc}} = 1$. Lagos et al. (2013) found that once outflows are followed throughout their evolution in the interstellar medium from the adiabatic expansion to the snow-plough phase, β_{disc} can take higher values, and in fact, SHARK-ref adopts $\beta_{\text{disc}} = 4.5$. Here, we vary the value of β_{disc} to examine the effect this has on the H I content of haloes.

In Fig. 6, we present the effect of varying β_{disc} on the $M_{\text{HI}}/M_{\text{vir}}-M_{\text{vir}}$ relation at $z = 0$. We change the value of β_{disc} from 0.5 to 5. The way β_{disc} affects the H I content of haloes is different at different halo masses. The H I content of haloes below the virial mass of $10^{11.2} M_{\odot}$ is affected the most, with higher β_{disc} values inducing a smaller amount of H I in the halo. A similar trend is seen in haloes above the mass $M_{\text{vir}} > 10^{12.6} M_{\odot}$.

These trends are caused by a higher value of β_{disc} driving higher outflow rates, and hence depleting the ISM of both centrals and satellites alike. In the transition region we see that a higher β_{disc} value is associated with *higher* H I-to-halo mass ratios. This at first appears counter-intuitive as more outflows should lead to a lower H I content. However, this can be reconciled by the fact that what drives this trend is the transition from H I being dominated by the central galaxy to the satellites moving towards lower halo masses as β_{disc} increases.

One interesting aspect of having no stellar feedback (SHARK-no-stellar-feedback), is seen in Fig. 6. The $M_{\text{HI}}/M_{\text{vir}}$ ratio is very similar to the $\kappa_{\text{AGN}} = 10$ run, i.e. very high AGN feedback efficiency (Section 4.1.1), for the H I content of the entire halo. With stellar feedback off, we end up with more elliptical galaxies at lower halo masses which is indicative of the galactic disc being unstable and unable to sustain itself. This leads to galaxies being bulge dominated at $M_{\text{stellar}} \gtrsim 10^{8.5} M_{\odot}$ compared to $M_{\text{stellar}} \gtrsim 10^{10} M_{\odot}$ in SHARK-ref. Because the BH mass scales with the bulge mass in SHARK, AGN feedback can now be effective in galaxies of much lower stellar masses compared to SHARK-ref. In short, AGN feedback becomes overly efficient in the absence of stellar feedback across the whole stellar mass range. A similar effect was noticed in the EAGLE hydrodynamical simulations (see Wright et al. 2020), where AGN feedback becomes much more efficient when there is no stellar feedback present. We also vary other parameters related to stellar feedback. In particular we tested varying $\epsilon_{\text{disc}} = 1, 3, 5, 7$, and 10. We find that the effect of changing the ϵ_{disc} has a similar effect on the $M_{\text{HI}}/M_{\text{vir}}-M_{\text{vir}}$ relation as varying β_{disc} .

In the lower panel of Fig. 6, the H I contribution of central and satellites is shown for the SHARK-ref and SHARK-no-stellar-feedback runs. The H I content of centrals decreases rapidly for the SHARK-no-stellar-feedback and starts at a lower halo mass of $M_{\text{vir}} \approx 10^{10.4} M_{\odot}$, whereas for the SHARK-ref centrals, the H I content starts decreasing at $M_{\text{vir}} \approx 10^{12} M_{\odot}$. We also find that the H I content of satellites in the SHARK-no-stellar-feedback run is more significant than in the SHARK-ref run relative to the total, with the satellites becoming a major H I contributors at lower halo masses.

Despite stellar feedback having a clear effect on the HIHM relation, it appears like AGN feedback has a more dramatic effect on the shape of the HIHM relation. This makes sense as stellar feedback hardly quenches galaxies but instead plays a role in the self-regulation of SF. AGN, on the contrary, is very efficient at quenching galaxies above a stellar mass threshold that in SHARK-ref happens roughly at $M_{\text{stellar}} \approx 10^{10.5} M_{\odot}$.

4.1.3 The effect from other physical mechanisms

In addition to stellar and AGN feedback, we explore other physical mechanisms in SHARK, which we present in Appendix A. These include photoionization feedback, ISM modelling, and environmental effects. These other mechanisms have a lesser effect on the HIHM relation compared to AGN and stellar feedback. Here, we provide short description of the main conclusions.

As stated in Section 2.3.4, we vary the value of v_{cut} , which directly affects the circular velocity (v_{thresh}) of the haloes under which the halo gas is not allowed to cool down and thus remains ionized (see equation 16). We find that changing v_{cut} does not have any effect on the drop seen in the *transition region*, which remains at the $M_{\text{vir}} \sim 10^{12} M_{\odot}$ mass scale for all the runs with varying v_{cut} . Though, a lower photoionization feedback efficiency does result in higher H I content for haloes of $M_{\text{vir}} > 10^{12.4} M_{\odot}$. This is caused by the fact

that with lower photoionization feedback smaller haloes are allowed to retain their H I content, and when they become satellites, their H I contribution to the total H I of a halo increases (see Fig. A1). See Appendix A1 for more details.

We also tested the effect of using different models for the molecular-to-atomic gas partition on the total H I content of the halo. We compared the BR06 (the default model of choice) and GD14 models for gas partition in the ISM (see Section 2.3.1). We find that the *transition region* for the model adopting the GD14 prescription occurs at lower halo masses, $M_{\text{vir}} \approx 10^{11.5} M_{\odot}$ compared to $M_{\text{vir}} \approx 10^{12} M_{\odot}$ for BR06. We find that this is due to the interplay between AGN feedback and the ISM model, as bigger BHs are produced in the GD14 run compared to BR06 at fixed halo mass in the transition region, again highlighting the complex interplay between the physical processes modelled in SHARK. We also find that using GD14 results in higher H I content for low- and high-mass haloes as opposed to BR06 (see Fig. A2). The latter is due to the fact that the centrals of low-mass haloes are more H I-rich in GD14 than BR06, which boosts the H I content of those, but also of high-mass haloes as they become satellites. We delve deeper into the ISM model effect in Appendix A2.

Finally, we test the ram-pressure stripping effect on the H I content of haloes, by switching between the stripping mode ‘on’ and ‘off’ (see Section 2.3.5). We find that the total amount of H I in either model is approximately the same, though the stripping ‘off’ model leads to a slightly lower H I in the transition region (see Fig. A3). More details on this effect are given in Appendix A3.

One major inference made through these tests was that despite the variations above, the shape of the HIHM relation essentially remained the same.

4.1.4 Summary

In conclusion, we find that several physical processes affect the shape of the $M_{\text{HI}}/M_{\text{vir}}-M_{\text{vir}}$ relation and therefore we cannot isolate a single process that is the sole contributor for this. We can, none the less, rank different processes by their apparent effect. By doing this we find that AGN feedback appears to have the strongest effect as the transition region changes shape dramatically with varying AGN feedback efficiency, and moreover, the existence of a transition region (regardless of its shape) seems to be solely determined by AGN feedback. We expect the exact way of modelling AGN feedback to also have an effect (though this is not tested explicitly here). Other physical processes, such as stellar feedback, the ISM modelling and photoionization feedback have a noticeable effect on the shape of the relation but qualitatively the relation continues to clearly have three distinct regions.

4.2 Physical drivers behind the scatter of the HIHM relation

The shape of the HIHM is only half the story. To fully characterize the HIHM scaling relation, we also need to understand the underlying scatter and its physical drivers. This is necessary for the purpose of Section 5, in which we aim to develop a numerical way of populating DM-only simulations with H I. For the latter, it is then important to explore how the scatter correlates with different halo properties which are accessible in these simulations. With this in mind, we explore how the scatter of the HIHM relation related to halo properties such as the halo mass assembly history, the halo’s spin parameter, etc, in the following sections. Here, we focus on the SHARK-ref model only.

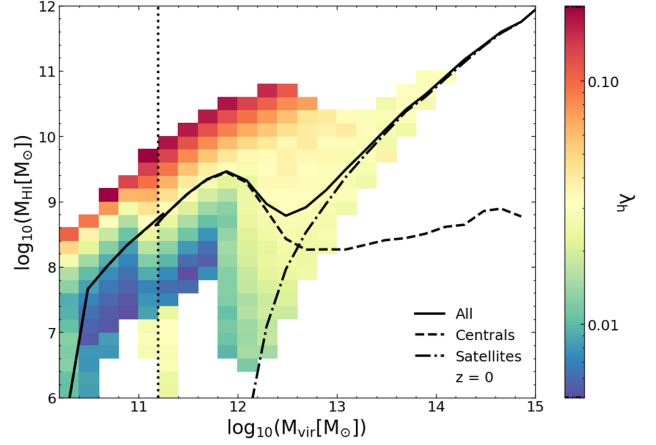


Figure 7. The HIHM relation of haloes in SHARK-ref at $z = 0$, with each bin being coloured by the median halo’s spin parameter, as labelled in the colour bar. The solid line represents the median H I mass of the halo as a function of M_{vir} , while the dashed and dash-dotted lines represent the central and satellite contributions, respectively. The vertical dotted line shows the transition from micro-SURFS to medi-SURFS at lower and higher halo masses, respectively. There is a strong correlation between the H I mass and the spin parameter at fixed halo mass for haloes with $M_{\text{vir}} < 10^{12} M_{\odot}$. Haloes with higher spin parameters are H I-rich than their counterparts. This trend becomes less prominent at the transition regions and completely disappears in the high-mass region.

4.2.1 Spin parameter effect

An intrinsic halo property that has recently been discussed in length in the literature in connection to the H I content of galaxies is the spin parameter. The spin parameter of a halo is normally quantified as follows (Peebles 1969):

$$\lambda = \frac{J\sqrt{|E|}}{GM^{5/2}}, \quad (18)$$

where J is the magnitude of the angular momentum vector of the particles within the virial radius, M is the virial mass, E is the total energy of the system, and G is the gravitational constant. Maddox et al. (2015) and Obreschkow et al. (2016) have suggested based on ALFALFA and THINGS (Walter et al. 2008) observations that the angular momentum of a galaxy regulates its H I mass and the atomic-to-baryon mass fraction; the idea being that a galaxy with high angular momentum can support a larger H I disc, thus sustaining more H I mass as well, compared to a lower angular momentum disc, which is subject to more instabilities. Empirically this has been observed as a correlation between the angular momentum, H I content, and physical extent (Lutz et al. 2018). Angular momentum in haloes scales steeply with mass, dependence that is removed when focusing instead on the spin parameter. Hence, for our purpose – studying what drives the scatter of H I content in haloes at fixed halo mass – the halo spin is a more natural property to focus on than angular momentum.

Fig. 7 shows the $M_{\text{HI}}-M_{\text{vir}}$ relation with bins in this space this time coloured by the median spin parameter of haloes. The halo’s spin parameter is very strongly correlated with the scatter in the HIHM relation at $M_{\text{vir}} < 10^{12} M_{\odot}$, with higher spin parameters being associated with more H I-rich haloes. The H I content in haloes at the low-mass region is primarily contributed by the central galaxy. Hence, the relation between the H I mass and spin parameter for haloes is pretty much a reflection of the relation between the H I content and angular momentum of the central galaxy.

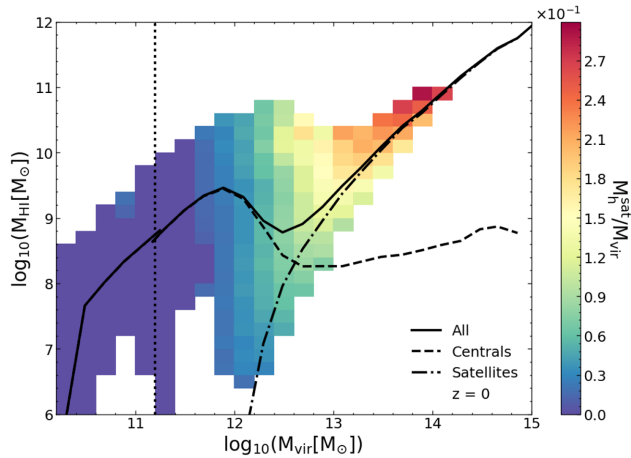


Figure 8. As in Fig. 7, but here bins are coloured by the median ratio between the total mass in subhaloes (M_h^{sat}) to the total halo mass (M_{vir}), as labelled in the colour bar. At $M_{\text{vir}} > 10^{12} M_{\odot}$, a correlation emerges with higher $M_h^{\text{sat}}/M_{\text{vir}}$ associated with a higher HI content at fixed halo mass.

We would like to caution our readers that we use the halo spin parameter as opposed to the spin of the galaxies and these can be very different. The cited observations have no access to the halo spin. The strong correlation seen in Fig. 7 could be exaggerated due to the simplistic model assumptions. For instance, SHARK assumes that the halo gas has the same specific angular momentum as the halo’s DM, with the specific angular momentum of the gas being conserved as it cools. SHARK also assumes the specific angular momentum of the galaxy’s components and halo to be aligned.

As we move towards the transition and high-mass regions, this correlation is no longer observed. This is because in these regions we see the emergence of the satellite population as the main contributors of HI in haloes and hence the relation between HI mass and angular momentum of the central galaxy is no longer relevant. Satellite galaxies on the other hand, have angular momenta which is largely uncorrelated with the host-halo’s spin. Satellite galaxies in SHARK have a specific angular momentum that is inherited from their hosthalo last time they were centrals. Due to the stochastic nature of the halo spin parameter, by $z = 0$ satellite galaxies have stellar spins, and therefore HI masses, that are uncorrelated with the central galaxy spin.

