

Galaxy bimodality due to cold flows and shock heating

Avishai Dekel[★] and Yuval Birnboim[★]

Racah Institute of Physics, The Hebrew University, Jerusalem, Israel

Accepted 2005 December 16. Received 2005 December 12; in original form 2004 December 13

ABSTRACT

We address the origin of the robust bimodality observed in galaxy properties about a characteristic stellar mass $\sim 3 \times 10^{10} M_{\odot}$. Less massive galaxies tend to be ungrouped *blue* star forming discs, while more massive galaxies are typically grouped *red* old-star spheroids. Colour–magnitude data show a gap between the red and blue sequences, extremely red luminous galaxies already at $z \sim 1$, a truncation of today’s blue sequence above L_* , and massive starbursts at $z \sim 2$ –4. We propose that these features are driven by the *thermal* properties of the inflowing gas and their interplay with the clustering and feedback processes, all functions of the dark matter halo mass and associated with a similar characteristic scale. In haloes below a critical shock-heating mass $M_{\text{shock}} \lesssim 10^{12} M_{\odot}$, discs are built by *cold streams*, not heated by a virial shock, yielding efficient early star formation. It is regulated by supernova feedback into a long sequence of bursts in blue galaxies constrained to a ‘fundamental line’. Cold streams penetrating through hot media in $M \geq M_{\text{shock}}$ haloes preferentially at $z \geq 2$ lead to massive starbursts in $L > L_*$ galaxies. At $z < 2$, in $M > M_{\text{shock}}$ haloes hosting groups, the gas is heated by a virial shock, and being dilute it becomes vulnerable to feedback from energetic sources such as active galactic nuclei. This shuts off gas supply and prevents further star formation, leading by passive evolution to ‘red-and-dead’ massive spheroids starting at $z \sim 1$. A minimum in feedback efficiency near M_{shock} explains the observed minimum in M/L and the qualitative features of the star formation history. The cold flows provide a hint for solving the angular momentum problem. When these processes are incorporated in simulations they recover the main bimodality features and solve other open puzzles.

Key words: shock waves – cooling flows – galaxies: evolution – galaxies: formation – galaxies: haloes – dark matter.

1 INTRODUCTION

Observations reveal a robust bimodality in the galaxy population, being divided into two classes, the ‘blue’ and ‘red’ sequences, at a characteristic stellar mass $M_{s,\text{crit}} \simeq 3 \times 10^{10} M_{\odot}$. This corresponds to a dark halo mass $M_{\text{crit}} \lesssim 10^{12} M_{\odot}$ and a virial velocity $V_{\text{crit}} \simeq 120 \text{ km s}^{-1}$ today. Less massive galaxies tend to be blue, star forming discs residing in the ‘field’. Their properties are correlated along a ‘fundamental line’ of decreasing surface brightness, internal velocity and metallicity with decreasing luminosity. Galaxies above M_{crit} are dominated by spheroids of red, old stars, with high surface brightness and metallicity independent of luminosity. They tend to reside in the high-density environments of groups and clusters and they preferentially host active galactic nucleus/nuclei (AGN/AGNs).

Current models of galaxy formation have difficulties in reproducing this bimodality and the broad colour distribution observed.

In particular, the extremely red bright ellipticals which start showing up already at $z \sim 1$ are not predicted. They require efficient star formation at earlier epochs, followed by an effective shutdown of star formation in massive galaxies. The observations also reveal very blue galaxies in excess of the predictions, indicating repeating starbursts over the lifetimes of galaxies. Today’s blue sequence is non-trivially truncated at the bright end, while at $z \geq 2$ there are indications for very luminous starbursts in big objects, both posing severe theoretical challenges.

1.1 The observed bimodality

The bimodality or transition in galaxy properties is observed in many different ways. We list here the main relevant observed features¹ which we address in this paper (Section 6).

(i) **Luminosity functions.** Blue galaxies dominate the stellar mass function below $M_{s,\text{crit}}$ while red galaxies take over above it

[★]E-mail: dekel@phys.huji.ac.il (AD); yuval@phys.huji.ac.il (YB)

¹ Quoting only sample references, making no attempt to be complete.

[Bell et al. 2003c; Baldry et al. 2004, in Sloan Digital Sky Survey (SDSS) and two-Micron All Sky Survey (2MASS)]. The transition occurs slightly below L_* , the characteristic luminosity of the brightest disc galaxies, beyond which the luminosity function drops.

(ii) **Colour–magnitude.** A colour bimodality shows robustly in Colour–magnitude diagrams, where the galaxies are divided into a blue sequence and a red sequence separated by a gap. In SDSS (Blanton et al. 2005; Baldry et al. 2004), the gap is at $u - r \sim 2$. The colour distribution is non-trivially *broad*, with the red tip stretching beyond $u - r = 2.5$ and the blue tail reaching well below $u - r = 1.0$.

(iii) **Star formation rate.** The current star formation rate (SFR), and the typical age of the stellar population, show a robust bimodality about $M_{s,crit}$. The less massive galaxies are dominated by young populations, while the more massive galaxies are dominated by old stars [Kauffmann et al. 2003b; Madgwick et al. 2003b, in SDSS and two-degree Field (2dF)], in agreement with the colour bimodality. A similar bimodality is seen in the gas-to-stellar mass fraction, which is high in the blue sequence and low in the red sequence, steeply increasing with stellar mass below $M_{s,crit}$, and only moderately so above it (Kannappan 2004, in SDSS+2MASS).

(iv) **Colour–magnitude at $z \sim 1$.** The colour bimodality is similar back to $z \sim 1.5$ (Bell et al. 2004, in COMBO17). Extremely red massive galaxies exist at the bright tip of the red sequence already at $z \sim 1$ (e.g. Moustakas et al. 2004). Very blue small galaxies indicating starbursts show in the blue sequence (e.g. Ferguson & Babul 1998; Fioc & Rocca-Volmerange 1999).

(v) **Massive starbursts at high z .** Very luminous and massive dusty objects are detected at $z \sim 2-4$, indicating an excessive activity of star formation in surprisingly big objects (Smail et al. 2002; Chapman et al. 2003, 2004; Shapley et al. 2004, LBG and SCUBA sources).

(vi) **Star formation history.** The cosmological history of SFR has a broad maximum near $z \sim 1-2$, followed by a sharp drop from $z \sim 1$ to $z = 0$ (e.g. Madau et al. 1996; Dickinson et al. 2003; Giavalisco et al. 2004; Hartwick 2004; Heavens et al. 2004). Still, about half the stars in today’s spirals seem to have formed after $z \sim 1$, e.g. in luminous infrared galaxies (LIRGs) near $M_{s,crit}$ (Hammer et al. 2005). Massive galaxies tend to form their stars earlier than smaller galaxies (‘downsizing’ Thomas et al. 2005).

(vii) **Bulge-to-disc ratio.** The galaxy bulge-to-disc ratio shows a transition from disc dominance in the blue sequence below $M_{s,crit}$ to spheroid dominance in the red sequence (Blanton et al. 2005; Kauffmann et al. 2003b, in SDSS).

(viii) **Environment dependence.** The distributions in colour and SFR depend strongly on the galaxy density in the ~ 1 Mpc vicinity: the blue and red sequence galaxies tend to populate low- and high-density environments, respectively (Blanton et al. 2005; Hogg et al. 2003; Balogh et al. 2004; Blanton et al. 2004; Kauffman et al. 2004, in SDSS). The colour–environment correlation is stronger than the morphology–density relation Dressler (1980).

(ix) **Halo mass and HOD.** The environment density is correlated with the mass of the host dark matter (DM) halo, where haloes less massive than $\sim 10^{12} M_\odot$ typically host one dominant galaxy each while more massive haloes tend to host groups and clusters of luminous galaxies, as quantified by the halo occupation distribution (HOD, Yan, Madgwick & White 2003; Kravtsov et al. 2004; Abazajian et al. 2005, in 2dF, SDSS and in simulations). The environment dependence thus implies that the galaxy properties are correlated with the host halo mass, with the bimodality at $M_{crit} \lesssim 10^{12} M_\odot$ (Blanton et al. 2004).

(x) **Hot halo gas.** Ellipticals of $L_B > 10^{10.5} L_\odot$ show a significant excess of X-ray flux plausibly associated with hot halo gas

(Ciotti et al. 1991; Mathews & Brighenti 2003). Intergalactic X-ray radiation is detected predominantly in groups where the brightest galaxy is an elliptical. Group properties have a transition near $\sigma_v \sim 140 \text{ km s}^{-1}$ (Helsdon & Ponman 2003; Osmond & Ponman 2004).

(xi) **Luminosity/mass functions.** The stellar mass function has a ‘knee’ near $M_{s,crit}$, where the shallow $dn/dM_s \propto M_s^{-1}$ on the faint side turns into an exponential drop. In contrast, the dark halo mass function is predicted by the standard Lambda cold dark matter (Λ CDM) model to be $dn/dM \propto M^{-1.8}$ everywhere below $\sim 10^{13} M_\odot$. A match at $M_{s,crit}$ requires a baryonic fraction $M_s/M \sim 0.05$, indicating gas loss, and associating $M_{s,crit}$ with $M_{crit} \simeq 6 \times 10^{11} M_\odot$. The halo mass-to-light function has a minimum near M_{crit} , varying as $M/L \propto M^{-2/3}$ and $\propto M^{+1/2}$ below and above it, respectively, thus implying increased suppression of star formation away from the critical mass on both sides (Marinoni & Hudson 2002; Bell et al. 2003a,b; Yang, Mo & van den Bosch 2003, in a B-mag sample and in SDSS, 2MASS, 2dF).

(xii) **Fundamental line versus plane.** A transition is detected in the galaxy structural scaling relations near $M_{s,crit}$, e.g. the surface brightness changes from $\mu_s \propto M_s^{0.6}$ at lower masses to $\mu_s \sim \text{constant}$ at the bright end (Kauffmann et al. 2003b, in SDSS). The correlation below $M_{s,crit}$ is part of the ‘fundamental line’ relating stellar mass, radius and velocity over five decades in M_s (e.g. Dekel & Woo 2003). The mean metallicity shows a transition near a similar mass scale from $Z \sim M_s^{0.4}$ to $Z \sim \text{constant}$ (Dekel & Woo 2003; Tremonti et al. 2004, in SDSS and the Local Group).

(xiii) **AGNs.** Black hole masses are correlated with their host spheroid properties (e.g. Tremaine et al. 2002). The optical AGN population, with high accretion rate and SFR, peaks near $M_{s,crit}$ with little AGN activity at smaller masses, and is associated with black hole masses $\lesssim 10^8 M_\odot$ (Kauffmann et al. 2003a, in SDSS). Radio-loud AGNs, uncorrelated with the optical activity and the SFR, dominate in larger haloes hosting $\sim 10^{8-9} M_\odot$ black holes (G. Kauffmann private communication).

1.2 Key physical processes

The bimodality imprinted on almost every global property of galaxies deserves a simple theoretical understanding. We propose that the main source of the bimodality is the transition from cold flows to virial shock heating at a critical scale, in concert with feedback processes and gravitational clustering that emphasize the same characteristic scale. We address the cross-talk between these processes, and integrate them into a scenario which attempts to address simultaneously the variety of observed phenomena. The key processes are:

(i) **Cold infall versus hot medium.** The thermal behaviour of the gas as it falls through the halo is qualitatively different below and above a critical mass scale of $M_{shock} \lesssim 10^{12} M_\odot$ (Birnboim & Dekel 2003; Kereš et al. 2005). In less massive haloes, the disc is built by cold flows ($\sim 10^{4-5}$ K), which are likely to generate early bursts of star formation. In more massive haloes, the infalling gas is first heated by stable shocks to near the virial temperature ($\sim 10^6$ K). Near and above M_{shock} at $z \geq 2$ (and preferentially in isolated galaxies), streams of dense cold gas penetrate through the dilute shock-heated medium (Fardal et al. 2001; Kravtsov 2003; Kereš et al. 2005) (discussed in Sections 2–4).

(ii) **Gravitational clustering.** Non-linear gravitational clustering of the DM occurs on a characteristic mass scale, M_* , marking the typical haloes forming at a given epoch and the lower bound

for groups of galaxies. We point out that the clustering scale, which varies rapidly with cosmological time, happens to coincide with M_{shock} at $z \sim 1$, and the interplay between these scales plays a role in determining the galaxy properties (Section 4).

(iii) **Feedback.** We argue that the feedback processes affecting galaxy evolution are relatively ineffective near M_{shock} , largely due to the shock-heating process itself, and they therefore help emphasizing the imprint of this scale on the galaxy properties. Supernova and other feedback processes regulate star formation in the blue sequence below M_{shock} . Feedback by AGNs, or other sources, becomes efficient in haloes more massive than M_{shock} , because it preferentially affects the dilute shock-heated medium and may prevent it from ever cooling and forming stars effects (Section 5).