We also study the evolution of the HI–halo mass– λ relation towards high redshift, up to $z = 2$ (see Appendix D). We find that the trend remains prominent throughout the whole redshift range. We also find evidence of the *transition region* shrinking in dynamic range due to the systematic effect of AGN feedback efficiency decreasing as we move to higher redshifts.

4.2.2 Substructure mass effect

As stated in previous sections, satellite galaxies are the primary source of HI in haloes in the high-mass region. Hence, we expect the amount of substructure to be a good predictor of the scatter in the HIHM relation at high halo masses. To explore this idea, Fig. 8 shows the HIHM relation with bins now coloured by the fraction of mass in a halo that is contained in subhaloes, $M_h^{\text{sat}}/M_{\text{vir}}$. Note that here we use subhalo and halo masses of the VELOCIRAPTOR catalogues of the micro-SURFS and medi-SURFS.

We note that already at the transition region the effect of substructure on the HI content of haloes is visible, but certainly becomes clearer in the high-mass region, in a way that haloes with higher $M_h^{\text{sat}}/M_{\text{vir}}$ also have more HI. This is largely due to the larger number of satellites a halo with a higher $M_h^{\text{sat}}/M_{\text{vir}}$ has compared to one with a lower $M_h^{\text{sat}}/M_{\text{vir}}$ at fixed halo mass. The fact that the trend is weaker in the transition region than at $M_{\text{vir}} > 10^{12.5} M_{\odot}$ is due to the fact that many of those haloes have very few or no satellites. The clear correlation we obtain between the HI mass and $M_h^{\text{sat}}/M_{\text{vir}}$ at high halo masses makes it a good candidate to be used to predict the HI content of massive haloes.

We explore the evolution of the HI–halo mass relation dependence on $M_h^{\text{sat}}/M_{\text{vir}}$ over the redshift range $0 \leq z \leq 2$ in Appendix D, and find the trend to remain prominent and continue to be the main parameter that correlates with the scatter of the HI–halo mass relation at the high halo mass end ($\gtrsim 10^{13} M_{\odot}$).

4.2.3 Other Halo Properties

In addition to the halo parameters analysed here, we also explored the halo concentration and the effect of formation age (redshift at which the halo has assembled 50 per cent of its present mass) of the halo on the HI content. We found no correlation between the HI content of haloes and its concentration. This is due to the fact that SHARK adopts the concentration model of Duffy et al. (2008), which only depends on halo mass and time. Hence, naturally, at fixed halo mass, we obtain no dependence of the HI content on concentration.

It has been speculated in previous studies that the formation age of haloes, hereafter referred to as z_{50} , is correlated to their HI content (see Guo et al. 2017; Spinelli et al. 2020). When testing the effect of z_{50} with SHARK, we find that a slight trend is noticeable in the *transition region*, with younger haloes having more HI than their counterparts of the same mass (see Fig. B1). We discuss more on the effects of formation age on the HI content in Appendix B, and its relation to AGN feedback in Section 4.2.4.

4.2.4 Baryon physics effects

As stated previously (see Section 4.1.1), the dip in the HI–halo scaling relation (at $M_{\text{vir}} \approx 10^{12} M_{\odot}$) is caused by AGN feedback, which becomes prominent at these masses. AGN feedback is also responsible for the flaring of the scatter in the transition region, which increases from about 0.5 dex in the low-mass region to almost 1.2 dex at the transition region. As pointed out above, the halo spin parameter and $M_h^{\text{sat}}/M_{\text{vir}}$ are promising second variables to reduce the scatter at the low- and high-mass end regions, respectively.

For the transition region, however, a combination of these two parameters is required, as in this region we get both types of haloes, those that have their HI content mostly in their central, and those that have most of their HI in satellite galaxies. But even when including both parameters, we still cannot reduce the residual scatter to below 0.9 dex (discussed in detail in Section 5).

This is to be expected, as the exact effect of AGN feedback cannot be trivially predicted from halo properties only but instead we require insight into the BH mass and cooling luminosity. To better illustrate the effect of AGN feedback at the transition region, Fig. 9 shows the HIHM relation with bins coloured by the median ratio between the BH mass, M_{BH} , and stellar mass of the central galaxy. We find a stronger correlation between the halo HI mass with M_{BH}/M_{\star} at fixed halo mass than that seen with z_{50} , halo spin parameter and

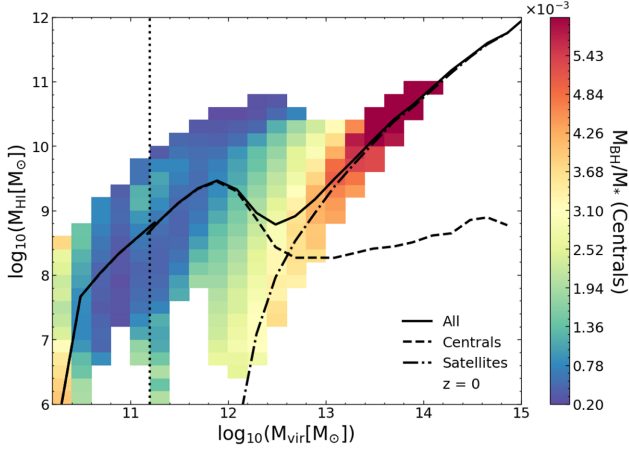


Figure 9. As in Figs 7 and 8, but bins here are coloured by the median of the fraction of BH mass (M_{BH}) to the central galaxy’s stellar mass (M_*) as labelled in the colour bar. A clear trend emerges in the transition region of more H I residing in haloes whose central has a low-mass BH relative to its stellar mass.

$M_{\text{h}}^{\text{sat}}/M_{\text{vir}}$. Haloes with low-mass BHs relative to the stellar mass of the central tend to have more H I mass compared to haloes with more massive BHs. However, there still is a causal relation between the AGN feedback efficiency and z_{50} . We find that at fixed halo mass, more massive BHs inhabit older haloes, and hence more powerful AGN feedback is possible in older haloes. We do, however, find the correlation between the scatter of the HIHM relation at fixed halo mass to be stronger with the BH mass than with z_{50} .

Despite the significance of the BH mass in reducing the residual scatter of the HIHM relation, we do not use it in Section 5 to build up our numerical model for how to populate haloes with H I. This is because we are interested in a model that can be applied to large-scale DM-only simulations. This analysis, however, serves to remind the reader that the complexity of baryon effects cannot be fully described with halo properties alone.

4.2.5 The H I content of subhaloes

In this section, we discuss how the H I mass inside the subhaloes is related to subhalo properties.

Section 4.2.1 showed that there is a strong correlation between the HIHM scatter and the spin parameter of the halo at fixed halo mass in the low-mass region. A possible interpretation of Fig. 7 is that the weakening of the correlation at $M_{\text{halo}} > 10^{11.8} M_{\odot}$ is due to the contribution of satellite galaxies becoming significant, and their subhalo’s spin being uncorrelated to the host halo’s spin. In this scenario, it is possible that the H I content of the underlying subhalo population is well correlated with the subhalo’s spin parameter instead. To test this idea, we plot the $M_{\text{HI}}-M_{\text{subhalo}}$ relation for the central subhaloes in Fig. 10, colouring by the spin of the central subhalo. Here, we only include galaxies $\text{type} = 0$ (centrals). We remind the reader that galaxies $\text{type} = 0$ are centrals of the central subhalo in a halo, while galaxies $\text{type} = 1$ are centrals of satellite subhaloes.

The solid line shows the median H I content of the central subhalo as a function of the subhalo mass, at $z = 0$. The dotted vertical line demarcates the micro- to medi-SURFS subhalo population transition. The central subhalo spin parameter is strongly correlated with the scatter in the $M_{\text{HI}}^{\text{subhalo}}-M_{\text{subhalo}}$ at $M_{\text{subhalo}} < 10^{11.5} M_{\odot}$, after which

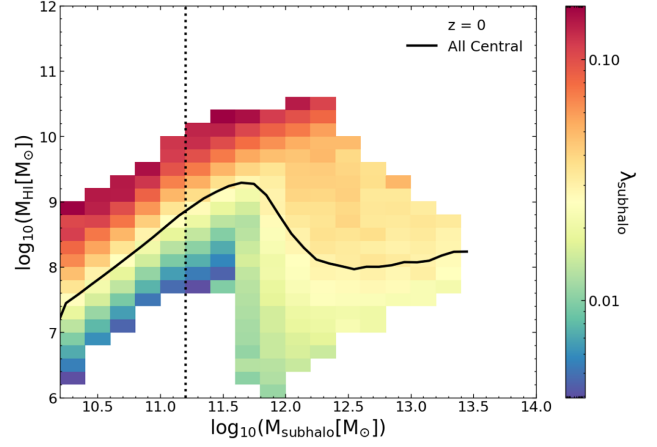


Figure 10. The H I content of central galaxies as a function of their subhalo mass at $z = 0$. Bins are coloured by the median subhalo’s spin parameter, as labelled in the colour bar. The solid black line shows the median H I mass as a function of the mass of the subhalo, M_{subhalo} . Subhaloes with higher spin parameters are H I-richer than their counterparts up to $M_{\text{subhalo}} \sim 10^{12} M_{\odot}$, after which the trend is almost completely lost at the transition and high-mass regions.

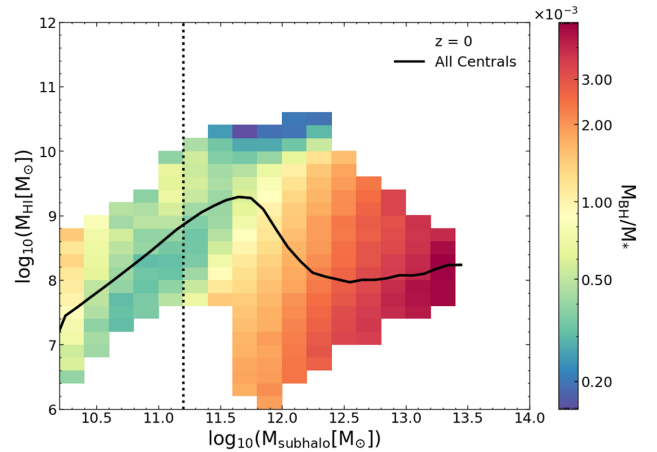


Figure 11. Similar to Fig. 10, but here the bins are coloured by the median ratio between of BH and the stellar mass of the central galaxy of central subhaloes, as labelled in the colour bar at $z = 0$. The solid line show the median H I in central subhaloes. A clear trend emerges at $M_{\text{subhalo}} \gtrsim 10^{11.4} M_{\odot}$, where we find that the subhaloes with higher BH-to-stellar mass ratio of the galaxy, lesser the H I abundance.

the correlation becomes much weaker, similar to the behaviour we obtained for the total halo mass. On the other hand, we find that satellite subhaloes³ do not show a correlation between the H I mass and the satellite subhalo’s spin at fixed subhalo mass. This shows that the weakening of the correlation between the HIHM and halo’s spin parameter is not driven by the effect of satellite galaxies, and instead central subhaloes display the same behaviour.

Fig. 11 explores the effect of AGN feedback in erasing the spin parameter dependency in the transition region at the subhalo level. We plot the $M_{\text{HI}}-M_{\text{subhalo}}$ relation explicitly for central subhaloes,

³We only use $\text{type} = 1$ satellites as they are associated with a satellite subhalo. Galaxies $\text{type} = 2$ are not included here as their host subhalo has been lost.

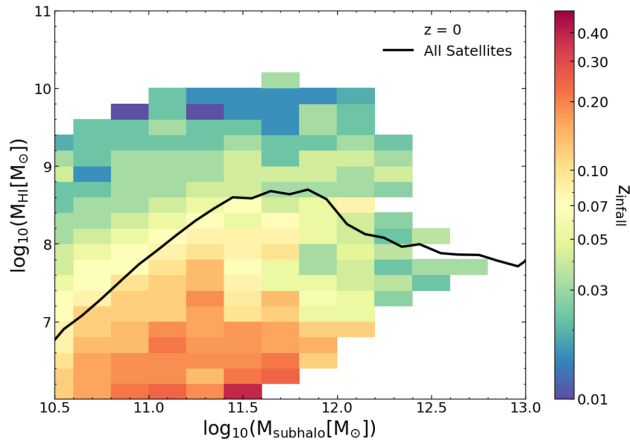


Figure 12. The HIHM relation of the satellite subhaloes ($\text{type} = 1$) in SHARK-ref at $z = 0$, with each bin coloured by the median z_{infall} of the subhalo. The solid line represents the median HI mass of all the satellite subhaloes, irrespective to their z_{infall} as a function of M_{subhalo} . A clear trend is seen between the HI of the satellite subhaloes and their z_{infall} , with later the z_{infall} , more HI-rich is the satellite subhalo.

colouring the bins by the median M_{BH}/M_* ratio, where M_{BH} and M_* are the BH and galaxy stellar masses, respectively, of the central galaxy of the central subhalo, at $z = 0$. We find that the AGNs do not show a strong correlation with the scatter of the HIHM relation for subhaloes for $M_{\text{subhalo}} \leq 10^{11.5} M_{\odot}$, but at higher subhalo masses a clear correlation emerges. This shows that the weakening of the λ_{subhalo} –HI mass correlation at fixed subhalo mass in Fig. 10 in the transition region is driven by the effect of AGN feedback. We also find a similar, albeit weaker trend in satellite subhaloes, meaning that AGN feedback is also playing a role in reducing the HI content of massive satellites $\text{type} = 1$. This is similar to what we saw for the entire haloes: the significant increase in the scatter of the HI mass–subhalo mass relation is driven by AGN feedback.

In order to understand the HI in satellite subhaloes and their lack of correlation with the subhalo’s spin parameter, we explore the correlation between the HI mass of the satellite subhalo and the redshift at which the subhalo became a satellite subhalo, z_{infall} . In Fig. 12, we plot the $M_{\text{HI}}-M_{\text{subhalo}}$ relation for satellite subhaloes, colouring by the median z_{infall} of the satellite subhaloes in each bin at $z = 0$. For this figure we limit ourselves to using medi-SURFS only, as there are not enough satellite subhaloes in micro-SURFS for a statistical study at $M_{\text{subhalo}} < 10^{11} M_{\odot}$. A clear trend emerges, where we see later infalling subhaloes being HI-rich than earlier infallers.

We remind the reader that here we are only including $\text{type} = 1$ satellites, as these quantities are not well defined for $\text{type} = 2$ satellites. This trend is expected as in SHARK we implement instantaneous stripping of the hot halo of subhaloes that become satellites, leaving the ISM to exhaust itself by continuing SF. This process of stripping plus starvation is the cause for the loss of correlation with the subhalo’s spin.

5 DEVELOPING A NUMERICAL MODEL TO POPULATE DARK MATTER HALOES WITH HI

The relation between HI and the underlying distribution of DM will be explored in significant detail over the coming years thanks to the advent of the SKA and its pathfinders. Hence, it becomes imperative that physical galaxy formation models explore the ways in which HI and DM trace each other in advance of these experiments.

Most atomic hydrogen is expected to reside in dense systems in or around galaxies, where HI is shielded from ionizing UV photons (Spinelli et al. 2020). Understanding this distribution and evolution opens up new avenues for cosmology and galaxy evolution. A significant challenge in HI cosmology applications is the requirement to produce thousands of mock observations to measure the statistical uncertainties in parameter determinations. The only plausible way of doing this is by approximate N -body, dark-matter only simulations (see Howlett et al. 2015b for an example in the optical and Howlett, Manera & Percival 2015a for an example of fast methods to produce N -body halo catalogues). Having a physical way of populating these simulations with HI is a crucial step.

As discussed previously, both the functional form and scatter of this relation can be described in terms of non-baryonic halo properties. This presents a unique advantage and the possibility to apply the phenomenological behaviour in which HI traces DM haloes we described above to large simulations. In this section, we present a numerical method to populate DM haloes with HI based on SHARK-ref. We perform exhaustive fits to the relations analysed in Section 4 in the same three halo-mass regimes presented there.

We develop our numerical model in the redshift range $0 \leq z \leq 2$, as SHARK predictions for the cosmic density of HI starts to deviate significantly from the observations at higher redshifts (see Lagos et al. 2018; Hu et al. 2019). Lagos et al. (2018) argue that the reason for this discrepancy is the fact that SHARK models only the HI content in the ISM of galaxies, while it does not explicitly model the HI content in the circumgalactic medium. Hydrodynamical simulations, e.g. van de Voort & Schaye (2012) and Diemer et al. (2019), show that at $z \gtrsim 2$ the majority of HI resides in the circumgalactic medium.

We caution the reader that the fits presented here are for one physical model of galaxy formation (SHARK-ref), though we do expect different models to behave differently (see Fig. 2). Hence, this should not be taken as a unique way of populating haloes in DM-only simulations with HI, but a way of doing it that reflects a physical model that matches a variety of observational constraints.

5.1 The total HI–halo mass scaling relation

To develop our numerical model for how to populate haloes with HI (here HI being the total HI content of a halo), we perform a fit to our simulation in two parts. We first fit the shape of the relation, $f_{M_{\text{HI}}}(M_{\text{vir}}, z)$, which depends solely on halo mass and redshift, and then a perturbation component, $\delta_{M_{\text{HI}}}$, which scales with halo properties other than mass

$$\log_{10}(M_{\text{HI}}) = f_{M_{\text{HI}}}(M_{\text{vir}}, z) + \delta_{M_{\text{HI}}}. \quad (19)$$

The median HIHM relation of SHARK is fitted with a polynomial function $f_{M_{\text{HI}}}(M_{\text{vir}}, z)$, with the fit done in bins of 0.1 dex of halo mass. We use different polynomial fits for different regions, which will be expanded upon later in this section. Our polynomial fit for the median can formally be written as

$$f_{M_{\text{HI}}}(M_{\text{vir}}, z) = \sum_{i=0}^n a_i(z) (\log_{10}(M_{\text{vir}}))^i. \quad (20)$$

The value of n differs between halo mass regions: $n = 2, 5$, and 1 , respectively, for the low-mass, transition, and high-mass regions. These were found upon iterating with different dimensions and finding the minimum n that provides a reasonable fit.

After fitting the median, we use the R HYPER-FIT package of Robotham & Obreschkow (2015) to fit a plane to the residual scatter ($\delta_{M_{\text{HI}}}$) around the HIHM relation. HYPER-FIT derives a general

likelihood function that is maximized to recover the best-fitting model describing a set of D -dimensional data points with a $(D - 1)$ -dimensional plane, with some intrinsic scatter. The secondary parameters involved in fitting the residual scatter vary according to regions. Sections 5.1.1–5.1.3 provide details of these fits for the low-mass, transition, and high-mass regions, respectively. We report the vertical scatter around the best-fitting plane provided by HYPER-FIT and use that to quantify the goodness of the fit.

5.1.1 *H I*-halo scaling relation: Low-mass region

For the *low-mass region*, we use a quadratic ($n = 2$) polynomial to fit the median H I–halo relation. A quadratic is needed to incorporate the slight downturn seen at the end of the low-mass region (around $M_{\text{vir}} \simeq 10^{11.8} M_{\odot}$). We find that the best-fitting coefficients of the median relation change with redshift. This redshift dependence can itself be fitted well with polynomials, as follows:

$$\begin{aligned} a_0^{\text{low}} &= -101.322 + 15.853 z, \\ a_1^{\text{low}} &= 17.982 - 2.757 z - 1.9808 z^2, \\ a_2^{\text{low}} &= -0.7725 + 0.2759 z, \end{aligned} \quad (21)$$

where a_{0-2}^{low} are the coefficients for the polynomial fit of equation (20) for the low-mass region, and z is redshift.

We have shown in Section 4.2.1 that for fixed M_{vir} in the low-mass region, the halo spin parameter is strongly correlated with the amount of H I contained in the halo. We therefore use that as our sole property to constrain the scatter in this region. When fitted, we find

$$\delta_{M_{\text{HI}}}^{\text{low}}(\lambda_{\text{h}}) = 1.433 (\log_{10}(\lambda_{\text{h}})) + 2.124. \quad (22)$$

Here, λ_{h} is the halo spin parameter. We get a vertical scatter of $\sigma = 0.19$ dex around our relation when we fit the residual scatter of the HIHM relation with λ_{h} using HYPER-FIT. By residual scatter we refer to the residual left after subtracting the fitted $f_{M_{\text{HI}}}(M_{\text{vir}}, z)$ from the intrinsic SHARK-ref M_{HI} values. We find that the residual scatter- λ_{h} fit for the low-mass region does not change significantly over the redshift range $0 \leq z \leq 2$ and hence, the above equation is at least valid for $z \lesssim 2$, which is the tested regime.

5.1.2 *H I*-halo scaling relation: Transition region

The fitting is the hardest at the *transition region*, as this region is dominated by AGN feedback, and the inherent scatter cannot be defined solely on halo properties. When we focus solely on halo properties, it is seen in Figs 7 and 8 that both halo spin parameter and $M_{\text{h}}^{\text{sat}}/M_{\text{vir}}$ play a role in defining the scatter of the transition region.

When fitting the median relation, $f_{M_{\text{HI}}}(M_{\text{vir}}, z)$, we use a quintic ($n = 5$) polynomial fit for our model, in order to incorporate the squiggle seen in the region from $M_{\text{vir}} = 10^{11.8} - 10^{13} M_{\odot}$. The coefficients for this fit have been tabulated in Table E1, as the parameters of the fit change with redshift in ways that are not easy to parametrize.

We note that although the halo spin parameter becomes less important in the *transition region*, haloes with a higher spin systematically retain more H I up to $M_{\text{vir}} \simeq 10^{12.5} M_{\odot}$. In this region ($M_{\text{vir}} < 10^{12.5} M_{\odot}$), the H I is still prominently contained in the central galaxies of these haloes, even though we see the beginning of the emergence of satellite population. At $M_{\text{vir}} \gtrsim 10^{12.5} M_{\odot}$, satellites become the dominant H I reservoirs of the halo and the host halo's spin parameter is not a meaningful property to define the H I content of satellite subhaloes. When we lose the spin parameter dependence,

the vertical scatter around the best-fitting plane in the transition region is captured almost entirely by $M_{\text{h}}^{\text{sat}}/M_{\text{vir}}$.

We find the HIHM relation's scatter to be reasonably well captured by

$$\begin{aligned} \delta_{M_{\text{HI}}}^{\text{TR}} \left(\lambda_{\text{h}}, \frac{M_{\text{h}}^{\text{sat}}}{M_{\text{vir}}} \right) &= b_{\text{frac}}(z) \log_{10} \left(\frac{M_{\text{h}}^{\text{sat}}}{M_{\text{vir}}} \right) \\ &+ b_{\lambda}(z) \log_{10}(\lambda_{\text{h}}) + b_0(z), \end{aligned} \quad (23)$$

with

$$\begin{aligned} b_{\text{frac}}(z) &= 0.25 e^{-z} + 0.2192, \\ b_{\lambda}(z) &= 2.77 e^{-z} + 0.7854, \\ b_0(z) &= 4.56 e^{-z} + 1.4041. \end{aligned} \quad (24)$$

We find that the scatter around the transition region changes considerably with redshift. This is due to both the AGN and subhalo populations being markedly different at earlier epochs. Despite our finding in Section 4.2.3 and Appendix B, that there is a slight correlation between the H I content of haloes and their z_{50} , we did not find z_{50} to be useful at reducing the vertical scatter compared to λ_{h} and $M_{\text{h}}^{\text{sat}}/M_{\text{vir}}$.

When we fit the residuals using HYPER-FIT, we obtain a vertical scatter of $\sigma = 0.91$ dex around the plane at $z = 0$. Although this is much larger than the 0.19 dex we achieve in the low-mass region, we find the vertical scatter decreasing to $\sigma = 0.27$ dex at $z = 2$, making it highly redshift dependent. This is due to the fact that as we move to higher redshift, the AGN influence decreases and so does the scatter dependence on it, thus making it easier to fit the relation with spin parameter and $M_{\text{h}}^{\text{sat}}/M_{\text{vir}}$ at those redshifts.

Another aspect which was discussed in Section 4.2.4 is that the scatter in this region can be better described with baryon properties; for example, using the BH-to-stellar mass ratio of the central galaxy instead of $M_{\text{h}}^{\text{sat}}/M_{\text{vir}}$ brings the HYPER-FIT vertical scatter down to $\sigma = 0.8$ dex at $z = 0$. But as the goal of this analysis is to define the HIHM scaling relation solely on the basis of halo properties, baryons are not included.

5.1.3 *H I*-halo scaling relation: High-mass region

In the *high-mass region*, the dependence of H I mass on the spin parameter or z_{50} becomes negligible. This is because haloes' H I content is almost entirely contained in satellite galaxies. Thus, in this region we see that at fixed halo mass the H I content is primarily correlated with the number of substructures present in that particular halo. In this region, we find that a linear function (i.e. polynomial fit of $n = 1$) is sufficient to describe the dependence of the median H I mass on halo mass.

The coefficients of this linear fit vary as a function of redshift as follows:

$$\begin{aligned} a_0^{\text{high}} &= -8.9448 + 8.7511 z - 5.153 z^2 + 0.891 z^3, \\ a_1^{\text{high}} &= 1.3918 - 0.4618 z + 0.1756 z^2. \end{aligned} \quad (25)$$

The scatter in the HIHM relation is then fitted as

$$\delta_{M_{\text{HI}}}^{\text{high}} \left(\frac{M_{\text{h}}^{\text{sat}}}{M_{\text{vir}}} \right) = b_{\text{frac}}(z) \log_{10} \left(\frac{M_{\text{h}}^{\text{sat}}}{M_{\text{vir}}} \right) + b_0(z), \quad (26)$$

with

$$\begin{aligned} b_{\text{frac}}(z) &= 0.498 e^{-z} + 0.11, \\ b_0(z) &= 0.669 e^{-z} + 0.1734. \end{aligned} \quad (27)$$

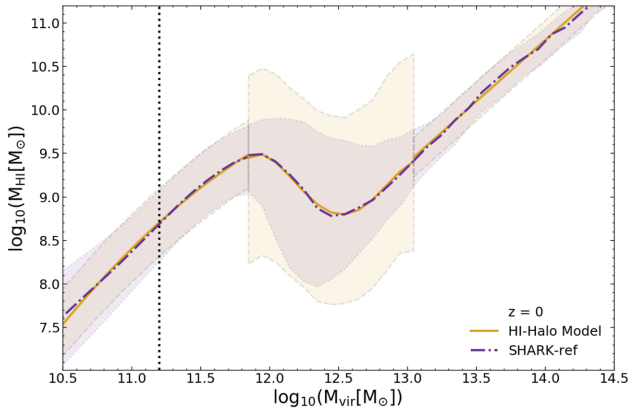


Figure 13. Overall H I content in a halo as a function of M_{vir} as predicted by our scaling relation (see Section 5), which was fitted to the output of SHARK-ref compared here. The purple (dot-dashed) and yellow (solid) line represent our the median relations for SHARK-ref and our model, respectively. The shaded region of corresponding colour around each relation shows the 16th–84th percentile range. Our scaling relation stays close to the values predicted by the SAM, but we see a slightly higher scatter in the transition region.