We show how the introduction of shock-stability physics crystallizes our understanding of the origin of the characteristic scales of galaxies. We argue that the combination of shock heating, feedback and clustering introduces a new feature in galaxy formation modelling – *a complete suppression of cold gas supply in haloes above a critical mass after a critical redshift*. This could be the key to solving many of the open questions posed by the observations, focusing on the bright-end truncation of the luminosity function, the appearance of very red bright galaxies already at $z \sim 1$ at the expense of big blue galaxies, and the indications for massive starbursts at higher redshifts. We note that some of the issues have been addressed in parallel, in a similar spirit and in different ways, by Birnboim & Dekel (2003), Benson et al. (2003), Kereš et al. (2005) and Binney (2004). We make here a more thorough investigation into the cold flows versus shock-heating phenomenon, relate it to the feedback and clustering processes, and attempt an integrated scenario that addresses simultaneously the variety of observed features.

1.3 Outline

In Section 2, we provide an improved presentation of the original analysis of spherical shock stability (Birnboim & Dekel 2003, hereafter BD03). In Section 3, we compute the associated critical mass scale in the cosmological context. In Section 4, we describe how the phenomena is demonstrated in cosmological simulations, and learn about cold filaments in massive hot haloes at high redshift. In Section 5, we address the cross-talk with the relevant feedback processes working alternatively below and above the critical scale. Then, in Section 6, we integrate the above processes into a scenario which attempts to explain the origin of the bimodality and related features, and report first results from simulations that implement the new ingredients. In Section 7, we briefly discuss possible implications on other open issues in galaxy formation, and in Section 8 we summarize our results, the proposed re-engineering of galaxy formation simulations, and the open theoretical issues.

2 SPHERICAL SHOCK-STABILITY ANALYSIS

The standard paradigm of disc formation (Rees & Ostriker 1977; Silk 1977; White & Rees 1978; Blumenthal et al. 1984; White & Frenk 1991; Mo, Mao & White 1998), which lies at the basis of all current models of galaxy formation, assumes that while a DM halo relaxes to a virial equilibrium, the gas that falls in within it is *shock heated* near the halo virial radius R_v to the halo virial temperature. It is then assumed to cool radiatively from the inside out. As long as the cooling time is shorter than a certain global free-fall time (or the Hubble time), typically inside a current ‘cooling radius’, the gas is assumed to accrete gradually on to a central disc and then form stars

in a quiescent way. The maximum halo mass for efficient cooling was estimated to be of the order of $\sim 10^{12-13} M_{\odot}$, and the common wisdom has been since then that this explains the upper bound for disc galaxies. However, early hints, based on one-dimensional simulations, indicated that this scenario cannot reproduce the sharp drop in the luminosity function above this scale (Thoull & Weinberg 1995). Even earlier studies, valid in the context of the pancake picture of galaxy formation, indicated that virial shock heating may not be as general as assumed (Binney 1977). More advanced cosmological simulations have started to reveal the presence of cold flows (Fardal et al. 2001). With the new data from big surveys such as SDSS, 2MASS and 2dF, and the detailed semi-analytic modelling (SAM) of galaxy formation, it is becoming clear that the observed scale is somewhat smaller and the drop is sharper than predicted by the original picture. It seems that the current models based on the standard paradigm have hard time trying to reproduce many of the observed bimodality features summarized in Section 1. This motivated us to attempt a closer look at the shock-heating mechanism.

2.1 Spherical simulations

Fig. 1 shows the time evolution of the radii of Lagrangian gas shells in a spherical gravitating system consisting of gas (in this case with

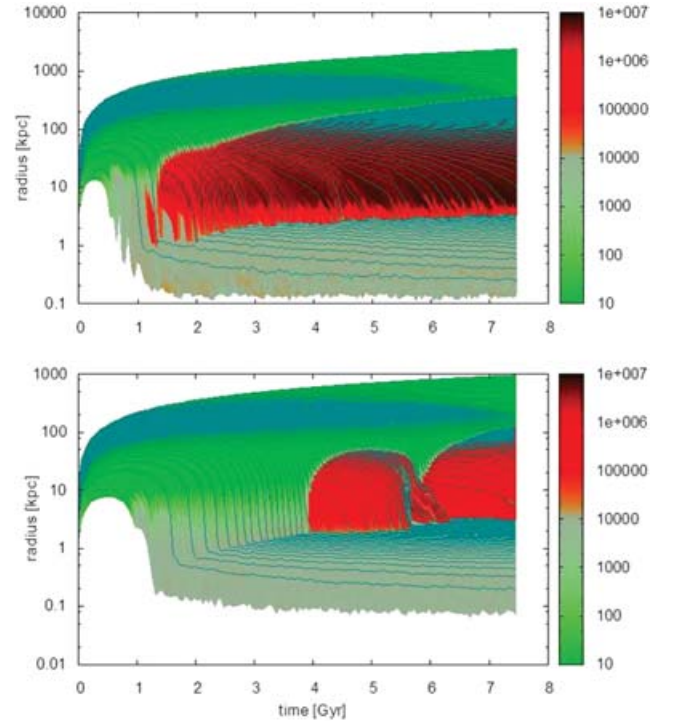


Figure 1. Time evolution of the radii of Lagrangian gas shells (lines) in a spherical simulation of a protogalaxy consisting of primordial gas ($Z = 0$) and DM. Temperature is marked by colour. A shock shows up as a sharp break in the flow lines, namely a sudden slowdown of the infall, associated with an abrupt increase in the temperature. The lower discontinuity where the inflow is brought to a final halt marks the ‘disc’ radius, formed due to an artificial centrifugal force. (a) A massive system, where the virialized mass grows from 10^{11} to $10^{13} M_{\odot}$. (b) A less massive system, growing from 10^{10} to $10^{12} M_{\odot}$. A virial shock exists only in systems more massive than a critical mass, while in smaller haloes the gas flows cold and unperturbed into the inner halo. With more realistic metallicities the critical mass becomes $\sim 10^{12} M_{\odot}$.

primordial composition) and DM, similar to the original simulations by BD03 using an accurate one-dimensional hydrodynamical code. Not shown are the dissipationless DM shells, which detach from the cosmological expansion, collapse and then oscillate into virial equilibrium such that they deepen the potential well attracting the dissipating gas shells. The initial density perturbation was assumed to have a profile proportional to the linear two-point correlation function of Λ CDM (justified in Dekel 1981), and the final density profile roughly mimics the universal NFW profile (Navarro, Frenk & White 1997) seen in cosmological simulations. The gas is cooling radiatively based on the atomic cooling function computed by Sutherland & Dopita (1993). The collapse of each gas shell is stopped at roughly $0.05 R_v$ by an artificial centrifugal force which mimics the formation of a central disc.

The upper panel focuses on massive haloes of $\sim 10^{12} M_\odot$. As expected in the common picture, a strong shock exists near the virial radius, namely at roughly half the maximum-expansion radius of the corresponding shell. The virial shock gradually propagates outwards, encompassing more mass in time. The hot post-shock gas is in a quasi-static equilibrium, pressure supported against gravitational collapse. The lower panel focuses on halo masses smaller by an order of magnitude, and shows an interesting new phenomenon. A stable shock forms and inflates from the disc towards the virial radius only after a total mass of more than a few times $10^{11} M_\odot$ has collapsed. In less massive systems, the cooling rate is faster than the compression rate required for restoring the pressure in the post-shock gas; had there been a shock, the post-shock gas would have become unstable against radial gravitational contraction and unable to support the shock. In the specific case shown, with zero metallicity, the critical mass is biased low; with more realistic metallicities it becomes $\sim 10^{12} M_\odot$ due to the more efficient cooling via metal lines (Section 3).

We demonstrate below (Section 4) that the behaviour seen in the spherical simulations is robustly reproduced in cosmological simulations. But first we wish to understand the spherical case, and use it for simple predictions.

2.2 Post-shock stability criterion

The existence or absence of a shock, as seen in the simulations, can be evaluated via a straightforward stability analysis of the post-shock gas (first introduced in BD03). We provide here a brief, improved presentation of this analysis, followed by a more detailed estimate of the predicted critical mass for shock stability as a function of redshift (Section 3).

A stable extended shock can exist when the pressure in the post-shock gas is sufficient to balance the gravitational attraction towards the halo centre. The standard equation of state for an ideal gas expresses the pressure as $P = (\gamma - 1)\rho e$ where ρ and e are the gas density and internal energy per unit mass, and $\gamma = 5/3$ for a monatomic gas. In the textbook case of no cooling, the adiabatic index is defined as $\gamma \equiv (\partial \ln P / \partial \ln \rho)_{\text{ad}}$, and the system is known to be gravitationally stable once $\gamma > 4/3$. When there is energy loss (e.g. by radiation) at a rate q per unit mass, we define a new quantity along the particle trajectories:

$$\gamma_{\text{eff}} \equiv \frac{d \ln P}{d \ln \rho} = \gamma - \frac{\rho q}{\dot{\rho} e}. \quad (1)$$

The second equality follows from energy conservation, $\dot{e} = -P\dot{V} - q$ (where $V = 1/\rho$), plugged into the equation of state. Note that $\gamma_{\text{eff}} = \gamma$ when $q = 0$. The difference between the two is a ratio of characteristic rates for the two competing processes: the cooling,

which reduces the pressure in the post-shock gas, and the compression due to the pattern of the post-shock infall, which tends to increase the pressure. If the compression rate is efficient compared to the cooling-loss rate, it restores the pressure necessary for supporting a stable extended shock, but otherwise the post-shock gas collapses inwards under gravity, failing to support the extended shock.

It is convenient to express the *compression rate* in the post-shock region as the inverse of a compression time, which we define by

$$t_{\text{comp}} \equiv \Gamma \frac{\rho}{\dot{\rho}} \quad \Gamma \equiv \frac{3\gamma + 2}{\gamma(3\gamma - 4)} = \frac{21}{5}, \quad (2)$$

with the factor Γ to be justified below, and the last equality referring to $\gamma = 5/3$. For a spherical shock at radius r_s , and a post-shock radial velocity u_1 , we assume that the radial flow pattern in the post-shock region is homologous, $u/r = u_1/r_s$. This is justified based on the spherical simulations described above, where the log-linear post-shock flow lines in Fig. 1 are nearly parallel straight lines. We then obtain using continuity

$$t_{\text{comp}} = \frac{\Gamma}{(-\nabla \cdot \mathbf{u})} = \frac{\Gamma r_s}{(-3u_1)}. \quad (3)$$

The competing *cooling rate* in the post-shock region is expressed as the inverse of the standard radiative cooling time defined by

$$t_{\text{cool}} \equiv \frac{e}{q}, \quad (4)$$

where $e = e(T)$ and $q \propto \rho \Lambda(T, Z)$, functions of temperature T and metallicity Z . Then in equation (1)

$$\gamma_{\text{eff}} = \gamma - \Gamma^{-1} \frac{t_{\text{comp}}}{t_{\text{cool}}}. \quad (5)$$

In order to test for stability, BD03 performed a perturbation analysis where the radius of a shell is perturbed by $r \rightarrow r + \delta r$ and the sign of the force, $\delta \ddot{r} / \delta r$, is computed. Writing $\delta r = u \delta t$, using the homology, and assuming further that the gravity and pressure forces balance each other near the transition state between stability and instability, $\rho^{-1} \nabla P = GM/r^2$, one obtains a restoring force, i.e. stability, for

$$\gamma_{\text{eff}} > \gamma_{\text{crit}} \equiv \frac{2\gamma}{\gamma + 2/3} = \frac{10}{7}. \quad (6)$$

The $\gamma_{\text{crit}} = 10/7$ replaces the standard $\gamma_{\text{crit}} = 4/3$ of the adiabatic case.²

Using equation (5) and the definitions of the time-scales above, the *shock stability criterion* of equation (6) becomes the simple condition that *the cooling rate should be slower than the compression rate*:

$$t_{\text{cool}} > t_{\text{comp}}. \quad (7)$$

Once the cooling rate is slower, the pressure gain by compression can balance the loss by radiative cooling, which allows the post-shock gas to be stable against global gravitational collapse and thus support the shock. The factor $\Gamma = 21/5$ has been introduced in the definition of t_{comp} , equation (2), in order to simplify this final expression.

² If the spherical symmetry assumed above is replaced by planar symmetry, both for the shock and the gravitational field, the stability criterion $\gamma_{\text{eff}} > 10/7$ is replaced by $\gamma_{\text{eff}} > 10/11$ (Birnbom et al., in preparation). One can therefore assume in general that the actual critical value lies somewhere between these two limits; if $\gamma_{\text{eff}} < 10/11$ there is no stable shock, if $\gamma_{\text{eff}} > 10/7$ the conditions allow a stable shock, and if $10/11 < \gamma_{\text{eff}} < 10/7$ the shock stability depends on the local geometry.

Note that the relevant quantity for stability is the *ratio* of rates associated with the two competing processes, independent of how slow each of them actually is in absolute terms. Each of the characteristic times could in principle be longer than the Hubble time – it is their ratio which determines whether a stable shock is possible or the gas falls in subject to gravity, cold and unperturbed.