The vertical scatter from HYPER-FIT for this fit comes out to be $\sigma = 0.3$ dex at $z = 0$, making it a good fit for the residuals. As we move from $z = 0$ to $z = 2$, we find σ changing from 0.3 to 0.23 dex, which is a weak change and could be driven by the decreasing number of haloes in the high-mass region at higher redshifts.

5.2 Assessing the effectiveness of the numerical model for the H I–halo mass scaling relation

Fig. 13 compares the actual H I content of SHARK-ref haloes at $z = 0$ with that from our fitted HIHM scaling relation as applied to the same underlying halo population (see equations 20–26).

We can see that our numerical model produces a comparable relation to that of the intrinsic model for the low-mass and high-mass regions, highlighting the fits approximately capture the correct amount of scatter. However, for the *transition region*, the 16th and 84th percentiles of our fit are higher than in the SHARK-ref model, though when comparing with the cumulative H I in the simulation boxes (see Section C2), this might not make a huge inconsistency as the percentage of H I contributed from this region is small.

This numerical model represents significant progress over previous work, which focused only on the median H I content of haloes, without considering the scatter around the relation. This is important as H I–selected surveys will always be preferentially biased towards the more gas-rich systems rather than the typical at fixed halo mass. To properly capture this effect in mock observations it is crucial to have an understanding on how much scatter the underlying relation displays. Our numerical model offers exactly this and hence we expect it will prove useful for future H I surveys planning.

5.3 H I evolution with redshift

It has been shown in previous sections (see Section 5.1) that the coefficients of f_{MHI} are redshift-dependent. We find that as we move to higher redshifts, AGN feedback is less efficient at preventing the halo gas from cooling. Also, haloes at higher redshifts have not had enough time to assemble all their mass, leading to a lesser number of substructures. Both of these factors significantly contribute to

the evolution in the scatter at the *transition region*. By $z = 2$, the transition is barely visible. We explore this in detail in Appendix C1.

We have seen in Fig. 13 that the transition region is characterized by a large scatter that is difficult to fully account for with halo properties alone. It is therefore informative to know how much H I in SHARK-ref resides in the transition region, to quantify the impact inaccurate estimates of the scatter can have on studies that focus on unveiling the total H I content of the Universe. We find at $z = 0$, about 25 per cent of the total H I resides in the transition region and 60 per cent in the low-mass region. At $z = 2$, we find that almost 80 per cent of the H I in the SHARK-ref resides in the low-mass region. We explore the evolution of the cumulative H I in SHARK-ref in Appendix C2.

6 CONCLUSIONS

Understanding the evolution of H I throughout cosmic time provides key insights into cosmology and galaxy evolution. Unlike the stellar–halo mass relation, the HIHM relation is not necessarily monotonic and is likely to be characterized by a large scatter (given the large scatter in the H I–stellar mass relation; Catinella et al. 2018). In this paper, we have used the state-of-the-art semi-analytic galaxy formation model SHARK, with the aim of understanding the physical processes behind the shape, scatter and evolution of the H I–halo mass relation at $0 \leq z \leq 2$.

We compared the H I–halo mass relation and the H I clustering of SHARK with available observations. These observations were not used as part of the tuning of the free parameters of SHARK, and can hence be considered predictions. We find the predicted H I clustering in SHARK to be in excellent agreement with the observations. However, when comparing with observational inferences of the H I–halo mass relation, coming mostly from H I stacking of groups, we found that SHARK reproduces well the H I abundance in haloes of masses $< 10^{12}$ and $> 10^{13.3} M_{\odot}$, but in the range 10^{12} – $10^{13.3} M_{\odot}$, SHARK underpredicts the abundance of H I in haloes. In an upcoming paper (Chauhan et al. in preparation), we show that these discrepancies are largely due to the uncertainty in group definition around that halo mass (that in current spectroscopic surveys have a small occupancy).

We then explored the effect of different physical processes in the shape of the HIHM relation, and what properties of haloes are the best secondary parameter that correlates with the scatter in the HIHM relation. Our key results can be summarized as follows:

(i) The HIHM relation is characterized by three mass regions that display distinct behaviours. At $z = 0$, we find that the total H I content of haloes with $M_{\text{vir}} < 10^{11.8} M_{\odot}$, aka low-mass region, increases monotonically with the halo mass. In haloes of masses $10^{11.8} M_{\odot} < M_{\text{vir}} < 10^{13} M_{\odot}$, aka the transition zone, the total H I content of haloes peaks at $M_{\text{vir}} = 10^{12} M_{\odot}$ and then declines with increasing halo mass. For haloes of masses $M_{\text{vir}} > 10^{13} M_{\odot}$, aka the high-mass region, the total H I content of haloes starts to increase again with increasing halo mass. The scatter around the HIHM varies significantly in the three mass regions, being ~ 0.5 , ~ 1.2 , and ~ 0.4 dex in the low-mass, transition, and high-mass regions, respectively.

(ii) We find the contribution to the total H I mass of the halo to be dominated by central galaxies for haloes of $M_{\text{vir}} < 10^{12.5} M_{\odot}$. At higher halo masses, satellite galaxies are the dominant contributor. The bump seen in the HIHM relation in the transition zone is caused by central galaxies, while the total H I mass contributed by satellites scales monotonically with halo mass. The latter is what produces the increasing H I mass with increasing halo mass in the high-mass zone.

(iii) The peak of the HIHM relation in the transition region and the halo mass at which this peaks happens are largely determined by the AGN feedback efficiency, with stellar feedback, photoionization feedback, and ISM modelling playing a lesser role. The dip in the HIHM relation is caused by the suppression of gas cooling in these haloes due to the influence of AGN feedback. At lower halo masses, AGN does not play an important role.

(iv) We isolate the main secondary parameter responsible for the scatter of the HIHM relation. In the low-mass region, the scatter at fixed halo mass is highly correlated with the spin parameter of the halo, whereas for the high-mass zone, the scatter is correlated with the fractional contribution from substructure to the total halo mass, $M_h^{\text{sat}}/M_{\text{vir}}$. As for the transition zone, we find the scatter to be highly dependent on the black hole-to-stellar mass ratio of the central galaxy, reflecting the importance of AGN feedback in this region. However, when we explored halo properties only, we find that a combination of halo's spin and $M_h^{\text{sat}}/M_{\text{vir}}$ is relatively successful at characterizing the scatter of the HIHM relation in the transition zone. Once these secondary dependencies are included, the vertical scatter of the two-dimensional plane (between the median-subtracted halo mass HI mass and the secondary parameter) at $z = 0$ is significantly tighter than the HIHM relation, with values of ~ 0.19 , ≈ 0.91 , and 0.3 dex, in the low-mass, transition and high-mass zones, respectively.

(v) As we move to higher redshifts, the transition zone starts to shrink, as AGN feedback becomes less efficient. The vertical scatter in the three-dimensional plane over the transition zone decreases significantly with redshift, from $\sigma = 0.91$ dex at $z = 0$ to $\sigma = 0.27$ dex at $z = 2$. The latter values for the scatter already consider the dependency on spin and $M_h^{\text{sat}}/M_{\text{vir}}$. In the low- and high-mass regions, the decrease in the scatter is not as significant as in the transition zone, with the low-mass region hardly seeing a decrease in the vertical scatter (remaining at ~ 0.19 dex once the halo spin is considered) and the high-mass region sees a decrease from ~ 0.3 dex at $z = 0$ to ~ 0.23 dex at $z = 2$, once $M_h^{\text{sat}}/M_{\text{vir}}$ is considered. By $z = 2$, the HIHM relation is monotonic over the whole halo mass range.

Finally, we use the lessons learned to develop a numerical model to populate haloes in DM-only simulations with HI, depending on their halo mass, spin parameter, $M_h^{\text{sat}}/M_{\text{vir}}$, and redshift. Obvious applications of this numerical model include HI intensity mapping, HI stacking, and modelling of HI clustering. This study also opens up avenues for exploring the role of different halo properties in the HIHM relation. With the upcoming SKA and its Pathfinders, we will be able to explore the role of halo properties in the HIHM relation observationally, providing better constraints and deeper insight into the HIHM relation.

ACKNOWLEDGEMENTS

We would like to thank Marta Spinelli, Hong Guo, Jian Fu, Kristine Spekkens, Cullan Howlett, Matías Bravo, Stéphane Courteau, and Rob Crain for their constructive comments and useful discussions. We also thank Aaron Robotham, Rodrigo Tobar, and Pascal Elahi for their contribution towards SURFS and SHARK, and Mark Boulton for his IT help. GC is funded by the Mobilising European Research in Astrophysics and Cosmology (MERAC) Foundation, through the Postdoctoral Research Award of CL, and the University of Western Australia. Parts of this research were carried out by the ARC Centre of Excellence for All Sky Astrophysics in 3 Dimensions (ASTRO 3D), through project number CE170100013. CL is funded by ASTRO 3D. ARHS acknowledges receipt of the Jim Buckee Fellowship at ICRAR-UWA. This work was supported by resources provided by

the Pawsey Supercomputing Centre with funding from the Australian Government and the Government of Western Australia.

DATA AVAILABILITY

The data that support the findings of this study are available upon request from the corresponding author, GC. The SURFS simulations used in this work can be freely accessed from <https://tinyurl.com/y4pvra87> (micro-SURFS) and <https://tinyurl.com/y6ql46d4> (medi-SURFS).

REFERENCES

- Abdalla F. B., Rawlings S., 2005, *MNRAS*, 360, 27
Amarantidis S. et al., 2019, *MNRAS*, 485, 2694
Asplund M., Grevesse N., Sauval A. J., Scott P., 2009, *ARA&A*, 47, 481
Baugh C. M., 2006, *Rep. Prog. Phys.*, 69, 3101
Baugh C. M. et al., 2019, *MNRAS*, 483, 4922
Behroozi P. S., Conroy C., Wechsler R. H., 2010, *ApJ*, 717, 379
Blanton M. R., Moustakas J., 2009, *ARA&A*, 47, 159
Blitz L., Rosolowsky E., 2006, *ApJ*, 650, 933 (BR06)
Bonatto C., Bica E., 2011, *MNRAS*, 415, 2827
Bondi H., 1952, *MNRAS*, 112, 195
Bower R. G., Benson A. J., Malbon R., Helly J. C., Frenk C. S., Baugh C. M., Cole S., Lacey C. G., 2006, *MNRAS*, 370, 645
Boylan-Kolchin M., Springel V., White S. D. M., Jenkins A., Lemson G., 2009, *MNRAS*, 398, 1150
Bravo M., Lagos C. D. P., Robotham A. S. G., Bellstedt S., Obreschkow D., 2020, *MNRAS*, 497, 3026
Brown T., Catinella B., Cortese L., Kilborn V., Haynes M. P., Giovanelli R., 2015, *MNRAS*, 452, 2479
Brown T. et al., 2017, *MNRAS*, 466, 1275
Cañas R., Elahi P. J., Welker C., del P Lagos C., Power C., Dubois Y., Pichon C., 2019, *MNRAS*, 482, 2039
Catinella B. et al., 2010, *MNRAS*, 403, 683
Catinella B. et al., 2018, *MNRAS*, 476, 875
Chauhan G., Lagos C. D. P., Obreschkow D., Power C., Oman K., Elahi P. J., 2019, *MNRAS*, 488, 5898
Cole S., Lacey C., 1996, *MNRAS*, 281, 716
Cole S., Lacey C., Baugh C., Frenk C., 2000, *MNRAS*, 319, 168
Crain R. A. et al., 2015, *MNRAS*, 450, 1937
Crain R. A. et al., 2017, *MNRAS*, 464, 4204
Croton D. J. et al., 2006, *MNRAS*, 365, 11
Croton D. J., Gao L., White S. D. M., 2007, *MNRAS*, 374, 1303
Croton D. J. et al., 2016, *ApJS*, 222, 22 (Croton16)
Davies L. J. M. et al., 2018, *MNRAS*, 480, 768
de Blok W. J. G., McGaugh S. S., van der Hulst J. M., 1996, *MNRAS*, 283, 18
Diemer B. et al., 2019, *MNRAS*, 487, 1529
Driver S. P. et al., 2018, *MNRAS*, 475, 2891
Duffy A. R., Schaye J., Kay S. T., Dalla Vecchia C., 2008, *MNRAS*, 390, L64
Duffy A. R., Meyer M. J., Staveley-Smith L., Bernyk M., Croton D. J., Koribalski B. S., Gerstmann D., Westerlund S., 2012, *MNRAS*, 426, 3385
Eckert K. D., Kannappan S. J., Stark D. V., Moffett A. J., Norris M. A., Snyder E. M., Hoversten E. A., 2015, *ApJ*, 810, 166
Eckert K. D. et al., 2017, *ApJ*, 849, 20
Elahi P. J., Welker C., Power C., Lagos C. D. P., Robotham A. S. G., Cañas R., Poulton R., 2018, *MNRAS*, 475, 5338
Elahi P. J., Cañas R., Poulton R. J. J., Tobar R. J., Willis J. S., Lagos C. D. P., Power C., Robotham A. S. G., 2019a, *PASA*, 36, e021
Elahi P. J., Poulton R. J. J., Tobar R. J., Cañas R., Lagos C. D. P., Power C., Robotham A. S. G., 2019b, *PASA*, 36, e028
Elmegreen B. G., 1989, *ApJ*, 338, 178
Fabello S., Catinella B., Giovanelli R., Kauffmann G., Haynes M. P., Heckman T. M., Schiminovich D., 2011, *MNRAS*, 411, 993