2.3 Pre-shock quantities

2.3.1 Compression rate

Using the standard jump conditions across a strong shock, we can express the characteristic time-scales (or γ_{eff}) in terms of the pre-shock gas quantities. The jump condition for the radial velocity is

$$u_0 - u_s = \frac{\gamma + 1}{\gamma - 1}(u_1 - u_s), \quad (8)$$

where u_s is the radial shock velocity and u_0 is the radial velocity of the pre-shock gas. Then

$$t_{\text{comp}} = \frac{\Gamma(\gamma + 1)}{3(\gamma - 1)} \frac{r_s}{|u_0|} \left(1 - \frac{2}{\gamma - 1} \frac{u_s}{|u_0|}\right)^{-1} \quad (9)$$

$$= \frac{28}{5} \frac{r_s}{|u_0|} (1 - 3\tilde{u}_s)^{-1} \quad (10)$$

$$\simeq 5.48 \text{ Gyr} \frac{r_s}{|u_0|} (1 - 3\tilde{u}_s)^{-1}, \quad (11)$$

where $\tilde{u}_s \equiv u_s/|u_0|$ (see Section 2.3.4) and the last expression assumes $\gamma = 5/3$, r_s in 100 kpc, and u_0 in 100 km s⁻¹. If $u_s = 0$, say, then t_{comp} is about six times larger than $r_s/|u_0|$, a typical free-fall time from r_s into the halo centre. At the virial radius, t_{comp} is comparable to the Hubble time at the corresponding epoch, but at inner radii it becomes significantly shorter.

2.3.2 Cooling rate

The cooling time (e.g. based on Sutherland & Dopita 1993) is

$$t_{\text{cool}} \equiv \frac{e}{q} = \frac{1 + 2\epsilon}{1 + \epsilon} \frac{3}{2} kT \left[\frac{\chi^2}{m} \rho \Lambda(T, Z) \right]^{-1}, \quad (12)$$

where $\Lambda(T, Z)$ is the cooling function, k is the Boltzmann constant, $\epsilon \equiv n_{\text{He}}/n_{\text{H}}$, the mass per particle is $m \equiv \mu m_{\text{p}}$ with $\mu = (1 + 4\epsilon)/(2 + 3\epsilon)$, and the number of electrons per particle is $\chi = (1 + 2\epsilon)/(2 + 3\epsilon)$. For 25 per cent He in mass, one has $\epsilon = 1/12$, yielding $\mu = 0.59$. If we express the post-shock temperature as $T_6 \equiv T/10^6$ K, the post-shock baryon density as $\rho_{-28} \equiv \rho/10^{-28}$ g cm⁻³, and the cooling function as $\Lambda_{-22}(T, Z) \equiv \Lambda(T, Z)/10^{-22}$ erg cm³ s⁻¹, we have

$$t_{\text{cool}} \simeq 2.61 \text{ Gyr} \rho_{-28}^{-1} T_6 \Lambda_{-22}^{-1}(T, Z). \quad (13)$$

The post-shock gas density is related to the pre-shock density by the jump condition

$$\rho_1 = \frac{\gamma + 1}{\gamma - 1} \rho_0 = 4\rho_0, \quad (14)$$

and the post-shock temperature entering the cooling time is related to the pre-shock radial velocity u_0 via

$$\frac{kT_1}{m} = \frac{2(\gamma - 1)}{(\gamma + 1)^2} (u_0 - u_s)^2 = \frac{3}{16} u_0^2 (1 + \tilde{u}_s)^2. \quad (15)$$

We note in passing that for a virial shock, where $u_0 = -V_{\text{v}}$ (BD03), the post-shock temperature is actually

$$T_1 \gtrsim \frac{3}{8} T_{\text{v}}, \quad (16)$$

comparable to but somewhat smaller than the virial temperature as defined in equation (A9).

2.3.3 Stability criterion

Using equations (11) and (13), the critical stability condition becomes

$$\frac{t_{\text{cool}}}{t_{\text{comp}}} \simeq 0.48 \frac{\rho_{-28}^{-1} T_6 \Lambda_{-22}^{-1}(T, Z)}{r_s |u_0|^{-1} (1 - 3\tilde{u}_s)^{-1}} \simeq 1, \quad (17)$$

with r_s in 100 kpc and u_0 in 100 km s⁻¹. Recall that equation (14) relates ρ to ρ_0 , and equation (15) relates T to u_0 . Thus, for given shock radius r_s , shock velocity relative to infall $u_s/|u_0|$, and pre-shock gas density ρ_0 , once the metallicity Z is given and the cooling function $\Lambda(T, Z)$ is known, one can solve equation (17) for the critical values of T and the corresponding u_0 . When put in a cosmological context (Section 3), this solution is associated with a unique critical halo mass.

The stability criterion derived above, equation (6) or equation (7), is found to work very well when compared to the results of the spherical simulations shown in Fig. 1. When γ_{eff} (or $t_{\text{cool}}/t_{\text{comp}}$) is computed using pre-shock quantities at a position just outside the ‘disc’, we find that as long as the halo is less massive than a critical scale, before the shock forms, the value of γ_{eff} is indeed well below γ_{crit} and is gradually rising, reaching γ_{crit} almost exactly when the shock starts propagating outwards. The value of the γ_{eff} computed using the quantities just outside the shock then oscillates about γ_{crit} with a decreasing amplitude, following the oscillations in the shock radius seen in Fig. 1. As the shock eventually settles at the virial radius, γ_{eff} approaches 5/3, larger than $\gamma_{\text{crit}} = 10/7$, where the cooling is negligible. The same stability criterion is found to be valid to a good approximation also in cosmological simulations (Section 4).

2.3.4 Shock velocity

What value of \tilde{u}_s is relevant for evaluating stability? In the inner halo, we use $\tilde{u}_s = 0$. This is because, as the halo is growing in mass, the shock first forms in the inner halo and then propagates outwards (Fig. 1b). The onset of shock stability is therefore marked by its ability to develop a velocity outwards.

During the stable phase when the shock is expanding with the virial radius, the spherical simulations indicate roughly $u_1 \simeq -u_s$ (Fig. 1a), namely $\tilde{u}_s \simeq 1/7$ (equation 8). This indicates that a small shock velocity of such a magnitude is appropriate for evaluating stability at the virial radius.

Note that stability is harder to achieve when the shock is expanding relatively fast. In particular, in the extreme case $\tilde{u}_s = 1/3$, the post-shock velocity vanishes, $u_1 = 0$ (equation 8). The compression rate becomes infinitely slow (equation 3), implying that the shock cannot be stabilized.

3 SHOCK-HEATING SCALE IN COSMOLOGY

3.1 Haloes in cosmology

We wish to translate the critical stability condition, equation (7) or equation (17), into a critical post-shock temperature, and the

corresponding critical halo virial velocity and mass as a function of redshift. Equation (17) has a unique solution when combined with the two virial relations between halo mass, velocity and radius (equation A6), and the relation between post-shock temperature and pre-shock infalling velocity (equation 15).

As summarized in Appendix A, the time dependence of the virial relations can be expressed in terms of the convenient parameter

$$A \equiv \left(\Delta_{200} \Omega_{m0.3} h_{0.7}^2 \right)^{-1/3} a, \quad (18)$$

where $a \equiv 1/(1+z)$ is the cosmological expansion factor and the other parameters are of order unity. The parameters $\Omega_{m0.3}$ and $h_{0.7}$ correspond to today's values of the cosmological mass density parameter and the Hubble expansion parameter, respectively, and for the standard Λ CDM cosmology adopted here they are both equal to unity. The parameter Δ_{200} is the virial density factor given approximately in equation (A8): at redshifts $z > 1$ it is $\Delta_{200} \simeq 1$, but at lower redshifts it becomes somewhat larger, reaching $\Delta_{200} \simeq 1.7$ at $z = 0$.

3.2 Compression rate

For a shock at the virial radius, $r_s = R_v$, we approximate $u_0 = -V_v$, as predicted by the spherical collapse model in an Einstein–de Sitter cosmology (BD03, appendix B).

When the shock is at an arbitrary inner radius r , where the infalling velocity is $|u|$, we multiply R_v and V_v by appropriate factors $f_r \equiv r/R_v$ and $f_u \equiv |u|/V_v$ (discussed in Section 3.5). Then equation (11) becomes

$$t_{\text{comp}} \simeq 14.3 \text{ Gyr } A^{3/2} f_r f_u^{-1} (1 - 3\tilde{u}_s)^{-1}. \quad (19)$$

3.3 Cooling rate: gas density

In order to express the cooling time of equation (13) in terms of cosmological quantities, we first evaluate the pre-shock baryon density, which we write as

$$\rho_b = 4 f_b (\rho/\bar{\rho})_{\text{vir}} \Delta \rho_u f_\rho. \quad (20)$$

Here, ρ_u is the universal mean mass density (equation A4), and Δ is the top-hat mean overdensity inside the virial radius (equation A8). The factor $(\rho/\bar{\rho})_{\text{vir}}$ translates $\bar{\rho}$, the mean total density interior to R_v , to ρ , the local total density at R_v . The effective baryonic fraction f_b turns it into a pre-shock baryonic density. The factor 4 stands for the ratio between the post-shock gas density and the pre-shock gas density.³ The factor $f_\rho \equiv \rho(r)/\rho(R_v)$ reflects the ratio of the actual gas density at some radius r within the halo to its value at the virial radius (see Section 3.5).

The ratio $\rho/\bar{\rho}$ at the virial radius is derived for the universal NFW halo density profile revealed by cosmological simulations (Navarro, Frenk & White 1997). For a virial concentration parameter c , this ratio is

$$\left(\frac{\rho}{\bar{\rho}} \right)_{\text{vir}} = \frac{c^2}{3(1+c)^2} \left[\ln(1+c) - \frac{c}{(1+c)} \right]^{-1}. \quad (21)$$

A typical concentration of $c = 12$ is associated with $\rho/\bar{\rho} \simeq 0.17$; we therefore express the approximate results below using the factor $f_{\bar{\rho},0.17} \equiv (\rho/\bar{\rho})/0.17$. In our more accurate evaluation of the critical

scale (Section 3.7), we model the dependence of the mean concentration on mass and time using the fit of Bullock et al. (2001) for the Λ CDM cosmology:

$$c(M, a) = 18 M_{11}^{-0.13} a. \quad (22)$$

The effective baryon fraction f_b may in principle be as large as the universal fraction $\simeq 0.13$, but it is likely to be smaller because of gas loss due to outflows. For the approximate expressions we define $f_{b,0.05} \equiv f_b/0.05$.

The gas density at r , equation (20), thus becomes

$$\rho_{-28} = 0.190 A^{-3} f_{b,0.05} f_\rho f_{\bar{\rho},0.17}. \quad (23)$$

Inserting this baryon density into equation (13), the cooling time becomes

$$t_{\text{cool}} = 13.7 A^3 f_{b,0.05}^{-1} f_\rho^{-1} f_{\bar{\rho},0.17}^{-1} T_6 \Lambda_{-22}^{-1}(T, Z) \text{ Gyr}. \quad (24)$$

The cooling function that we use below (based on Sutherland & Dopita 1993) neglects two physical processes: Compton scattering off the cosmic microwave background and the possible effect of external radiation on the cooling rate through the re-ionization of hydrogen. Based on the more complete cooling function as implemented by Kravtsov & Gnedin (2004), one learns that these processes become important only for densities below $\sim 10^{-28}$ and $10^{-26} \text{ g cm}^{-3}$ at $z \sim 0$ and 4, respectively. Using equation (23), we conclude that while these processes may have a certain effect on the cooling rate near the virial radius, they should be negligible once the analysis is applied inside the inner half of the halo, where the critical scale for shock heating is determined in practice. We address these effects in more detail elsewhere (Birnbom, Dekel & Loeb in preparation).

3.4 Metallicity

The metallicity near the virial radius and in the inner halo, which also enters the cooling rate, is one of our most uncertain inputs. For the mean metallicity Z (in solar units) as a function of redshift z , we use the two-parameter functional form

$$\log(Z/Z_0) = -s z, \quad (25)$$

where Z_0 is today's metallicity and the slope s governs the rate of growth.

An upper limit may be imposed by the hot, X-ray emitting intra-cluster medium (ICM) at low redshifts, which indicate $Z_0 \sim 0.2$ – 0.3 . The ICM metallicity evolution in SAMs, assuming a range of different feedback recipes, yields consistently an average enrichment rate of roughly $s \simeq 0.17$ (R. Somerville private communication; De Lucia, Kauffmann & White 2004). We adopt this enrichment rate s in our modelling below.

A realistic estimate of the metallicity near the virial radius (or perhaps a lower limit for the inner halo) may be provided by C IV absorbers in the intergalactic medium at $z \sim 2$ – 4 (Schaye et al. 2003). At densities appropriate to typical NFW haloes at $z = 3$ (with $c = 3$), namely $\rho_{\text{vir}} \simeq 53 \rho_u$, they measure an average of $[C/H] = -2.47$. Silicon measurements, on the other hand, seem to indicate a metallicity that is about five times larger (A. Aguirre private communication). If one takes the geometrical mean between the metallicities indicated by C IV and by Si one has $Z(z = 3) \simeq 0.0075$. This translates to $Z_0 = 0.025$ if $s = 0.17$.

We note that another popular indicator, Mg II, indicates consistently $Z < 0.01$ within 50–100 kpc of galaxy centres at $z \sim 1$ (private communication with J. Charlton; e.g. Ding et al. 2003).

³ In the spherical simulations, the relevant factor relating the baryon density to the DM density in equation (20) is actually closer to ~ 3 because of a 'bump' in the DM density just inside the virial radius.

The damped Ly α systems (DLAS) are believed to sample cold gas deeper inside the haloes, and can thus provide another interesting limit. Observations in the range $z = 1\text{--}4$ (Prochaska et al. 2003) can be fitted on average by equation (25) with $Z_0 \simeq 0.2$ and a somewhat steeper slope $s \simeq 0.26$. However, a fit with $s = 0.17$ (and then $Z_0 = 0.1$) is not ruled out.

Based on the above estimates, we adopt as our fiducial metallicities $Z_0 = 0.03$ at R_v and $Z_0 = 0.1$ at the ‘disc’ radius $\sim 0.1R_v$, both with an enrichment rate $s = 0.17$.

3.5 Inside the halo

For a shock in the inner halo we wish to estimate the factors f_r , f_u and f_ρ .

Empirically from the spherical simulation of BD03, for a shell encompassing a mass just shy of the critical mass (as well as from the toy model of BD03 of gas contracting in a static isothermal sphere), we estimate for $f_r \equiv r/R_v$

$$f_\rho \equiv \frac{\rho_0(r)}{\rho_0(R_v)} \simeq \begin{cases} f_r^{-1.6}, & r \lesssim R_v \\ f_r^{-2.1}, & r \sim 0.1R_v \end{cases}. \quad (26)$$

We adopt below $f_\rho = f_r^{-2}$ at $f_r = 0.1$.

Energy conservation assuming pure radial motion inside a static singular isothermal sphere yields

$$f_u \equiv \frac{u_0(r)}{u_0(R_v)} = [1 + 2f_b (f_r^{-1} - 1) + 2(1 - f_b) \ln f_r^{-1}]^{1/2}. \quad (27)$$

For $f_r = 0.1$ and $f_b = 0.05$ this gives the estimate $f_u \simeq 2.5$.

Based on the definition of f_u , the temperature behind a virial shock is related to the temperature obtained from the stability condition at radius r by

$$T(R_v) = f_u^{-2} T(r). \quad (28)$$

3.6 Crude explicit estimates

The critical temperature for stability, as obtained by comparing t_{cool} and t_{comp} in the cosmological context, equations (24) and (19), is

$$T_6 \lambda_{-22}^{-1}(T, Z) = 1.04 A^{-3/2} F, \quad (29)$$

where

$$F \equiv f_r f_u^{-1} f_\rho f_{b,0.05} f_{\bar{\rho},0.17} (1 - 3\tilde{u}_s)^{-1}. \quad (30)$$

The cooling function as computed by Sutherland & Dopita (1993) can be crudely approximated in the range $0.1 < T_6 < 10$ by

$$\Lambda_{-22} \simeq 0.12 Z_{0.03}^{0.7} T_6^{-1} + 0.02 T_6^{1/2}, \quad (31)$$

where $Z_{0.03} \equiv Z/0.03$, with Z in solar units. The above expression is valid for $-2.5 \leq \log Z \leq 0$, and at lower metallicities the value of Λ is practically the same as for $\log Z = -2.5$. This fit is good near $T_6 \sim 1$ for all values of Z . The first term refers to atomic cooling, while the second term is due to bremsstrahlung. For an approximation relevant in haloes near M_{shock} we ignore the bremsstrahlung term, which becomes noticeable only at higher temperatures. One can then obtain in equation (29) an analytic estimate for the critical temperature:

$$T_6 \simeq 0.36 A^{-3/4} (Z_{0.03}^{0.7} F)^{1/2}. \quad (32)$$

Using equations (15) and (28), with $|u_0| = V_v$ just outside the virial radius, we then obtain for the critical velocity and mass

$$V_{100} \simeq 1.62 A^{-3/8} (Z_{0.03}^{0.7} F)^{1/4} f_u^{-1} (1 + \tilde{u}_s)^{-1}, \quad (33)$$

$$M_{11} \simeq 25.9 A^{3/8} (Z_{0.03}^{0.7} F)^{3/4} f_u^{-3} (1 + \tilde{u}_s)^{-3}. \quad (34)$$

A comment regarding the \tilde{u}_s dependence of our results. The critical temperature depends on the shock velocity \tilde{u}_s via F , $T \propto (1 - 3\tilde{u}_s)^{-1/2}$, reflecting the \tilde{u}_s dependence of t_{comp} . The critical temperature is thus monotonically increasing with \tilde{u}_s . An additional \tilde{u}_s dependence enters when the temperature is translated to a critical virial velocity using the jump condition, $V \propto (1 + \tilde{u}_s)^{-1} T^{1/2}$, and then to a critical mass, $M \propto (1 + \tilde{u}_s)^{-3} T^{3/2}$. For a slowly moving shock, $\tilde{u}_s \ll 1/3$, the combined \tilde{u}_s dependence of the critical mass is $M \propto [1 + (9/4)\tilde{u}_s](1 - 3\tilde{u}_s) \simeq 1 - (3/4)\tilde{u}_s$ – a decreasing function of \tilde{u}_s . This means that at a given radius in a halo of a given mass, when everything else is equal, a slowly expanding shock is actually more stable than a shock at rest. For example, if the shock is expanding with $\tilde{u}_s = 1/7$ rather than $\tilde{u}_s = 0$, the critical mass is smaller by about 24 per cent. However, recall that stability is harder to achieve when the shock is expanding relatively fast, and the compression completely vanishes if $\tilde{u}_s \geq 1/3$ (Section 2.3.4).

For actual crude estimates of the critical scales at $z = 0$, we assume $f_{b,0.05} \simeq f_{\bar{\rho},0.17} \simeq 1$. For a shock at the virial radius, $f_r = f_u = f_\rho = 1$, we assume $Z_0 \simeq 0.03$ and $\tilde{u}_s \simeq 1/7$, and obtain

$$T_6 \simeq 0.5, \quad V_{100} \simeq 1.6, \quad M_{11} \simeq 26. \quad (35)$$

At an inner radius closer to the disc vicinity, say $f_r = 0.1$, we estimate $f_u \simeq 2.5$ and $f_\rho \simeq 100$ (Section 3.5). Assuming $Z_0 \simeq 0.1$ and $\tilde{u}_s \simeq 0$, we obtain

$$T_6 \simeq 1.1, \quad V_{100} \simeq 1.1, \quad M_{11} \simeq 8.8. \quad (36)$$

We see that for a shock at $r \sim 0.1R_v$, the expected critical mass is smaller than at R_v , somewhat below $\sim 10^{12} M_\odot$.

The above estimates are useful for exploring the qualitative dependences of the critical values on redshift, metallicity and gas fraction. For example, in equation (34), the explicit redshift dependence and the decrease of metallicity with redshift tend to lower the critical mass towards higher z . On the other hand, the decrease of halo concentration with z (i.e. increase of $f_{\bar{\rho},0.17}$), and the possible increase of the effective gas fraction with z (Section 6), tend to push the critical mass up at higher z .

3.7 More accurate estimates

We now obtain a better estimate of the critical temperature (and then critical mass and virial velocity) by solving equation (29) using the exact cooling function of SD93 and adopting specific models for the time evolution of metallicity and halo structure. The results are presented graphically.

The baryon density is computed assuming an NFW profile whose concentration evolves in time as in equation (22). The effective fraction of cold gas is assumed to be $f_b = 0.05$, motivated by best fits of SAMs to the Milky Way (Klypin, Zhao & Somerville 2002) and by fitting the Λ CDM halo mass function to the observed luminosity function near L_s (Bell et al. 2003a). The metallicity evolution is parametrized as in equation (25) with $s = 0.17$ for today’s metallicities in the range $Z_0 = 0.03\text{--}0.3$. Upper and lower estimates for the critical scales are evaluated at the virial radius and at $r = 0.1R_v$, respectively, using the crude estimates of Section 3.5. In the following figures the shock is assumed to be at rest, $u_s = 0$.

Fig. 2 shows the critical mass as a function of redshift. At a typical inner halo radius, $r = 0.1R_v$, we plot the curves for three different current metallicities: $Z_0 = 0.03, 0.1, 0.3$. The critical halo mass,

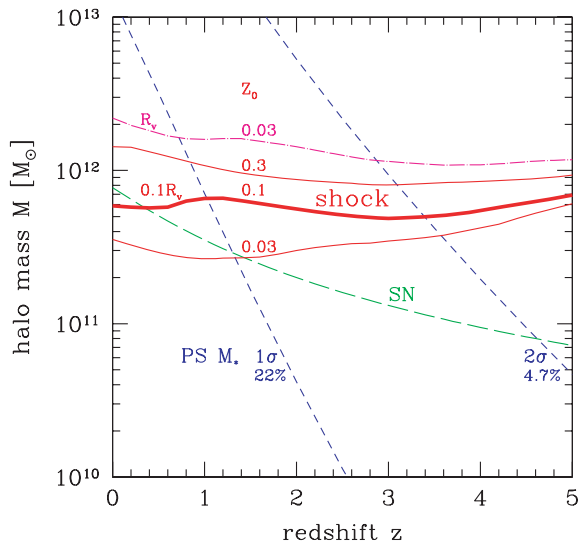


Figure 2. Critical shock-heating halo mass as a function of redshift. The three solid (red) curves refer to a shock at the inner halo, $r = 0.1 R_v$, with different metallicities as indicated; the middle curve ($Z_0 = 0.1$) is our best estimate. The dash-dotted (magenta) curve refers to a shock at the virial radius with $Z_0 = 0.03$. The other parameters used are: $f_b = 0.05$, $u_s = 0$, $s = 0.17$ (see text). Shown for comparison (short dash, blue) are the Press–Schechter estimates of the forming halo masses, corresponding to 1σ (M_*) and 2σ , where the fractions of total mass in more massive haloes are 22 and 4.7 per cent, respectively. Also shown is the critical mass for SN feedback discussed in Section 5 (long dash, green).

for $Z_0 = 0.1$, is $\simeq 6 \times 10^{11} M_\odot$ quite independent of redshift. The uncertain metallicity introduces a scatter by a factor of 2 up and down (for $z < 2.5$).

An upper limit of $\sim 2 \times 10^{12} M_\odot$ is obtained for a shock at R_v when a correspondingly low metallicity is assumed, $Z_0 = 0.03$. When the assumption of $\tilde{u}_s \simeq 0$ is replaced by $\tilde{u}_s \simeq 1/7$, allowing the shock to expand with the virial radius as seen in Fig. 1(a), the critical mass at R_v with $Z_0 = 0.03$ becomes comparable to that at $0.1 R_v$ with $Z_0 = 0.3$.

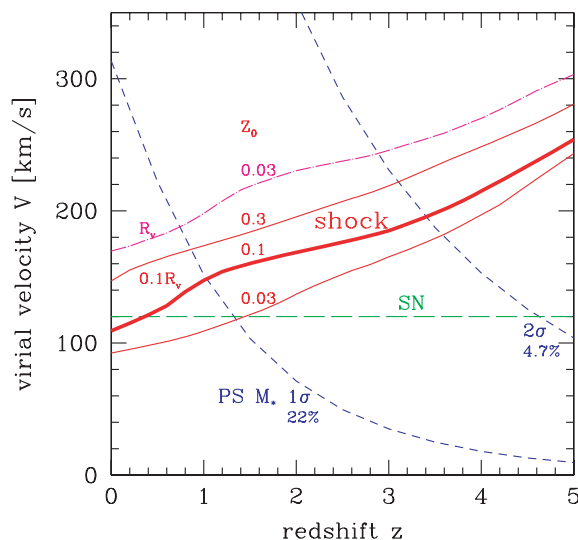


Figure 3. Same as Fig. 2 but for the corresponding halo virial velocity.

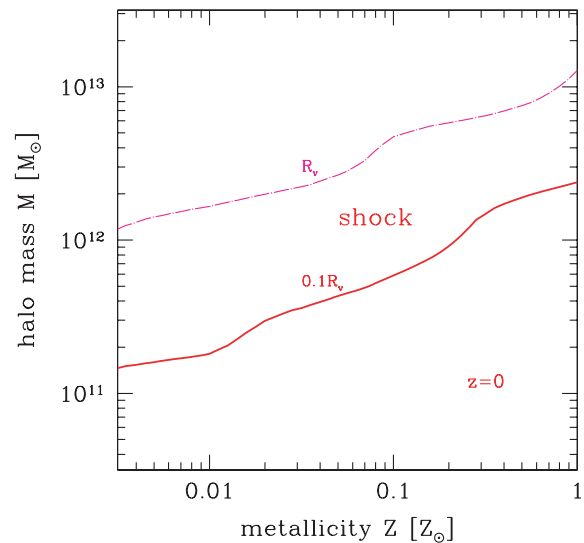


Figure 4. Critical shock-heating halo mass as a function of metallicity at redshift $z = 0$. The solid (red) curve refers to a shock at the inner halo, $r = 0.1 R_v$. The dash-dotted (magenta) curve refers to a shock at the virial radius.