- Fabello S., Kauffmann G., Catinella B., Li C., Giovanelli R., Haynes M. P., 2012, *MNRAS*, 427, 2841
- Fukugita M., Hogan C. J., Peebles P. J. E., 1998, *ApJ*, 503, 518
- Gao L., Springel V., White S. D. M., 2005, *MNRAS*, 363, L66
- Genzel R. et al., 2010, *MNRAS*, 407, 2091
- Giovanelli R. et al., 2005, *AJ*, 130, 2598
- Gnedin N. Y., 2012, *ApJ*, 754, 113
- Gnedin N. Y., Draine B. T., 2014, *ApJ*, 795, 37 (GD14)
- Guo H., Li C., Zheng Z., Mo H. J., Jing Y. P., Zu Y., Lim S. H., Xu H., 2017, *ApJ*, 846, 61
- Guo H., Jones M. G., Haynes M. P., Fu J., 2020, *ApJ*, 894, 92
- Haynes M. P. et al., 2018, *ApJ*, 861, 49
- Holwerda B. W., Blyth S., Baker A. J., MeerKAT Deep HI Survey Team, 2011, in American Astronomical Society Meeting Abstracts #217. p. 433.17
- Howlett C., Manera M., Percival W. J., 2015a, *Astron. Comput.*, 12, 109
- Howlett C., Ross A. J., Samushia L., Percival W. J., Manera M., 2015b, *MNRAS*, 449, 848
- Hu W. et al., 2019, *MNRAS*, 489, 1619
- Kim H.-S., Wyithe J. S. B., Power C., Park J., Lagos C. D. P., Baugh C. M., 2015, *MNRAS*, 453, 2316
- Kim H.-S., Wyithe J. S. B., Baugh C. M., Lagos C. D. P., Power C., Park J., 2017, *MNRAS*, 465, 111
- Koribalski B. S. et al., 2020, *ApSS*, 365, 118
- Kregel M., Van Der Kruit P. C., Grijs R. D., 2002, *MNRAS*, 334, 646
- Krumholz M. R., Dekel A., 2012, *ApJ*, 753, 16
- Kulier A., Padilla N., Schaye J., Crain R. A., Schaller M., Bower R. G., Theuns T., Paillas E., 2019, *MNRAS*, 482, 3261
- Lagos C. D. P., Padilla N. D., Cora S. A., 2009, *MNRAS*, 395, 625
- Lagos C. D. P., Lacey C. G., Baugh C. M., 2013, *MNRAS*, 436, 1787
- Lagos C. D. P., Davis T. A., Lacey C. G., Zwaan M. A., Baugh C. M., Gonzalez-Perez V., Padilla N. D., 2014, *MNRAS*, 443, 1002
- Lagos C. D. P., Tobar R. J., Robotham A. S. G., Obreschkow D., Mitchell P. D., Power C., Elahi P. J., 2018, *MNRAS*, 481, 3573
- Lagos C. D. P. et al., 2019, *MNRAS*, 489, 4196
- Lagos C. D. P., da Cunha E., Robotham A. S. G., Obreschkow D., Valentino F., Fujimoto S., Magdis G. E., Tobar R., 2020, preprint (arXiv:2007.09853)
- Leroy A. K., Walter F., Brinks E., Bigiel F., de Blok W. J. G., Madore B., Thornley M. D., 2008, *AJ*, 136, 2782
- Lim S. H., Mo H. J., Lu Y., Wang H., Yang X., 2017, *MNRAS*, 470, 2982
- Lutz K. A. et al., 2018, *MNRAS*, 476, 3744
- Maddox N., Hess K. M., Obreschkow D., Jarvis M. J., Blyth S.-L., 2015, *MNRAS*, 447, 1610
- Martin A. M., Giovanelli R., Haynes M. P., Guzzo L., 2012, *ApJ*, 750, 38
- Matthee J., Schaye J., Crain R. A., Schaller M., Bower R., Theuns T., 2017, *MNRAS*, 465, 2381
- Meyer M. J. et al., 2004, *MNRAS*, 350, 1195
- Meyer M. J., Zwaan M. A., Webster R. L., Brown M. J. I., Staveley-Smith L., 2007, *ApJ*, 654, 702
- Mitchell P. D., Lacey C. G., Baugh C. M., Cole S., 2016, *MNRAS*, 456, 1459
- Moster B. P., Somerville R. S., Maulbetsch C., van den Bosch F. C., Macciò A. V., Naab T., Oser L., 2010, *ApJ*, 710, 903
- Muratov A. L., Kereš D., Faucher-Giguère C.-A., Hopkins P. F., Quataert E., Murray N., 2015, *MNRAS*, 454, 2691
- Navarro J. F., Frenk C. S., White S. D. M., 1997, *ApJ*, 490, 493
- Nelson D. et al., 2018, *MNRAS*, 475, 624
- Nelson D. et al., 2019, *MNRAS*, 490, 3234
- Nulsen P. E. J., Fabian A. C., 2000, *MNRAS*, 311, 346
- Obreschkow D., Glazebrook K., Kilborn V., Lutz K., 2016, *ApJ*, 824, L26
- Obuljen A., Alonso D., Villaescusa-Navarro F., Yoon I., Jones M., 2019, *MNRAS*, 486, 5124
- Padmanabhan H., Kulkarni G., 2017, *MNRAS*, 470, 340
- Padmanabhan H., Refregier A., 2017, *MNRAS*, 464, 4008
- Papastergis E., Giovanelli R., Haynes M. P., Rodríguez-Puebla A., Jones M. G., 2013, *ApJ*, 776, 43
- Peebles P. J. E., 1969, *ApJ*, 155, 393
- Pillepich A. et al., 2018, *MNRAS*, 475, 648
- Planck Collaboration XIII, 2016, *A&A*, 594, A13
- Poulton R. J. J., Robotham A. S. G., Power C., Elahi P. J., 2018, *PASA*, 35
- Power C., Knebe A., 2006, *MNRAS*, 370, 691
- Power C. et al., 2015, in Proc. Sci., Galaxy Formation & Dark Matter Modelling in the Era of the Square Kilometre Arra. SISSA, Trieste, PoS#133
- Pritchard J. R., Loeb A., 2012, *Rep. Prog. Phys.*, 75, 086901
- Prochaska J. X., Wolfe A. M., 2009, *ApJ*, 696, 1543
- Putman M., Peek J., Joung M., 2012, *ARA&A*, 50, 491
- Rhee J., Lah P., Briggs F. H., Chengalur J. N., Colless M., Willner S. P., Ashby M. L. N., Le Fèvre O., 2018, *MNRAS*, 473, 1879
- Robotham A. S. G., Obreschkow D., 2015, *PASA*, 32, e033
- Robotham A. S. G. et al., 2011, *MNRAS*, 416, 2640
- Sargent M. T. et al., 2014, *ApJ*, 793, 19
- Schaye J. et al., 2015, *MNRAS*, 446, 521
- Schombert J. M., McGaugh S. S., Eder J. A., 2001, *AJ*, 121, 2420
- Sheth R. K., Tormen G., 2002, *MNRAS*, 329, 61
- Sinha M., Garrison L. H., 2020, *MNRAS*, 491, 3022
- Sobacchi E., Mesinger A., 2013, *MNRAS*, 432, L51
- Somerville R. S., Davé R., 2015, *ARA&A*, 53, 51
- Spinelli M., Zoldan A., Lucia G. D., Xie L., Viel M., 2020, *MNRAS*, 493, 5434
- Springel V. et al., 2005, *Nature*, 435, 629
- Stevens A. R. H., Brown T., 2017, *MNRAS*, 471, 447
- Stevens A. R. H. et al., 2019, *MNRAS*, 483, 5334
- van de Voort F., Schaye J., 2012, *MNRAS*, 423, 2991
- van de Voort F., Schaye J., Booth C. M., Haas M. R., Dalla Vecchia C., 2011, *MNRAS*, 414, 2458
- Villaescusa-Navarro F. et al., 2018, *ApJ*, 866, 135
- Walter F., Brinks E., de Blok W. J. G., Bigiel F., Kennicutt R. C., Thornley M. D., Leroy A., 2008, *AJ*, 136, 2563
- Wechsler R. H., Tinker J. L., 2018, *ARA&A*, 56, 435
- Wolfire M. G., McKee C. F., Hollenbach D., Tielens A. G. G. M., 2003, *ApJ*, 587, 278
- Wright R. J., Lagos C. D. P., Power C., Mitchell P. D., 2020, *MNRAS*, preprint (arXiv:2006.00924)
- Xie L., De Lucia G., Hirschmann M., Fontanot F., Zoldan A., 2017, *MNRAS*, 469, 968
- Xie L., De Lucia G., Hirschmann M., Fontanot F., 2020, *MNRAS*, preprint (arXiv:2003.12757)
- Zhang W., Li C., Kauffmann G., Xiao T., 2013, *MNRAS*, 429, 2191

APPENDIX A: UNDERSTANDING THE SHAPE – CONT.

In this section, we explore a bit more on the impact using different models and parameters for a physical process have on the shape of the HIHM relation, which we have only briefly discussed in the main paper.

A1 Photoionization effect

We discuss the implementation of photoionization feedback in SHARK in Section 2.3.4, and briefly touched on the effect changing its parameters has on the HIHM relation in Section 4.1.3. Here, we show the effect of photoionization feedback on the overall HI content of the haloes at $z = 0$. We vary the value of v_{cut} , which from equation (16) directly affects the circular velocity (v_{thresh}) of haloes under which the halo gas is not allowed to cool down.

The effect of varying v_{thresh} is presented in Fig. A1, where haloes below a certain mass (which correspond to the circular velocity v_{thresh}) do not have HI in them, as the halo gas is kept ionized. As expected, increasing v_{cut} has the effect of shifting the steep decline of the HI fraction–halo mass relation to higher halo masses. Though, changing v_{cut} value does not have any effect on the *transition region* – the drop essentially remains at the same M_{vir} value ($10^{12} M_{\odot}$) for all the variations. We find that photoionization feedback becomes

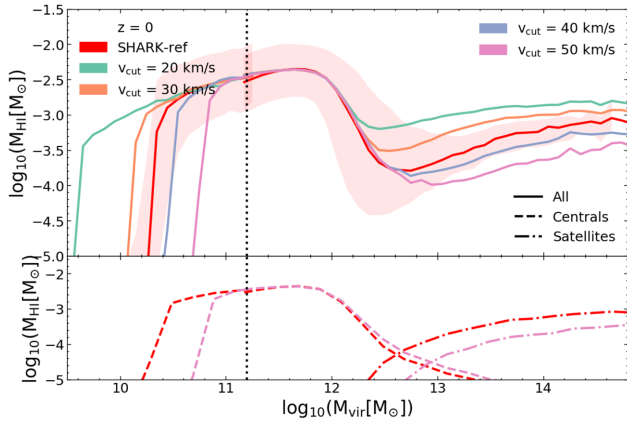


Figure A1. As in Fig. 5 but for different values of v_{cut} , which represents the virial velocity threshold under which the gas in haloes is assumed to be kept ionized by the UV background, and is hence not allowed to cool down and replenish the interstellar medium of the central galaxy (see equation 16). Different colour lines represent different v_{cut} values, with red representing the default SHARK-ref model and the shaded region being the 16th–84th percentile range. Photoionization heating does not affect the knee of the HIHM relation, though it does affect the amount of H I contained in haloes in the low- and high-mass regions. *Lower panel:* The median H I contribution from centrals and satellites to the total H I of the halo. For clarity, we only show SHARK-ref and the $v_{\text{cut}} = 50 \text{ km s}^{-1}$ variation. Unlike previous figures, changing the photoionization feedback leads to a change in the H I content of satellites, with higher the feedback the lesser the amount of H I in satellites.

more prominent for the H I content of haloes after $M_{\text{vir}} > 10^{12.4} M_{\odot}$, with a smaller v_{cut} driving a higher H I content in haloes. This effect is due to smaller haloes being allowed to cool down their halo gas under smaller v_{cut} values, increasing their H I content. These centrals of low-mass haloes can then become satellites of larger haloes and contribute to the total H I content of that halo.

In the lower panel of Fig. A1, we compare the H I contributions from satellites and centrals for the SHARK-ref and $v_{\text{cut}} = 50 \text{ km s}^{-1}$ runs. We find that the central H I contribution remains almost unchanged in both the runs, except for the halo mass below which the H I content sharply decreases, which is at $M_{\text{vir}} \approx 10^{10.4}$ and $M_{\text{vir}} \approx 10^{11} M_{\odot}$ for SHARK-ref and $v_{\text{cut}} = 50 \text{ km s}^{-1}$ runs, respectively. On the contrary, the contribution from satellite galaxies is different in these runs, with SHARK-ref having higher H I content in satellites than the other extreme run. The latter is due to the galaxies that become satellites being more H I-rich with smaller v_{cut} values.

A2 Interstellar medium model effect

In SHARK, stars form from molecular gas, and different models are implemented for how to split the ISM into ionized, atomic and molecular gas phases. Here, we compare two models for the molecular-to-atomic gas partition, specifically the BR06 (the default model of choice) and the GD14 model. In both cases, stars are formed from the molecular gas with a fixed efficiency (see equation 7). A brief overview of the effect of changing the ISM model had been given in Section 4.1.3, here, we delve into more details.

Fig. A2 shows $M_{\text{HI}}/M_{\text{vir}}$ as a function of M_{vir} for different H₂-to-H I partition models that are implemented in SHARK, with the top panel showing the total $M_{\text{HI}}/M_{\text{vir}}$ ratio, and the bottom panel showing the centrals and satellite contributions at $z = 0$. When using the GD14 prescription, the overall H I content of haloes is higher than when adopting the BR06 prescription, except at halo masses between

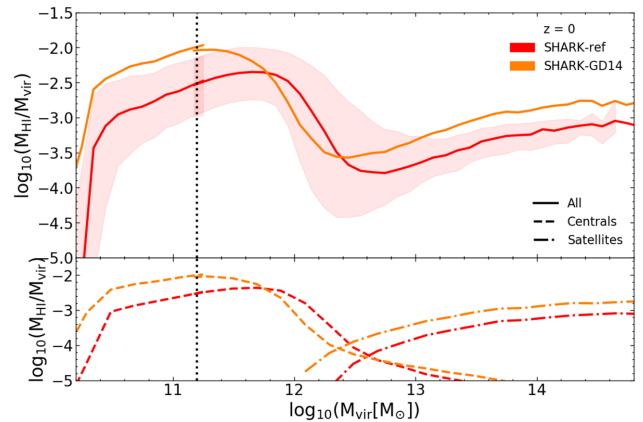


Figure A2. As in Fig. 5 but for two variations of the molecular-to-atomic interstellar gas partition in SHARK. The models being compared are the default SHARK model as shown in Lagos et al. (2018), which incorporates the Blitz & Rosolowsky (2006) prescription (SHARK-ref) to split atomic and molecular gas in the interstellar medium of galaxies, with a variant adopting the Gnedin & Draine (2014) atomic-to-molecular transition prescription (SHARK-GD14). In both variants, stars from the molecular gas with the same efficiency. The top panel shows the entire H I fraction whereas the bottom panel shows the central and satellite contributions.

10^{12} and $10^{12.7} M_{\odot}$. The *transition region* for the model adopting the GD14 prescription is at a lower halo masses, $M_{\text{vir}} \approx 10^{11.5} M_{\odot}$ against $M_{\text{vir}} \approx 10^{12} M_{\odot}$ for BR06. In this transition region, BR06 predicts a slightly higher abundance of H I. However, at lower and higher halo masses, GD14 results in higher H I content. The fact that centrals of low-mass haloes are more H I-rich in GD14 than BR06 is the cause for the higher H I abundance at high halo masses, as many of the low-mass centrals become satellites as time progresses.

When we compare the H I contributions of centrals and satellites to the overall H I of the halo (bottom panel in Fig. A2), we find that the H I contribution of centrals in GD14 is higher than BR06 for haloes $M_{\text{vir}} < 10^{12} M_{\odot}$, while at higher masses there is virtually no difference. This happens due to the fact that H I-H₂ partition in GD14 depends on the gas metallicity (among other parameters). This is not the case for BR06, which is a purely pressure-based model. SHARK-ref predicts low metallicities for low-mass galaxies, which in turn makes the H I value for low-mass haloes to be higher in GD14, as the H I in low-mass haloes is dominated by the centrals. As for the satellite contribution, we see that GD14 consistently predicts more H I than BR06 throughout all virial masses, again due to the gas metallicity effect. The fact that the transition region happens at lower halo masses in GD14 than in SHARK-ref is, however, unrelated to the SF law. We find that BH masses are slightly bigger at intermediate mass galaxies (around the break of the stellar mass function) in GD14, causing AGN feedback to be more efficient than in SHARK-ref in those galaxies. As seen before, more efficient AGN feedback shifts the transition region to lower halo masses, which is effectively what happens in the GD14 run. This again highlights the complex interplay between the different baryon physics in models such as SHARK.

A3 Gas stripping effect

The last effect we want to test is the environmental effect, which we do by comparing the effect turning ‘off’ ram-pressure stripping has on the overall H I content of the haloes (see Section 2.3.5).

In Fig. A3 (top panel), we compare $M_{\text{HI}}/M_{\text{vir}}-M_{\text{vir}}$ relation for stripping mode ‘on’ and ‘off’ as a function of M_{vir} . We find that the

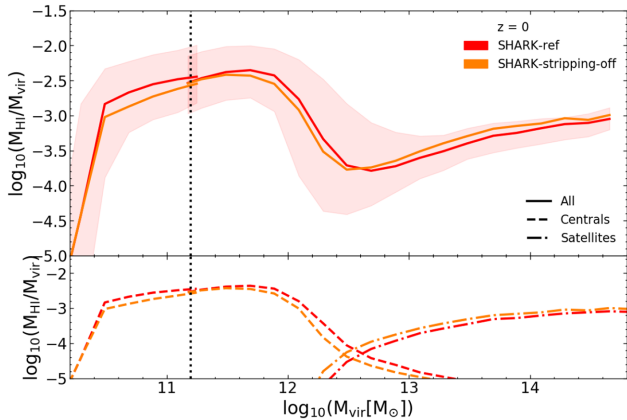


Figure A3. Similar to Fig. A2, but for the default model (red) versus a model with no gas stripping (yellow). The top panel shows the total H I fraction, whereas the bottom panel shows the central and satellite contributions. Differences between the two models are clear when we decompose the H I contribution between centrals and satellites, but these differences compensate each other so that the total H I in haloes is barely affected. The top panel shows the entire H I fraction whereas the bottom panel shows the central and satellite contributions.

total amount of H I in either model is approximately the same, though stripping ‘off’ tends to lead to a slightly lower H I in the transition region and higher H I in the high-mass region.

When looking at the central–satellite galaxies contribution to the total H I mass of the halo (bottom panel), we find centrals to reduce their H I content when stripping is off, while satellites become more important. This happens because when stripping is off, satellites are able to hold on their hot haloes for longer, which means that the hot halo of the central is now less massive than in the run with stripping. This leads to central galaxies accreting less gas (due to the smaller overall reservoir of gas), while satellite can continue to accrete gas for longer. Clearly these two competing effects compensate relatively well as to lead to small differences in the total H I content of haloes at $10^{12} M_{\odot} < M_{\text{vir}} < 10^{13.5} M_{\odot}$.

APPENDIX B: FORMATION AGE EFFECT

Here, we discuss the effect halo formation age has on the scatter in the HIHM relation. This was briefly discussed in Section 4.2.3.

We define formation age (z_{50}) as the redshift at which the halo accreted 50 percent of its present mass. It has been speculated that z_{50} is correlated to the amount of H I contained in a halo. Guo et al. (2017) found from their clustering measurements of ALFALFA galaxies that a way of describing the clustering bias dependence on scale was to assume H I-rich galaxies to live in preferentially young haloes. Under this assumption, they developed a subhalo abundance matching model (SHAM) which was used to derive a strong correlation between the H I content of the haloes and its z_{50} . A suitable explanation for this effect is the fact that young haloes would be expected to contain H I-rich galaxies, as they had not had enough time to lose their cold gas via ‘ram-pressure stripping’ or other environmental effect. Spinelli et al. (2020), using the semi-analytic model of galaxy formation GAEA, found that in low-mass haloes there was no difference between young and old haloes in terms of their H I content; but as the M_{vir} increased, a segregation appeared between young and old haloes, with the former being more H I-rich in agreement with Guo et al. (2017) inferences.

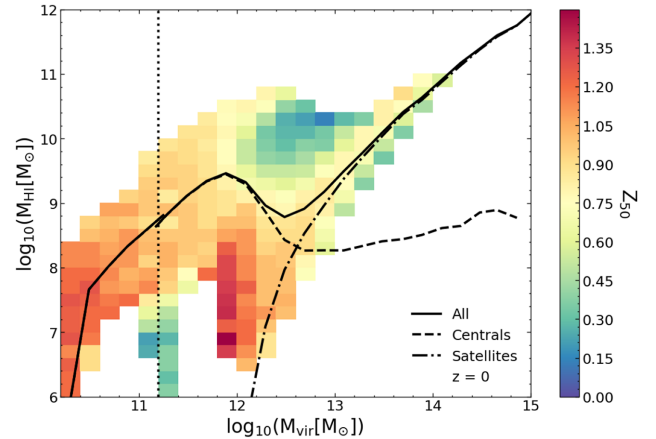


Figure B1. The HIHM relation of haloes in SHARK-ref at $z = 0$, with each bin being coloured by the median z_{50} of the haloes in that bin, as labelled in the colour bar. The solid line represents the median H I mass of the halo as a function of M_{vir} , while the dashed and dash-dotted lines represent the central and satellite galaxies contributions, respectively. The vertical dotted line shows the transition from micro-SURFS to medi-SURFS at lower and higher halo masses, respectively. A slight trend with z_{50} is seen at the transition region so that younger haloes tend to be more H I-rich. This trend reverses though at higher halo masses.

We test the effect of z_{50} here. Fig. B1 shows the HIHM relation at $z = 0$ colouring each bin by the median formation age of haloes in that bin. We find that in SHARK, z_{50} does not show a significant trend in the low-mass region, though a slight trend is noticeable in the transition region. We see that younger haloes (closer to $z = 0$) tend to have more H I than their counterparts of the same mass. We think the trend emerges here because it is in this region that satellites start to become a more prominent reservoir of H I compared to the central galaxy. We see a slight opposite trend in the high-mass region, where later forming haloes tend to be H I poorer which contradicts the conclusion in Guo et al. (2017). This is due to older haloes having on average more substructure and therefore more satellite galaxies at fixed halo mass (see Croton, Gao & White 2007; Wechsler & Tinker 2018), which contribute to the total H I content of the halo. We discuss this in more detail in Section 4.2.2. Fig. B2 shows the relation between the halo mass, z_{50} and the number of substructure

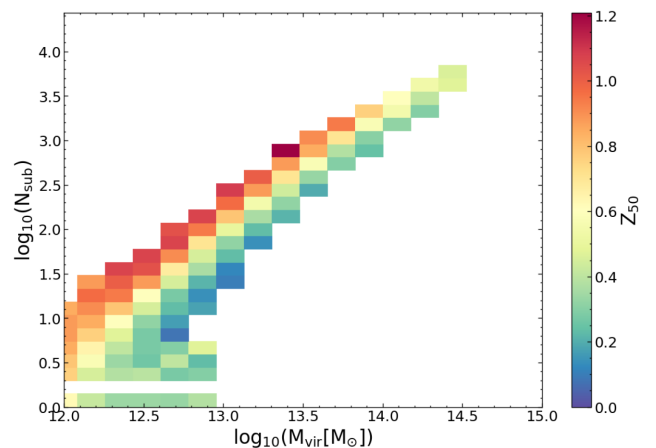


Figure B2. The number of subhaloes in a halo as a function of M_{vir} , with bins coloured according to the median z_{50} of that bin. Older haloes have more substructure than their younger counterparts at fixed halo mass.

per halo. We find that haloes formed earlier have more substructure as compared to their younger counterparts in the same mass bin.

Appendix D presents the redshift evolution of the HI halo mass– z_{50} relation up to $z = 2$. We find that the trend we see at $z = 0$ holds at high redshift with the main difference being the expected lack of massive haloes.

APPENDIX C: DEVELOPING NUMERICAL MODEL TO POPULATE DARK MATTER HALOES WITH HI – CONT.

After the brief overview given in Section 5.3, here we explore a bit more on the evolution of the HIHM relation through different redshifts.

C1 Redshift dependence

As noted in Section 5.1, the coefficients for $f_{M_{\text{HI}}}(M_{\text{vir}})$ are dependent on redshift. We find that as we move towards higher redshifts, the *transition region* shrinks, with the noticeable bump (around $M_{\text{vir}} \simeq 10^{12} M_{\odot}$) becoming flatter (see Appendix D). By the time we reach $z = 2$, the HI–halo scaling relation becomes a monotonically increasing function of M_{vir} . One of the key reasons behind this outcome is that for higher redshifts AGN feedback is less efficient than at $z = 0$ and therefore by $z = 2$ AGN feedback does not play a significant role at keeping the halo gas hot and preventing gas cooling and accretion on to galaxies. In addition, as the haloes have not had enough time to assemble all of their mass, they do not have enough substructures yet to contribute to increasing the scatter in the transition region. In short, we find that there are no distinctive regions at high redshifts, i.e. the *transition region* effectively disappears.