Fig. 3 shows the corresponding virial velocity. At $z = 0$, the critical virial velocity for a shock in the inner halo is $\sim 120 \text{ km s}^{-1}$, with a $\pm 30 \text{ km s}^{-1}$ scatter due to metallicity. The critical virial velocity increases monotonically with redshift, to $\sim 200 \text{ km s}^{-1}$ near $z \sim 3$ (a crude fit to the redshift dependence is $V_v = 120 + 28z$).

The dependence on metallicity at $z = 0$ is highlighted in Fig. 4. The metallicity enters strongly through the cooling function $\Lambda(T, Z)$. The critical mass grows roughly like $Z^{1/2}$, as predicted in equation (34), so it spans about an order of magnitude over the whole metallicity range.

The analytic estimates of equations (33) and (34), based on the approximate cooling function, turn out to provide good estimates in most cases, and can therefore be used for extending the results analytically to any desired choice of the relevant parameters.

We learn that the critical halo mass for shock stability at the disc vicinity, $M_{\text{shock}}(r_{\text{disc}})$, is somewhat smaller than for a shock at the virial radius, $M_{\text{shock}}(R_v)$. This result is robust: it is true even if the metallicity at the virial radius is smaller by an order of magnitude than the metallicity at the disc, and even when \tilde{u}_s at R_v is as large as $1/7$. This means that as the halo is growing in mass, the conditions for a stable shock develop first in the inner halo and somewhat later in the outer halo. Thus, in haloes of mass $M < M_{\text{shock}}(r_{\text{disc}})$, we expect cold flows with no shock heating throughout the halo. In the other extreme of haloes of mass $M > M_{\text{shock}}(R_v)$, we expect shock heating of most of the gas by a shock near the virial radius. In haloes of mass in the narrow intermediate range $M_{\text{shock}}(r_{\text{disc}}) < M < M_{\text{shock}}(R_v)$, we expect shock heating somewhere inside the halo, preventing most of the gas from falling in and giving rise to a hot medium. This predicted mass range, of a factor of 2–3 in mass, is consistent with the range of transition from all cold to mostly hot seen in cosmological simulations (Section 4).

Also shown in Fig. 2 are the typical masses of haloes forming at different redshifts, the 1σ (termed M_*) and 2σ halo masses according to the Press–Schechter formalism, equation (A18). According to the improved Sheth–Tormen version, the corresponding fractions of the total mass encompassed in haloes exceeding the mass M are 22 and 4.7 per cent, respectively. One can see in Fig. 2 that M_{shock} coincides with M_* at $z \sim 1$, and with the 2σ mass at $z \sim 3.4$. By

$z \sim 2$, say, most of the forming haloes are significantly less massive than M_{shock} . When embedded in a large-scale high σ -density peak, the distribution of forming haloes at a given z may shift towards more massive haloes. In fact, the most massive halo in a volume of co-moving size ~ 100 Mpc is likely to be more massive than $10^{12} M_{\odot}$ at all relevant redshifts ($z < 6$, say). Nevertheless, the qualitative result concerning the majority of the haloes remains valid. We conclude that *in the vast majority of forming galaxies the gas has never been shock heated to the virial temperature* – it rather flows cold all the way to the disc vicinity.

We note that the values obtained for M_{shock} at low redshifts are compatible with the observed bimodality/transition scale summarized in Section 1. The estimates in the inner halo, using the lower and upper limits for Z_0 , indeed border the observed characteristic halo mass of $\sim 6 \times 10^{11} M_{\odot}$. The upper-limit estimate at R_v corresponds to a halo mass similar to that of the Milky Way.

4 COLD STREAMS IN HOT HALOES

4.1 Cosmological simulations

Cosmological hydro simulations indicate that the phenomenon of cold flows is a general phenomenon not restricted to spherical symmetry. Fig. 5 displays snapshots of an Eulerian simulation from Birnboim et al. (in preparation).⁴ Shown are maps of gas temperature in two epochs in the evolution of a protogalaxy: one at $z \simeq 4$, when the halo is already relative massive, and the other at $z \simeq 9$, when the halo is still rather small. While the more massive halo, near the critical scale, shows a hot gas component near the virial temperature behind a virial shock, the smaller halo shows only cold flows inside the virial radius.

Similar results have been obtained earlier from smoothed particle hydrodynamics (SPH) simulations by Fardal et al. (2001), who emphasized the feeding of galaxies by cold flows preferentially at early epochs. Based on our spherical analysis, we understand that this redshift dependence mostly reflects the smaller masses of haloes at higher redshifts. Kereš et al. (2005) have analysed similar SPH simulations and presented the case for the two modes of infall, cold and hot, in more detail.⁵ For example, in their figs 1 and 2 they demonstrate the two-mode phenomenon very convincingly by showing the distribution of particles and their trajectories in temperature–density diagrams. Their most informative fig. 6 shows the fractions of cold and hot infall as a function of halo mass at different redshifts. For all haloes of masses below a critical mass the infall is predominantly cold. Near the critical mass there is a relatively sharp transition into a hot mode, which becomes dominant above the critical mass. The transition from 100 per cent cold to more than 50 per cent hot occurs across a range of only a factor ~ 2 in halo mass. In this simulation, where zero metallicity is assumed, the transition mass is $M \sim 3 \times 10^{11} M_{\odot}$ at all redshifts in the range 0–3. A similar transition as a function of halo mass, and the constancy of the critical mass as a function of redshift, are both reproduced in the Eulerian cosmological simulations studies in Birnboim et al. (in preparation).

The spherical simulations described above (Section 2.1), and the corresponding analytic analysis (Sections 2 and 3), yield very sim-

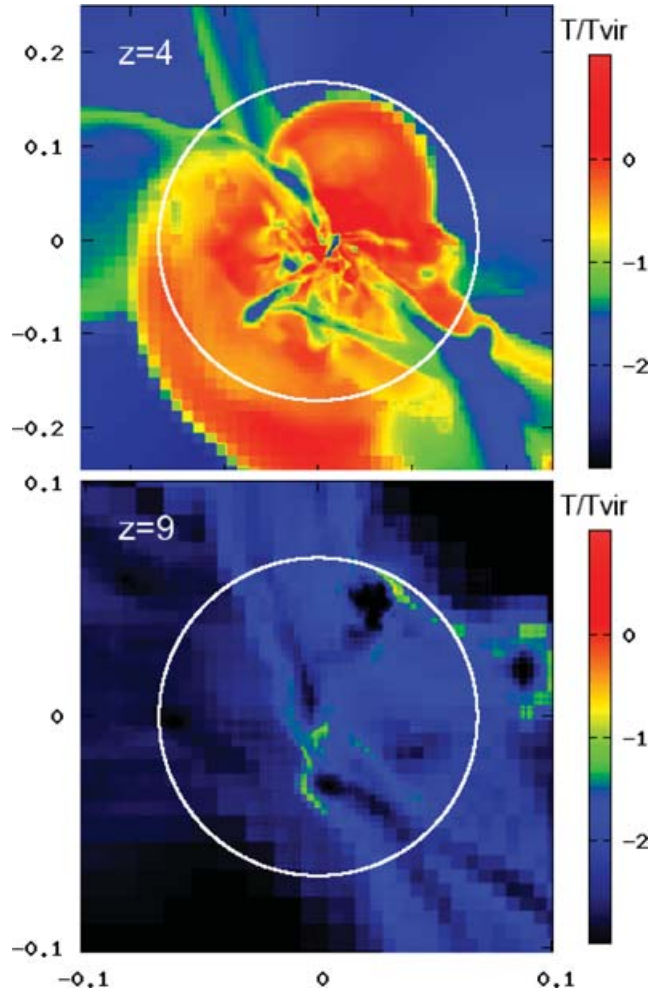


Figure 5. Snapshots from a cosmological hydrodynamical simulation (Birnboim et al. in preparation) showing the gas temperature in a slice of a protogalaxy at two different epochs, when it has two different masses. The temperature is relative to the virial temperature of the halo at that time. The side of each slice is scaled to be $3R_v$ (numbers in comoving h^{-1} Mpc). Top: At $z \simeq 4$, when the halo is already relatively massive, $M \simeq 3 \times 10^{11} M_{\odot}$. Much of the gas is heated by a strong shock near the virial radius (white circle). Cold streams penetrate through the hot medium deep into the halo. Bottom: At $z \simeq 9$, when the halo is still small, $M \simeq 2 \times 10^{10} M_{\odot}$. The gas flows in cold ($T \ll T_v$), showing no evidence for shock heating inside the virial radius (circle).

ilar results. In fact, our analytic predictions for the case of zero metallicity match the critical mass measured by Kereš et al. (2005) remarkably well.

The stability criterion derived in the spherical case turns out to be valid locally in the cosmological simulations where the non-spherical features are pronounced. Birnboim et al. (in preparation) use this criterion to identify the cold streams and hot media in the simulations without explicit information concerning the presence or absence of actual shocks. When testing the criterion in these simulations, in which the hot and cold phases may be present in the same halo, the local gas properties at each position is first transformed to post-shock quantities, as if there was a shock there, and the stability is evaluated based on the derived value of γ_{eff} there. The resultant maps of γ_{eff} resemble quite well the temperature maps of the

⁴ A description of the simulation technique can be found in Kravtsov (2003), where it was used for other purposes.

⁵ The hot phase becomes an ‘infalling’ mode in this simulation after the gas cools, but in reality it may be kept hot and be prevented from falling in by feedback effects, Section 5.

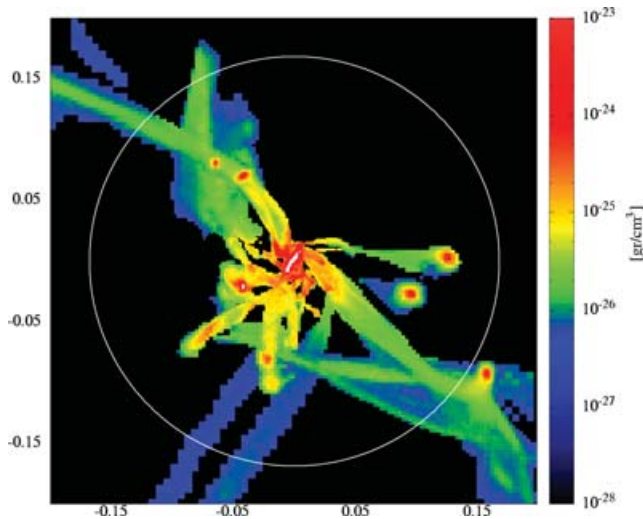


Figure 6. A snapshot from a cosmological hydrodynamical simulation (Zinger et al. in preparation) showing the gas density of the cold flows ($T < 10^{4.5}$ K, $\rho > 10^{-26.3}$ g cm $^{-3}$) within the virial radius of the same massive galaxy shown in Fig. 5 at $z = 4$. The cold phase is filamentary. In the outer radii, the gas filaments tend to ride on the large-scale DM filaments feeding the halo. About half the mass of the cold phase is in dense clumps.

actual simulation. This demonstrates that the wisdom gained by the spherical analysis is applicable in the general case.

4.2 Filaments in the simulations

The cosmological simulation of Fig. 5 also reveal that massive haloes at high redshift can have cold streams embedded in the hot medium. While the hot medium is rather spherical, the cold flows are *filamentary* and sometimes clumpy. Fig. 6 highlights the filamentary nature of the cold gas in the same big halo. This phenomenon is consistent with the findings of Kereš et al. (2005) that the cold mode may in some cases coexist with the hot mode even above M_{shock} , especially at $z > 2$ (their fig. 6) and preferentially in relatively isolated galaxies. We wish to understand the origin of this phenomenon, and learn about its dependence on cosmological time.

The simulation results of Kereš et al. (2005) provide several additional relevant clues. First, their figs 17 and 18 indicate that the cold infall indeed tends to be filamentary, especially at high redshifts, while the hot mode is more spherical. They report that the directional signal measuring filamentary infall in the cold accretion mode is stronger for haloes above M_{shock} while the infall is more isotropic below it.

Secondly, Kereš et al. (2005) display in their fig. 13 the *environment* dependence of the gas-infalling modes, showing that the cold and hot modes dominate at low and high neighbourhood densities, respectively. We learn from this figure that at $z < 2$ the cold mode dominates for galaxy densities below $n_{\text{gal}} \sim 0.3$ (h^{-1} Mpc) $^{-3}$ and becomes negligible at larger environment densities, while at $z = 3$ the cold mode is more pronounced than the hot mode for all neighbourhood densities up to $n_{\text{gal}} \sim 10$ (h^{-1} Mpc) $^{-3}$. The correlation between the environment density and the host halo mass implies that this z variation of the environment dependence could be partly attributed to the finding of a significant residual cold mode in massive haloes at high z (e.g. their fig. 6).