In the low-mass region we find that, while the shape of the median HIHM relation carries a redshift dependence, the fits to the residuals ($\delta_{M_{\text{HI}}}$, which captures the scatter) do not. That is to say, for example, the influence that halo spin has on the total HI in a halo of fixed virial mass is the same at all epochs.

In Fig. C1, we compare the HI mass calculated by our model with the intrinsic HI output from SHARK-ref, at each snapshot out to $z = 2$, to assess the performance of our numerical model. We show this for the individual halo mass regions as well as the total

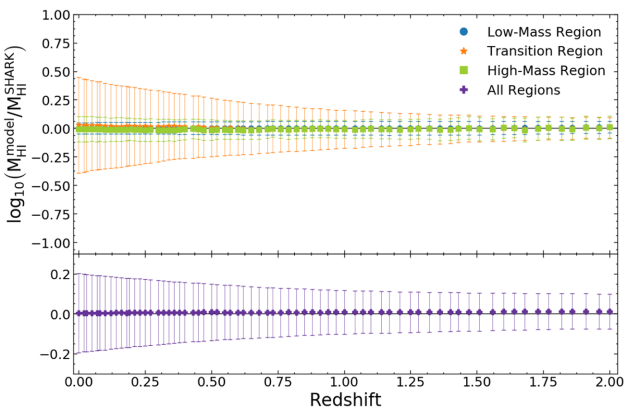


Figure C1. The ratio between the true HI masses of haloes in SHARK-ref and the derived masses from equations (21) to (27), i.e. the HI-mass residuals, as a function of redshift. Symbols with errorbars show the median and 16th–84th percentile range. This is presented for all haloes in the simulation (lower panel), and for each halo mass region separately (top panel), as labelled. For reference, the horizontal lines show equality.

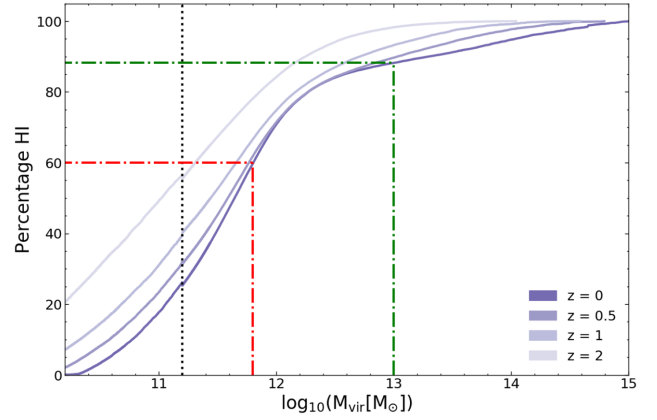


Figure C2. The cumulative fraction of cosmic HI mass contained in haloes as a function of virial mass at four different redshifts, as labelled. At $z = 0$, ~ 60 percent of the HI is contained in haloes with $M_{\text{vir}} < 10^{12} M_{\odot}$, with about ~ 25 per cent lying in the transition region of $10^{12} M_{\odot} \leq M_{\text{vir}} < 10^{13} M_{\odot}$ and the rest in haloes with $M_{\text{vir}} > 10^{13} M_{\odot}$. For reference, these halo mass thresholds are shown with dot–dashed lines. At higher redshift, the contribution from the lower mass region becomes even greater.

halo population (however, by sheer number, the low-mass region dominates the latter). It can be seen that as we move from low to high redshifts, the median of the residuals stays around 0, with small deviations of $\lesssim 0.02$ dex. We also find that the 16th–84th percentile range decreases as we move to higher redshift. This shows that our numerical model is able to successfully capture the dependence of HI mass on halo properties, within certain limits.

In Appendix D, we show the how the HIHM relation changes at higher redshifts. We find that the scatter around the median relation significantly changes for the transition region, as we move to higher redshifts, and this can be encapsulated in equations (E2)–(E4).

C2 Cumulative HI

As seen in Fig. 13, the scatter is well constrained for the low- and high-mass regions, by invoking secondary parameters (the halo’s spin parameter and $M_{\text{h}}^{\text{sat}}/M_{\text{vir}}$, respectively) but at the transition region we find this to be more difficult. It is therefore informative to ask how much of the total HI in SHARK-ref resides in the transition region. In Fig. C2, we plot the cumulative HI mass as a function of halo mass in SHARK-ref. We find that at $z = 0$ –60 percent of the HI is contained in haloes with $M_{\text{vir}} < 10^{11.8} M_{\odot}$, with about ~ 25 per cent lying in the transition region of $10^{11.8} M_{\odot} \leq M_{\text{vir}} < 10^{13} M_{\odot}$. The rest, ~ 15 per cent, is in haloes with masses $M_{\text{vir}} > 10^{13} M_{\odot}$.

As we move to higher redshifts, we find that the low-mass region becomes more important, with contributions that increase from 60 percent at $z = 0$ to 80 percent at $z = 2$.

This shows that, even if our numerical model is less reliable around the transition region, the majority of HI lies in regions that are very well modelled by our numerical method. This is particularly important in, for example, HI stacking or intensity mapping experiments, when the relevant quantity is the aggregated HI mass at a given redshift.

APPENDIX D: REDSHIFT DEPENDENCE OF THE HIHM RELATION

As stated in Section C1, as we move to higher redshifts we find the *Transition Region* getting noticeably smaller in dynamic range, with

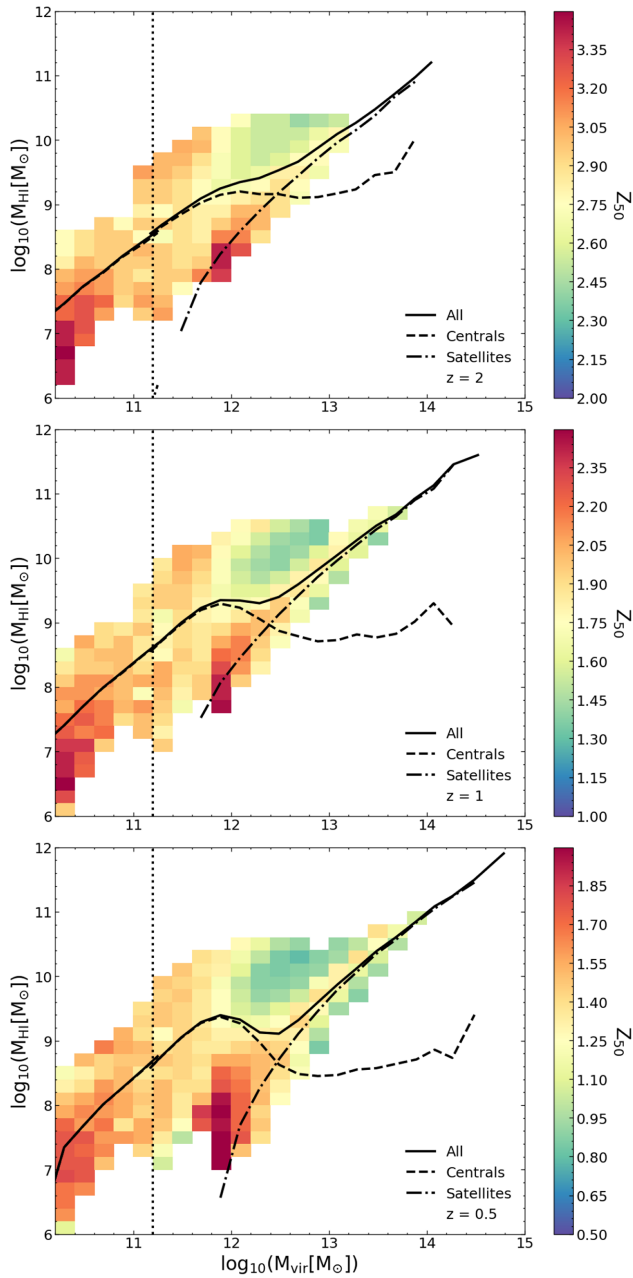


Figure D1. The HIHM relation of haloes in SHARK-ref at $z = 0.5, 1,$ and $2,$ with each bin being coloured by the median z_{50} of the haloes in that bin, as labelled in the colour bar. The solid line represents the median HI mass of the halo as a function of $M_{\text{vir}},$ while the dashed and dash-dotted lines represent the central and satellite galaxies contributions, respectively. The vertical dotted line shows the transition from micro-SURFS to medi-SURFS at lower and higher halo masses, respectively. Though not much can be seen, there is a slight trend with the younger formed haloes being more HI rich than the older ones.

the scatter around the relation decreasing as well. We have showed earlier in Section 5, that the residual fits for the HIHM relation are redshift dependent, though the halo properties comprising the residual fits remain the same throughout the redshift range in consideration.

Fig. D1 shows the $M_{\text{HI}}-M_{\text{vir}}$ relation at $z = 0.5, 1,$ and $2,$ colouring each bin with the median formation age. As had been seen for the

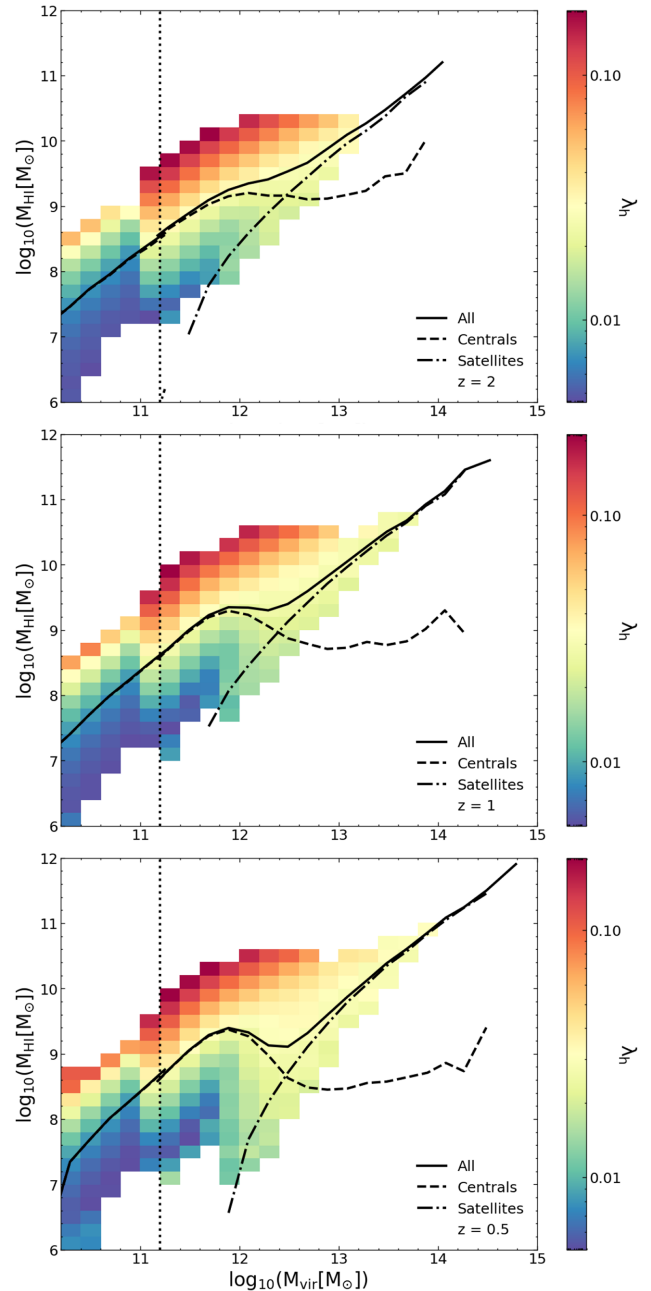


Figure D2. As in Fig. D1 but here bins are coloured by the median halo’s spin parameter, as labelled in the colour bar. There is a strong correlation between the HI mass and the spin parameter at fixed halo mass for haloes with $M_{\text{vir}} < 10^{12} M_{\odot}$ at $z = 0.5,$ with the M_{vir} threshold being $10^{12.5}$ and $\sim 10^{13} M_{\odot}$ for $z = 1$ and $2,$ respectively. Haloes with higher spin parameters are HI-rich than their counterparts.

$z = 0$ case, z_{50} does not show a very strong trend at low-mass region, though a slight trend is noticeable in transition and higher mass regions, with (relatively) younger haloes having higher HI than their older counterparts, throughout the redshift range in consideration.