Thirdly, Kereš et al. (2005) show in their fig. 16 that the cold accretion is on average of *higher density* than the hot mode. This is by

only a factor of ~ 2 (probably underestimated because they mix small and large haloes), but since the shock is responsible for a density increase by a factor of 4, the actual overdensity of the hypothetical post-shock (or pre-shock) gas is more like ~ 8 . Similarly, Nagai & Kravtsov (2003) find in their simulation of a massive halo that the filamentary structure is associated with gas entropy ($\propto T/\rho^{2/3}$) *far below* that of the surrounding halo gas. In the high-resolution simulation shown in Fig. 6, we find that the density in the cold streams is actually higher than the surrounding gas density by two orders of magnitude or more. The higher gas density in the filaments is associated with a more efficient cooling which prevents a shock from forming along the filaments (Section 2).

We find that the cold gas filaments at the halo outskirts are strongly correlated with the DM filaments that feed this halo (reported in detail in Seleson & Dekel in preparation). These filaments are part of the large-scale cosmic web. They enter massive haloes at high redshift as narrow streams with a density higher than the halo average by a factor of a few. The initial overdensity of the gas flowing along the DM filaments scales with the DM density, while its inflow velocity is comparable to the halo virial velocity. As a result, the initial cooling time in the thin filaments is shorter by a factor of a few than in the surrounding spherical halo, while the compression time is comparable in the filaments and the host halo. Equation (7) then implies that the thin filaments have a harder time supporting a stable shock. The gas filaments remain cold, and become denser as the stream penetrates through the hot medium into the halo centre. The result is that in massive haloes at high redshift the critical halo mass for shock heating in the filaments feeding them is larger than the estimate for a spherical virial shock derived in Section 3. We provide below a crude estimate for this revised critical mass.

4.3 Interplay with the clustering scale

What is the reason for the appearance of cold streams in massive haloes above M_{shock} at high z ? First, recall that there is another characteristic scale in the problem – the scale of *non-linear clustering* M_* , determined by the shape and amplitude of the initial fluctuation power spectrum and its growth rate. The masses for 1σ haloes (M_*) and 2σ haloes, based on equation (A18) with $\nu = 1$ and 2, respectively, are shown again in Fig. 7. We note that $M_{\text{shock}} \sim M_*$ at $z \leq 1$, while $M_{\text{shock}} \gg M_*$ at $z > 2$. This means that $\sim 10^{12} M_{\odot}$ haloes are typical at $z < 1$ but they are the highest rare peaks at $z > 2$. We argue that this is responsible for the difference in the cold-filament behaviour of $M \gtrsim M_{\text{shock}}$ haloes in the two epochs. Since the large-scale structure of DM is roughly self-similar in time (when measured in terms of M_* and the background universal density), we can learn about the difference between typical and rare haloes by comparing $M \sim M_*$ and $M \gg M_*$ haloes in a single simulation snapshot. One can see in any high-resolution cosmological N -body simulation (e.g. the ‘Millennium Run’, visualized in www.mpa-garching.mpg.de/galform/millennium) that the rare massive haloes tend to be nodes fed by a few intersecting relatively narrow filaments which are denser than the virial density of these haloes. On the other hand, a typical M_* halo is commonly embedded in a single filament of the cosmic web, and this halo thus sees a wide-angle inflow pattern in which the matter density is comparable to the virial density [this is quantified in Seleson & Dekel (in preparation)]. This explains why $\gtrsim 10^{12} M_{\odot}$ haloes are fed by narrow dense streams at $z > 2$ but not at $z < 1$.

A crude way to estimate the maximum halo mass for cold streams at a given redshift is as follows. Recall that the critical ratio for shock

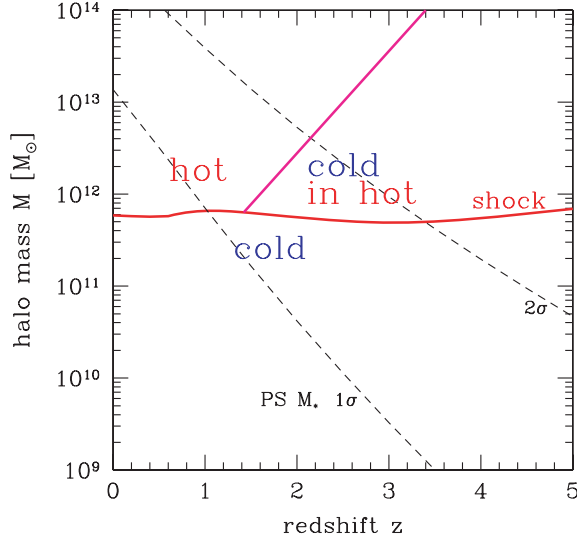


Figure 7. Cold streams and shock-heated medium as a function of halo mass and redshift. The nearly horizontal curve is the typical threshold mass for a stable shock in the spherical infall from Fig. 2, below which the flows are predominantly cold and above which a shock-heated medium is present. The inclined solid curve is the upper limit for cold streams from equation (40) with $f = 3$; this upper limit is valid at redshifts higher than $z_{\text{crit}} \sim 1-2$, defined by $M_{\text{shock}} > fM_*$. The hot medium in $M > M_{\text{shock}}$ haloes at $z > z_{\text{crit}}$ hosts cold streams which allow disc growth and star formation, while haloes of a similar mass at $z < z_{\text{crit}}$ are all hot, shutting off gas supply and star formation.

stability (equation 17) scales roughly like

$$\frac{t_{\text{cool}}}{t_{\text{comp}}} \propto \frac{\rho^{-1} T \Lambda^{-1}}{RV^{-1}}, \quad (37)$$

with T , R and V the halo virial quantities and ρ the gas density. At a given epoch, the typical halo ρ and R/V are roughly independent of halo mass (based on the definition of the virial radius), so with the virial relation $T \propto M^{2/3}$, and approximating the cooling function with $\Lambda \propto T^{-1}$ (equation 31), the critical ratio for spherical infall in virialized haloes is

$$\left(\frac{t_{\text{cool}}}{t_{\text{comp}}} \right)_{\text{halo}} = \left(\frac{M}{M_{\text{shock}}} \right)^{4/3}. \quad (38)$$

The analogous critical ratio in the dense streams inside a halo of mass M , assuming that RV^{-1} in the streams is the same as in the halo, is inversely proportional to the density enhancement $\rho_{\text{stream}}/\rho_{\text{halo}}$ (equation 37). Our estimates from N -body simulations indicate that $\rho_{\text{stream}}/\rho_{\text{vir}} \sim (fM_*/M)^{-2/3}$ with $f \sim 3$ (Seleson & Dekel in preparation). With equation (37) one obtains

$$\left(\frac{t_{\text{cool}}}{t_{\text{comp}}} \right)_{\text{stream}} = \left(\frac{fM_*}{M} \right)^{2/3} \left(\frac{M}{M_{\text{shock}}} \right)^{4/3}. \quad (39)$$

For this ratio to equal unity in the streams, the critical halo mass is

$$M_{\text{stream}} \sim \frac{M_{\text{shock}}}{fM_*} M_{\text{shock}}, \quad fM_* < M_{\text{shock}}. \quad (40)$$

This maximum mass for cold streams is shown in Fig. 7. At low z , where $fM_* > M_{\text{shock}}$, cold streams exist only for $M < M_{\text{shock}}$. At high z , where $fM_* < M_{\text{shock}}$, cold streams appear even in $M > M_{\text{shock}}$ haloes where shocks heat part of the gas, as long as $M < M_{\text{stream}}$. The critical redshift z_{crit} separating these two regimes is defined by

$$fM_*(z_{\text{crit}}) = M_{\text{shock}}. \quad (41)$$

This scenario is consistent with the cosmological hydrodynamical simulations. The shock-heating mass explains the transition from cold to hot at a given mass roughly independent of z , and the presence of cold streams above M_{shock} at $z > z_{\text{crit}}$ explains the dependence of the cold mode on redshift and environment. Besides its dependence on halo mass, the environment effect (e.g. Kereš et al. 2005) may also be due to the survivability of cold streams in different environments. While streams could survive unperturbed in relatively isolated galaxies, they are likely to be harassed by the active intergalactic environment in dense groups. The environment dependence may therefore also reflect variations in the HOD at a given halo mass. The properties of cold flows in haloes as a function of halo mass, redshift and grouping deserve a detailed analysis using high-resolution cosmological hydro simulations.

5 FEEDBACK AND LONG-TERM SHUTDOWN

Once the halo gas is shock heated in massive haloes, what is the process that keeps it hot and maintains the shutdown required by the bimodality? Is it also responsible for the rise of M/L with mass above $M_{s,\text{crit}}$ (and the absence of cooling flows in clusters)? Several feedback mechanisms can heat the gas. We suggest that they have a minimum effectiveness in haloes $\sim M_{\text{shock}}$. This can be largely induced by the shock heating itself, and in turn it can amplify the bimodality features. Some of the feedback mechanisms are limited to smaller haloes, while others, such as AGN feedback, are likely to be important in more massive haloes. The latter can be triggered by the shock heating and then help maintaining the gas hot for a long time.

5.1 Below the shock-heating scale

(i) **Supernova feedback.** Based on the physics of supernova (SN) remnants, the energy fed to the gas in haloes of $T \sim 10^5$ K is proportional to the stellar mass despite significant radiative losses (Dekel & Silk 1986). When compared to the energy required for significantly heating the gas, one obtains a maximum halo virial velocity for SN feedback, $V_{\text{SN}} \simeq 120 \text{ km s}^{-1}$. This is only weakly dependent on the gas fraction, density or metallicity (Dekel & Silk 1986, equation 49), and is therefore insensitive to redshift. Only in potential wells shallower than V_{SN} can the SN feedback significantly suppress further star formation and regulate the process. Fig. 3 shows V_{SN} and Fig. 2 shows the corresponding mass versus redshift. With an effective $f_b \sim 0.05$, the corresponding stellar mass at $z = 0$ is $\sim 3.5 \times 10^{10} M_{\odot}$, practically coinciding with the bimodality scale. The similarity of the SN and shock-heating scales is partly a coincidence, because the nuclear origin of the initial SN energy has little to do with galactic cooling or dynamics. However, there is an obvious similarity in the cooling processes and in the asymptotic behaviour of an SN remnant, which is not a strong function of its initial energy. The distinct correlations between the properties of galaxies below $M_{s,\text{crit}}$ indeed point at SN feedback as its primary driver. These correlations define a ‘fundamental line’, $V \propto M_s^{0.2}$, $Z \propto M_s^{0.4}$, $\mu \propto M_s^{0.6}$, where μ is surface brightness (Kauffmann et al. 2003b; Tremonti et al. 2004). SN feedback can explain the origin of the fundamental line (Dekel & Woo 2003) based on (1) the above energy criterion, which implies $M_s/M \propto V^2$; (2) the virial relations (equation A6); (3) the instantaneous recycling approximation, $Z \propto M_s/M_{\text{gas}}$ and (4) angular momentum conservation, $R_* \propto \lambda R$, with λ a constant spin parameter (Fall & Efstathiou 1980).

(ii) **UV-on-dust feedback.** Also working below M_{shock} are momentum-driven winds due to radiation pressure on dust grains,

arising from the continuum absorption and scattering of ultraviolet (UV) photons emitted by starbursts or AGNs Murray, Quataert & Thompson (2005). The dust survives and can provide sufficient optical depth if the gas is cold and dense, e.g. in the cold flows below M_{shock} , which can also provide the starbursts responsible for the required UV flux and metals. Since dust grains cannot survive above $\sim 10^6$ K, M_{shock} imposes an upper bound for this feedback.

(iii) **Photoionization feedback.** The UV radiation from stars and AGNs ionizes most of the gas after $z \sim 10$ (Bullock, Kravtsov & Weinberg 2000; Loeb & Barkana 2001), heats it to $\gtrsim 10^4$ K and prevents it from falling into haloes below the Jeans scale $V_v \sim 30 \text{ km s}^{-1}$ (Gnedin 2000). As the ionization persists for cosmological epochs, the hot gas evaporates via steady winds from haloes smaller than a similar scale (Shaviv & Dekel 2003). While this is important in dwarf galaxies, it cannot be very relevant to the bimodality at M_{shock} .

5.2 Above the shock-heating scale

(i) **AGN feedback.** The fact that AGNs exist preferentially in haloes above M_{crit} may be due to a lower limit for haloes hosting black holes (Koushiappas, Bullock & Dekel 2004), or starvation of AGNs by SN feedback in haloes below M_{crit} , or to another reason. The power emitted from AGNs, e.g. in their radio jets, seems to be more than necessary for keeping the gas hot. Given a black hole mass $M_{\text{BH}} \sim 10^7 M_{\odot} V_{100}^4$ in a galaxy of velocity dispersion V , and assuming that a fraction ϵ of this mass is radiated away, the ratio of energies is $E_{\text{AGN}}/E_{\text{gas}} \sim 7 \times 10^3 \epsilon f_{b,0.05}^{-1} V_{100}^{-1}$. For $\epsilon > 10^{-3}$, there seems to be enough AGN energy for affecting most of the halo gas. If this energy is released during relatively quiet, long phases of self-regulated AGN activity, it can keep the gas hot. However, black hole physics does not seem to explain the characteristic scale of bimodality. Furthermore, had the energy ratio been a measure of feedback effectiveness, it would have implied a decline with mass, in conflict with the trend of M/L .