As we move towards halo spin parameter in Fig. D2, we find that the spin parameter is strongly correlated with the scatter of low-mass region in the HIHM relation. One interesting aspect of the correlation seen is that as we move to higher redshifts, we find the spin parameter correlation extending to higher halo masses than seen in the lower

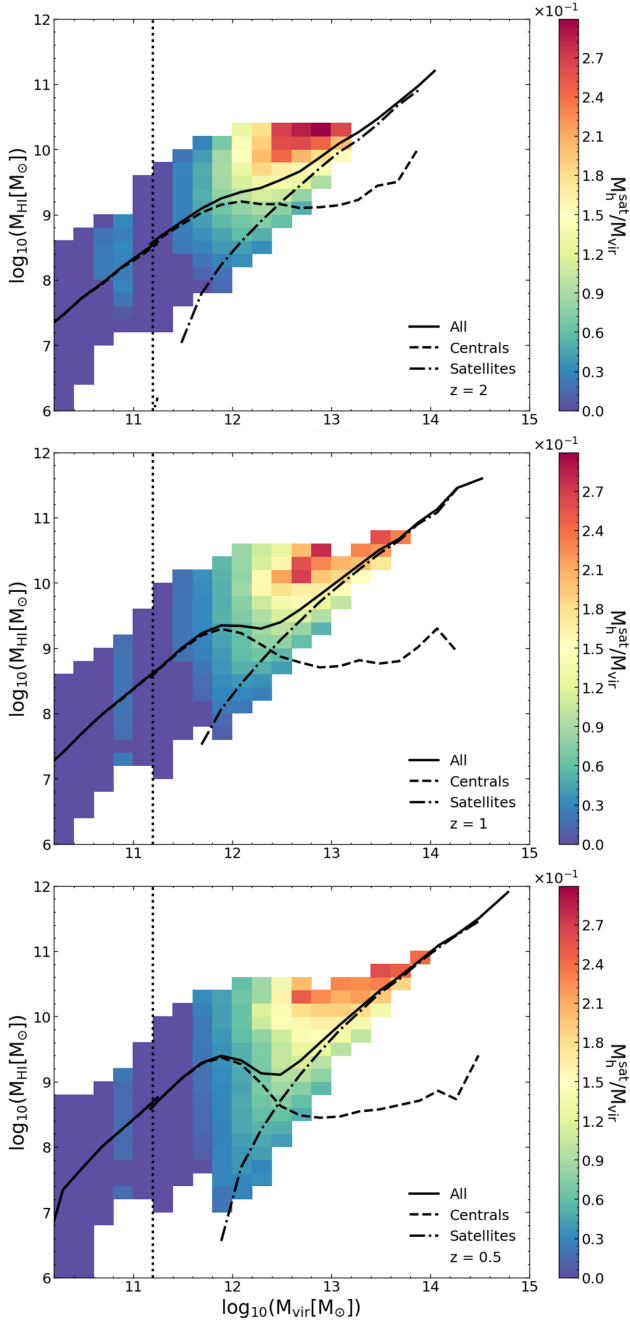


Figure D3. Similar to the earlier plots (see Fig. D1 and D2), here the contribution of HI contained in satellites to the total HI in the halo containing them. As we reach to higher virial masses, we can see that satellites contain most of the HI in the haloes, irrespective to which redshift it is being observed at.

redshift range. As opposed to halo spin parameter showing strong correlation with haloes of masses $M_{\text{vir}} < 10^{12} M_{\odot}$ at $z = 0$, we find the correlation goes as far as halo mass range of $M_{\text{vir}} < 10^{13} M_{\odot}$ at $z = 2$. This is in agreement to our assessment that as we move

to higher redshifts, the *Transition Region* gets smaller and move towards higher halo masses.

Similar to Figs D1 and D2, when we look at the evolution of the $M_{\text{h}}^{\text{sat}}/M_{\text{vir}}$ trend with redshift in Fig. D3, we find it more or less similar to what was seen at $z = 0$: the higher the value of $M_{\text{h}}^{\text{sat}}/M_{\text{vir}}$ the higher the HI mass in the halo. This is due to the fact that, as we move to higher halo masses, the number of satellites in those haloes increases, and thus does the total HI contribution of the satellites.

The evolution of the scatter around the HIHM relation, especially for the *transition region*, through the redshift points to the fact that the flaring of scatter in the transition region at $z = 0$ can be related to the AGN feedback efficiency adopted by the model. As we go higher in redshift, AGN feedback becomes less important leading to a decrease in the scatter around the transition region. This effect is also evident in the noticeable bump that is prominent in the $z = 0$ and 0.5, is smoothed out by the time we reach $z = 2$.

Therefore, from Figs D1, D2, and D3, it is clear that the trends of $z = 0$ persist towards at higher redshifts, which means that we can use the same secondary parameters to fit the scatter around the HIHM relation at different redshifts.

APPENDIX E: PARAMETER FITS

In Section 5.1.2, we pointed out that the dependence of the median relation parameters of the quintic polynomial fit for the transition region is hard to parametrize as a function of redshift, and thus we tabulate the coefficients in Table E1. The equation for estimating the HI in the transition region is as follows:

$$f_{\text{M}_{\text{HI}}}(M_{\text{vir}}, z) = 9 + \sum_{i=0}^n a_i(z) (\log_{10}(M_{\text{vir}}) - 11.8)^i, \quad (\text{E1})$$

where $n = 5$, with $a_1^{\text{TR}} = 0$.

Fig. E1 compares the true HI content of SHARK-ref haloes at $z = 2$, $z = 1$, and $z = 0.5$ with the outcome of applying our numerical HIHM scaling relation to the same underlying halo population (see equations 19–27). Fig. E1 showcases that as we move towards higher redshift the scatter around the HI relation decreases considerably for the transition region, and the shape also evolves into a monotonically increasing relation by the time we reach $z = 2$.

We find that the vertical scatter around the HIHM relation obtained from our numerical model decreases in a similar manner, and can be described by the following functions of redshift, with parameters that depend on the mass region

$$\sigma_{\text{low}} = 0.189 - 0.017z, \quad (\text{E2})$$

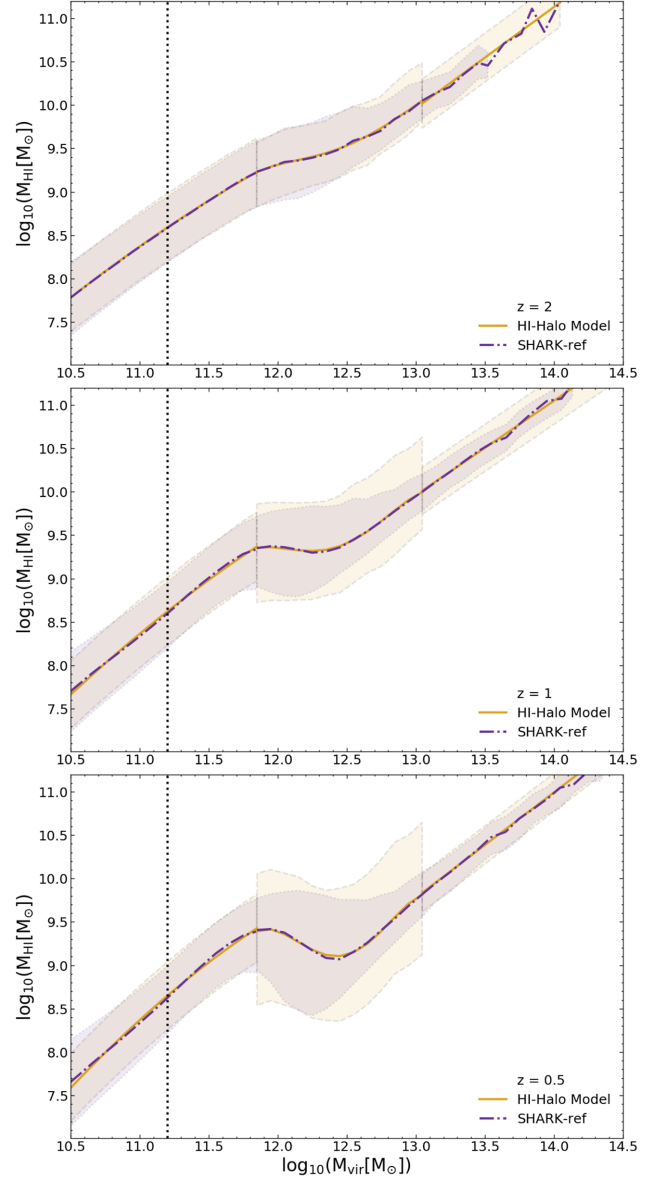
$$\sigma_{\text{TR}} = 0.138 + 0.771e^{-z}, \quad (\text{E3})$$

$$\sigma_{\text{high}} = 0.185 + 0.142e^{-z}, \quad (\text{E4})$$

with z being the redshift. Here, ‘low’, ‘TR’, and ‘high’ refer to the low-mass, transition, and high-mass regions, respectively. This also shows that our numerical model becomes more reliable in the transition region as the redshift increases.

Table E1. Parameters for the quintic polynomial fit for the transition region.

z	a_0^{TR}	a_2^{TR}	a_3^{TR}	a_4^{TR}	a_5^{TR}
2	0.247	2.033	-4.390	4.303	-1.409
1.96	0.249	1.810	-3.651	3.517	-1.140
1.91	0.260	1.319	-2.065	1.758	-0.496
1.86	0.261	1.448	-2.718	2.604	-0.826
1.77	0.269	1.393	-2.848	3.008	-1.054
1.73	0.275	1.180	-2.365	2.641	-0.962
1.68	0.286	0.768	-1.120	1.303	-0.476
1.64	0.297	0.454	-0.309	0.493	-0.183
1.6	0.304	0.235	0.170	0.168	-0.117
1.56	0.311	0.072	0.448	0.029	-0.103
1.51	0.316	-0.237	1.390	-0.908	0.198
1.43	0.328	-0.376	1.411	-0.582	-0.017
1.4	0.329	-0.228	0.683	0.394	-0.415
1.36	0.331	-0.155	0.265	0.930	-0.615
1.32	0.335	-0.023	-0.363	1.675	-0.876
1.28	0.338	-0.112	-0.127	1.436	-0.791
1.21	0.347	-0.635	1.419	-0.158	-0.237
1.17	0.354	-0.763	1.753	-0.536	-0.080
1.14	0.357	-0.796	1.725	-0.426	-0.138
1.1	0.369	-1.179	2.698	-1.265	0.101
1.07	0.370	-1.169	2.527	-0.984	-0.020
1	0.376	-1.128	1.968	-0.080	-0.413
0.97	0.380	-1.332	2.515	-0.598	-0.246
0.94	0.386	-1.410	2.483	-0.415	-0.335
0.91	0.389	-1.606	2.973	-0.868	-0.189
0.88	0.388	-1.479	2.227	0.162	-0.603
0.85	0.393	-1.666	2.628	-0.170	-0.500
0.82	0.398	-1.852	2.961	-0.399	-0.441
0.79	0.403	-1.862	2.682	0.075	-0.643
0.76	0.408	-2.031	2.939	-0.051	-0.625
0.73	0.411	-2.108	2.989	-0.005	-0.661
0.71	0.412	-2.037	2.434	0.782	-0.976
0.68	0.417	-2.144	2.568	0.711	-0.954
0.65	0.424	-2.457	3.239	0.202	-0.823
0.62	0.428	-2.428	2.926	0.611	-0.970
0.6	0.434	-2.317	2.216	1.536	-1.317
0.57	0.438	-2.307	1.885	2.024	-1.508
0.55	0.442	-2.393	1.881	2.166	-1.580
0.52	0.442	-2.144	0.902	3.222	-1.932
0.5	0.442	-1.976	0.071	4.227	-2.292
0.47	0.448	-2.020	-0.083	4.526	-2.416
0.45	0.451	-2.040	-0.175	4.677	-2.475
0.43	0.455	-2.059	-0.333	4.918	-2.565
0.4	0.460	-2.223	0.039	4.573	-2.448
0.38	0.466	-2.177	-0.460	5.301	-2.742
0.36	0.466	-2.233	-0.294	5.070	-2.640
0.34	0.473	-2.352	-0.178	5.045	-2.642
0.32	0.476	-2.309	-0.602	5.616	-2.855
0.3	0.476	-2.214	-1.053	6.138	-3.039
0.27	0.479	-2.395	-0.664	5.812	-2.942
0.25	0.481	-2.364	-0.927	6.117	-3.039
0.23	0.480	-2.152	-1.763	7.053	-3.368
0.21	0.486	-2.135	-1.968	7.281	-3.435
0.19	0.489	-2.004	-2.542	7.935	-3.665
0.18	0.490	-1.951	-2.868	8.322	-3.798
0.16	0.492	-2.209	-2.055	7.386	-3.441
0.14	0.493	-1.922	-3.068	8.454	-3.801
0.12	0.495	-1.711	-4.016	9.611	-4.235
0.1	0.501	-1.919	-3.551	9.211	-4.111
0.08	0.506	-2.079	-3.089	8.683	-3.909
0.07	0.508	-1.973	-3.669	9.425	-4.192
0.05	0.514	-2.043	-3.732	9.636	-4.293
0.03	0.517	-2.372	-2.641	8.373	-3.818
0.02	0.518	-2.231	-3.149	8.923	-4.013
0	0.524	-2.529	-2.355	8.173	-3.774


Figure E1. Overall $H\text{I}$ content of haloes as a function of halo mass for SHARK-ref (dot-dashed line), and predicted by our numerical model (solid line) at $z = 0.5, 1$ and 2 , as labelled. The shaded regions represent the 16th–84th percentile ranges of the distributions. A decrease in the scatter around the transition region is seen as we move towards higher redshifts.

 This paper has been typeset from a $\text{T}_{\text{E}}\text{X}/\text{L}^{\text{A}}\text{T}_{\text{E}}\text{X}$ file prepared by the author.