The shock heating of the gas into a dilute medium is likely to make it vulnerable to heating and pushing by the central energy source, thus providing the *trigger* for effective AGN feedback. Simulations of winds in a two-phase medium demonstrate that the dilute phase is pushed away while the dense clouds are hardly affected (Slyz et al. 2005). This behaviour is likely to be generic, though the mechanism by which the energy released near the black hole is spread in the halo gas is an open issue (Ruszkowski, Bruggen & Begelman 2004; Scannapieco & Oh 2004). If so, the feedback efficiency may be driven by the relative fraction of hot gas rather than the actual AGN energy. Fig. 6 of Keres et al. (2004) shows that the hot fraction varies roughly $\propto M^{1/2}$, implying $M/L \propto M^{1/2}$ near and above M_{shock} , in qualitative agreement with the observed trend (Section 1, item k). In this scenario, the shutdown scale arises naturally from the shock heating.

(ii) **Two-phase medium.** Given that the cooling function peaks near $\gtrsim 10^4$ K, the virialized gas at $\gtrsim 10^6$ K develops a two-phase medium, with cold, dense clouds pressure confined within the hot, dilute medium (Field 1965; Fall & Rees 1985). The cloud sizes and evolution are affected by thermal conductivity and dynamical processes (Voigt & Fabian 2004). This can help explaining the bright-end truncation of the blue sequence (Maller & Bullock 2004). Some of the gas is locked in the orbiting clouds and the density ρ_{hot} of the hot gas is reduced, slowing the cooling and the infall. Approximating equation (17) with $t_{\text{cool}}/t_{\text{comp}} \propto \rho^{-1} T_v^2$, and recalling that $T_v \propto M_v^{2/3}$, the longer cooling time makes the critical mass for further shock heating ($t_{\text{cool}}/t_{\text{comp}} \sim 1$) smaller by a similar factor, $M_{\text{shock}} \propto$

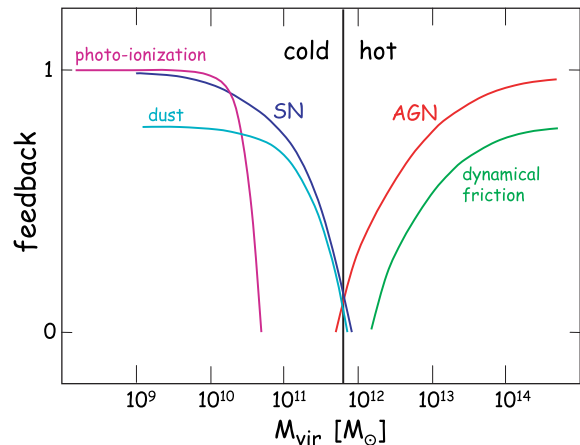


Figure 8. The ‘strength’ of the various feedback processes at $z = 0$, e.g. referring to the fraction of the initial gas that has been heated or removed (schematic). Different feedback processes are effective below and above the critical shock-heating scale $\lesssim 10^{12} M_{\odot}$, and the feedback efficiency is at a minimum near this scale, giving rise to a minimum in M/L there.

$\rho_{\text{hot}}^{3/4}$, namely it enhances the shock stability. The gas may be kept hot over longer periods by repeating shocks due to continuous accretion into the halo, which may alleviate the need for AGN feedback. Still, a necessary condition for hot gas is the initial shock heating, i.e. being in a halo above M_{shock} .

(iii) **Dynamical-friction feedback.** Another heating source is the dynamical friction acting on galaxies in a halo core. The energy transferred is comparable to that required for preventing cooling flows in cluster centres (El-Zant, Kim & Kamionkowski 2004b). The gas response to dynamical friction, unlike the DM response, has a sharp peak near a Mach number of unity (Ostriker 1999, fig. 3), namely when the gas is heated to near the virial temperature in $M > M_{\text{shock}}$ haloes and not in smaller haloes hosting cooler gas. As groups occur above a critical halo mass that roughly coincides with M_{shock} at $z \leq 1$, this feedback appears almost simultaneously with the hot medium, which then serves as the vulnerable victim of the same feedback (as for AGN feedback).

Fig. 8 is an illustration of the strength of the different feedback processes, referring to the gas fraction that could have been heated at $z \sim 0$. The figure highlights the fact that different processes dominate below and above $M_{\text{shock}} \lesssim 10^{12} M_{\odot}$. The transition from cold to hot infall has a crucial role in determining the feedback efficiencies near the critical mass; it induces a minimum in feedback efficiency at a critical scale $M_{\text{fdbk}} \sim M_{\text{shock}}$ and drives the shapes of the curves about this minimum. At higher redshifts this minimum becomes wider and deeper but centred on a similar critical mass.

6 THE ORIGIN OF BIMODALITY

6.1 A scenario from the assumed ingredients

We propose that the cold flows and shock heating play a key role in producing the observed bimodality features. These features are emphasized by the similarity between the scales associated with shock heating, feedback and clustering. Based on our current understanding of these physical processes, we assume the validity of the following.

(i) **A new mode of star formation.** The collisions of the (partly clumpy) cold streams with each other and with the inner disc are

assumed to produce starbursts, analogous to collisions of cold gaseous discs or clouds. These collisions are expected to produce isothermal shocks, behind which the rapid cooling generates dense, cold slabs where the Jeans mass is small. While the details are yet to be worked out, we assume that such a mode of star formation may be responsible for much of the stars in the universe. In some cases, this can be an enhanced quiescent mode, leaving the disc intact without producing a big spheroid, in other cases it may resemble the starbursts associated with mergers.

(ii) **Hot forever.** Once the gas in a massive halo is shock heated to near the virial temperature, it is assumed to be hot forever. This is based on the slow cooling time of the dilute hot medium and its vulnerability to AGN feedback, while cold, dense clouds and streams could be better shielded against winds and ionizing radiation. The shock heating is thus assumed to *trigger* a shutdown of all modes of star formation in haloes where cold streams do not prevail.

(iii) **Cold streams in a hot medium.** Cold streams in haloes above M_{shock} (Section 4) are assumed to supply cold gas for further disc growth and star formation, preferentially before $z_{\text{crit}} \sim 2$ and in isolated galaxies. After z_{crit} , especially in groups, cold streams are suppressed and a complete shutdown of star formation is assumed to follow.

The *bimodality* in colour (or stellar age or SFR) versus mass, the correlation with the environment and the evolution with redshift, all emerge naturally from the early efficient star formation followed after $z \sim 2$ by the *abrupt shutdown* in haloes above M_{shock} , which typically host groups. This is illustrated in Fig. 9.

(i) **The blue sequence.** It is dominated by galaxies in haloes below M_{shock} , as they grow by accretion/mergers. Cold flows lead to early disc growth and star formation, which is regulated by SN feedback over cosmological times. Galaxies can get very blue because of repeating starbursts due to the clumpy gas supply and the interplay between infall, starburst and outflow.

(ii) **Bright blue extension.** Some galaxies continue to be fed by cold streams even when they are more massive than M_{shock} , extending the blue sequence beyond L_* . This occurs especially at $z > z_{\text{crit}} \sim 2$, when the streams feeding high- σ haloes are relatively narrow and dense, resulting in massive starbursts.

(iii) **The red sequence.** Once a halo is more massive than M_{shock} , halo gas is shock heated; it becomes dilute and vulnerable to AGN feedback. At $z < z_{\text{crit}}$, where cold streams are suppressed, gas supply from the host halo shuts off, preventing any further growth of discs and star formation. If residual cold gas has been consumed in earlier mergers, there is a total shutdown of all modes of star formation above M_{shock} , allowing the stellar population to passively turn ‘red and dead’ into the red sequence. The massive tip of the blue sequence at $z > z_{\text{crit}}$ becomes the massive tip of the red sequence at $z < z_{\text{crit}}$, explaining the bright-end truncation of the blue sequence at low z , and the appearance of brighter, very red galaxies already at $z \sim 1$. Subsequent growth along the red sequence is induced by gas-poor mergers uncontaminated by new blue stars.

(iv) **Faint red extension.** When colour is plotted against stellar mass, the bimodality extends to smaller galaxies which are typically satellites of the central galaxies in common haloes. In haloes below M_{shock} , accretion on to satellite galaxies can keep them on the blue sequence for a while. In haloes above M_{shock} , where gas supply stops and the environment density is high, the satellites too become red and dead.

In an associated paper (Cattaneo et al. 2006, hereafter C05) we have implemented the proposed new physics in a hybrid SAM/N-

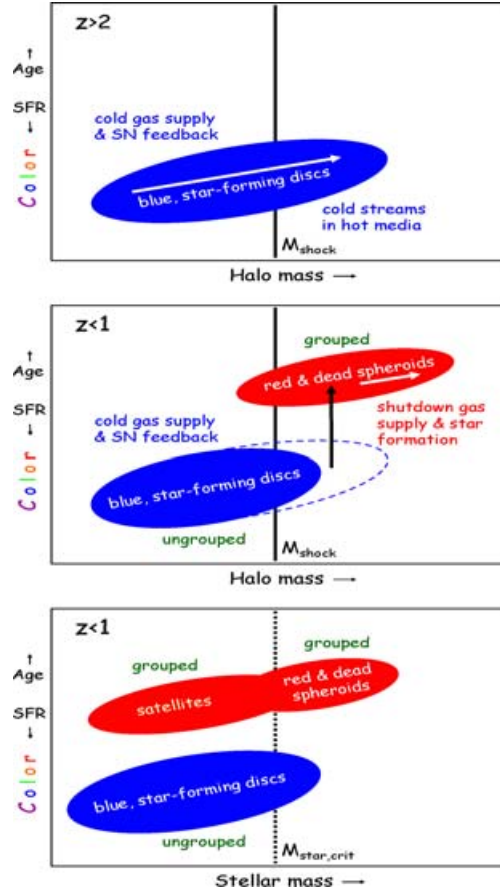


Figure 9. Schematic illustration of the origin of the bimodality in colour (or SFR or stellar age) as a function of halo mass. At $z > 2$ (top), continuous gas supply, regulated by SN feedback, makes the galaxies evolve along the blue sequence, which extends beyond the shock-heating mass due to cold streams in hot media. At $z < 1$ (middle), in the absence of cold streams above M_{shock} , the shock-heated gas is kept hot by AGN feedback, gas supply and star formation shut down, and the stellar population passively turns red and dead. Gas-poor mergers stretch the red sequence towards larger masses. When the halo mass is replaced by stellar mass (bottom), the red sequence is stretched towards small stellar masses due to satellite galaxies sharing a common halo. The colour is correlated with the environment density via the halo mass, with the minimum group mass being comparable to M_{shock} .

body simulation (GalICS; Hatton et al. 2003). The main new ingredient is a shutdown in gas cooling and star formation above $M_{\text{shock}} \sim 10^{12} M_{\odot}$ after $z_{\text{crit}} \sim 3$, while allowing for efficient star formation by cold streams in more massive haloes prior to z_{crit} . The revision yields excellent fits to the observed features at low and high redshifts (C05). Fig. 10 shows one example of the results – colour–magnitude diagrams which demonstrate the success of this model along the lines envisioned in Fig. 9. The top panel highlights the main deficiencies of the ‘standard’ model at $z = 0$: an excess of bright blue galaxies accompanied by a shortage of red-enough galaxies compared to the SDSS data (Baldry et al. 2004). The ‘new’ model puts the blue and red sequences where they should be at $z = 0$, with a proper truncation at the bright-blue end and appropriately red colours in the red sequence. The $z = 3$ diagram shows the predicted bright blue galaxies (which were absent in the ‘standard’ model, C05). The colours distinguish between galaxies in haloes below (blue) and above $2 \times 10^{12} M_{\odot}$, comparable to the mass separating haloes hosting field galaxies and groups. This indicates that

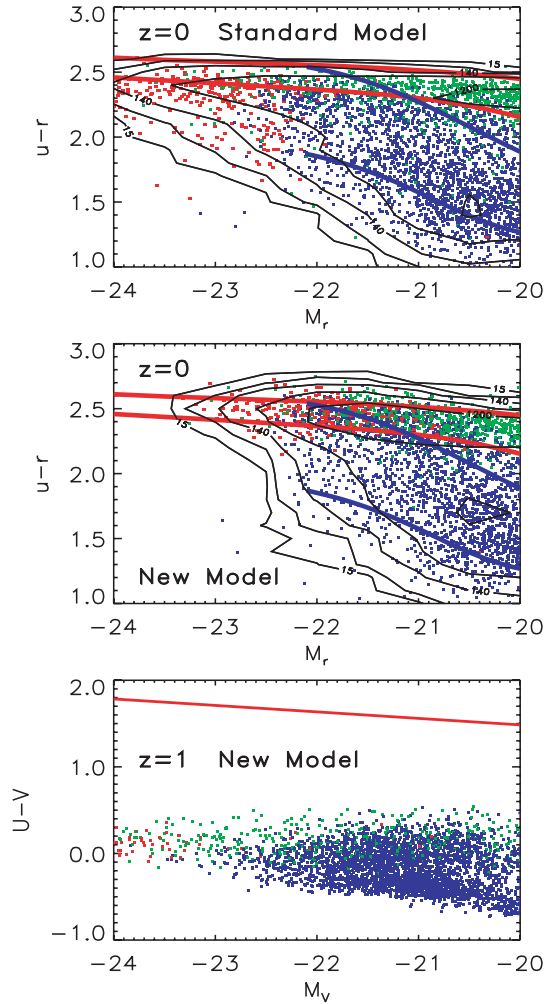


Figure 10. Colour–magnitude diagrams from a hybrid semi-analytic/ N -body simulation (C05) demonstrating the success of the proposed scenario in reproducing the main bimodality features along the lines of Fig. 9. Top: Based on the ‘standard’ version of the GalICS SAM. Middle: The result of incorporating a shutdown in star formation above $M_{\text{shock}} \sim 10^{12} M_{\odot}$ after $z_{\text{crit}} \sim 3$, and allowing for star formation even in $>M_{\text{shock}}$ haloes prior to z_{crit} . Bottom: Same at $z = 3$. Blue dots refer to galaxies in haloes below M_{shock} , while red and green dots are galaxies in haloes above M_{shock} , central and satellite galaxies, respectively. The halo mass correlates with number density in the environment. Contours mark number density of galaxies in the diagram. The colour curves mark the main bodies of the blue and red sequences in the SDSS data (Baldry et al. 2004).

the model recovers the correlation between colour and environment density. It also recovers the bimodality in bulge-to-disc ratio (C05).

6.2 Origin of observed features

The colour–magnitude bimodality emerges from the colour–mass bimodality lying at the basis of the model. The red satellites extend the luminosity range over which the blue and red sequences coexist and highlight the *gap* in colour. This gap is amplified because the galaxies making the transition into the red sequence once their haloes become $>M_{\text{shock}}$ tend to be the merger remnants with big bulges and AGNs from the *red tip* of the blue sequence.

The correlated bimodality in *bulge-to-disc ratio* can be understood in the scenario where most of the big spheroids in the red

sequence are passively aged galaxies that have grown massive stellar components already in the blue sequence (C05). The transition to the red sequence is likely to be made by galaxies with big spheroids because (1) they have consumed their gas in the same mergers that produced the spheroids, (2) these mergers tend to occur in big haloes hosting groups where shock heating stops the gas supply and (3) these spheroids contain massive black holes that can keep the gas hot.

The strong anticorrelation of *SFR* (and blue colour) with the number density in the *environment* is a natural outcome of the mass dependence. The key is that the minimum halo mass for groups at $\sim 10^{12} M_{\odot}$ is comparable to M_{shock} at $z \leq 1$. The strong dependence of cold gas supply on host halo mass can therefore be responsible for the distinction between the *SFR* in field and clustered galaxies (which is thus predicted to be limited to late times). The galaxies dominating low-HOD haloes below M_{shock} enjoy cold gas supply leading to discs forming stars. The galaxies populating groups of subhaloes, typically in haloes above M_{shock} at $z \leq 1$, suffered starvation of external gas supply and lost their internal gas in mergers, thus stopped forming stars and passively evolved to the red sequence. The faint end of the red sequence, which is preferentially present in high environment densities, is due to the starvation of satellite galaxies in the high-HOD haloes above M_{shock} . The model naturally predicts the secondary bimodality seen along the red sequence (along the lines of Berlind et al. 2005).

The classic *morphology–environment* relation, traditionally attributed to the correlation of merger rate with the environment, may also be viewed as a result of the cold-flow phenomenon. New discs are predicted to form only in haloes below M_{shock} , namely in ‘field’ galaxies, and not in the group haloes above M_{shock} . On the other hand, the frequent major mergers in groups help build the big spheroids preferentially there. The abrupt shutdown of star formation above M_{shock} makes the colour–environment correlation stronger than the morphology–environment correlation.

The environment dependence highlights an interesting cross-talk between the *clustering* and gas processes. The distinction between haloes hosting a single dominant galaxy and haloes hosting groups has traditionally been attributed to gas cooling on a dynamical time-scale (Rees & Ostriker 1977). Our shock-stability analysis helps quantifying this idea. However, it seems from N -body simulations that the gravitational DM clustering process has a parallel role (Kravtsov et al. 2004). The HOD of subhaloes develops a transition from single to multiple occupancy near a comparable halo mass, associated with the current M_* . The HOD of galaxies is similar to the HOD of galaxies deduced from the observed correlation function for galaxies (Section 1, item ix). One can conclude that in haloes above $\sim 10^{12} M_{\odot}$ the potential wells associated with the DM subhaloes provide the sites for the fragmented gas collapse on the scales preferred by cooling, thus emphasizing this scale as the minimum scale for groups. The coincidence between these scales is behind the environment dependence of the bimodality features. We have discussed above the other possible role of M_* in the appearance of cold filaments in hot haloes at high z (Section 4) and in some of the feedback processes (Section 5).

The minimum in mean halo *mass-to-light* ratio M/L near $\sim 10^{11-12} M_{\odot}$ can be attributed to the maximum in gas supply near M_{shock} dictated by the shock heating above this mass and the associated minimum in feedback there (Fig. 8). While SN feedback gets stronger towards smaller haloes, AGN feedback pumps energy more effectively into the shock-heated medium in more massive haloes. The ‘fundamental line’ due to SN feedback below M_{shock} corresponds to $M/L \propto M^{-2/3}$ (Dekel & Woo 2003), while the transition from cold to hot infall indicates $M/L \propto M^{1/2}$ above

M_{shock} , in the ballpark of the findings from 2dF (Section 1, item k). This settles the discrepancy between the halo mass function and the galaxy luminosity function both below and above the bimodality scale. Our model predicts a similar behaviour of $M/L(M)$ at $z \sim 1$, to be revealed by spectroscopic surveys (such as the DEEP2/3 Madgwick et al. 2003a; Coil et al. 2004), which have already confirmed the predicted invariance of the bimodality mass in this redshift range (Fig. 7).

The massive starbursts at high z , primarily due to cold streams, are helped by an increase in M_{shock} itself. At high z , the upper bound for SN feedback becomes $\ll M_{\text{shock}}$ (Fig. 2), allowing for a more efficient gas accretion at $\lesssim M_{\text{shock}}$. If f_b is 0.13 instead of 0.05 at $z \sim 3$, the value of M_{shock} doubles (equation 34). Another factor proportional to f_b enters when translating from halo to stellar mass, yielding a total increase of ~ 5 in the critical stellar mass at high z . We thus expect a strong star formation activity at high z in galaxies with stellar masses exceeding $\sim 10^{11} M_{\odot}$.

The global *star formation history* could be derived from the predicted SFR as a function of mass and redshift, convolved with the time evolution of the halo mass function in the given cosmology. With the prediction that haloes $\lesssim M_{\text{shock}}$ are the most efficient star formers, and with the Press–Schechter estimate that haloes of such a mass typically form at $z \sim 1$ (Fig. 2), the star formation density is predicted to peak near $z \sim 1$, with a relatively flat behaviour towards higher redshifts and a sharp drop towards lower redshifts, as observed. In particular, the cumulative stellar density seems to stop growing quite abruptly at $z \sim 1$ (Dickinson et al. 2003), when the typical forming haloes become larger than M_{shock} and the cold flows are suppressed. The growth of the hot fraction as a function of halo mass (Kereš et al. 2005) can be translated to the drop in SFR after $z \sim 1$. The ‘downsizing’ of star formation in galaxies is also helped by the shutdown of star formation above M_{shock} while smaller galaxies can make stars also later, until they fall into bigger haloes.

7 DISCUSSION: OTHER IMPLICATIONS

Possible implications on open issues in galaxy formation where further study is desired are worth mentioning.

X-rays. The shock-heating scale may be detectable in soft X-rays, as the lower limit for galaxies and groups containing hot halo gas. In turn, the suppression of heating below M_{shock} explains the missing soft X-ray background (Pen 1999; Benson et al. 2000). We predict a noticeable suppression of X-ray emission in the range $5 \times 10^5 - 2 \times 10^6$ K.

Ly α emission. The cold ($\sim 10^4$ K) flows may instead be an efficient source of Ly α radiation, possibly associated with observed Ly α emitters at high redshift (Kurk et al. 2003). It has been argued based on SPH simulations (Fardal et al. 2001; Furlanetto et al. 2003) that the flows radiate their infalling energy mostly in Ly α before they blend smoothly into the discs. Our Eulerian simulations (Zinger et al. in preparation) indicate that the streams do eventually shock in the inner halo. This produces X-ray, but, given the high density there, the X-ray radiation is likely to be confined to an ionized Strömgren sphere of a few kiloparsecs. This energy eventually transforms into Ly α radiation, which could propagate out via thermal broadening and systematic redshifts. A study involving radiative transfer is required.

Damped Ly α systems. The possible association of massive cold flows with DLAS (Prochaska et al. 2003) should be addressed in cosmological simulations.

LIRGs. Cold flows may explain the massive starbursts associated with LIRGs at $z \lesssim 1$ (Hammer et al. 2005). If half the stars in today’s

discs were formed in such LIRGs, and the majority of galaxies today are fragile discs, than many of the LIRGs could not have been produced by violent major mergers. The cold streams may provide a less violent starburst mechanism not associated with the destruction of discs. Simulations that properly incorporate cold streams should be confronted with these data.

Angular momentum. The proposed scenario may set the stage for solving the angular momentum puzzle – the overproduction of low angular momentum spheroids in current cosmological simulations (Navarro & Steinmetz 2000). The solution should involve the removal of baryons with low angular momentum. In small galaxies, SN feedback can blow the gas away from their small building blocks, which are otherwise the main source of low angular momentum via many minor mergers (Maller & Dekel 2002; Maller, Dekel & Somerville 2002). In galaxies near M_{shock} , we find from cosmological simulations (Zinger et al. in preparation; see Fig. 6) that the low angular momentum gas is typically associated with the shock-heated medium, which can be prevented from cooling by AGN feedback. The cold streams come from larger distances with ~ 50 per cent higher specific angular momentum, appropriate for producing extended discs. The feedback effects, both below and above the critical scale, have not been properly simulated yet because of incomplete treatment of the microphysics.

Cold clouds. The formation of discs by a clumpy cold gas phase may have the following implications. (1) It may explain the starbursts responsible for very blue galaxies. (2) It may help explaining the bright-end truncation of the luminosity function (Section 5). (3) The dynamical friction bringing the clouds into the disc transfers energy into the halo, which may help explaining the discrepancy between the steep cusps predicted by N -body simulations and the flat cores indicated by rotation curves in low-surface-brightness galaxies (Dekel et al. 2003a; Dekel, Devor & Hetzroni 2003b; El-Zant et al. 2004a; Ma & Boylan-Kolchin 2004). (4) The same process may lower the predicted maximum rotation velocity in discs at a given luminosity, balancing the adiabatic contraction of the dark halo, and thus repair the zero-point offset in models of the Tully–Fisher relation (e.g. Klypin et al. 2002; Abadi et al. 2003; Dutton et al. 2005). (5) This may explain the lack of anticorrelation between the residuals in velocity and radius at a given luminosity (Courteau & Rix 1999), indicating comparable contributions of disc and dark halo to the gravitational potential at the effective disc radius (Dutton et al., in preparation).

Dust lane. Edge-on discs above $V_v \simeq 120 \text{ km s}^{-1}$ show a well-defined dust lane, while less massive discs show diffuse dust above and below the disc (Dalcanton, Yoachim & Bernstein 2004). A thick, turbulent, dusty gas phase is indeed expected when SN feedback is effective, and when cold streams shock and produce stars, both predicted below an appropriate scale.

Shock heating in dwarf haloes. The cold/hot infall and feedback processes are expected to give rise to two scales characterizing dwarf galaxies. The lower bound at $V_v \sim 10\text{--}15 \text{ km s}^{-1}$ (e.g. Dekel & Woo 2003, Fig. 3) is commonly attributed to the drop in atomic cooling rate below $\sim 10^4$ K. We propose that this involves shock heating, in analogy to M_{shock} discussed above. Fig. 11 (bottom) shows the quantity relevant to shock stability, $t_{\text{cool}}/t_{\text{comp}}$ (equation 17), as a function of halo mass, now stretched to low masses. The molecular hydrogen cooling rate, relevant below 10^4 K, is weaker and may actually be eliminated after $z \sim 10$ due to molecule dissociation by the UV background (Haiman, Rees & Loeb 1996). The stability is evaluated at $z = 0$ both in the disc vicinity ($r = 0.1R_v$, $Z = 0.1$, $\tilde{u}_s = 0$, and near the virial radius ($Z = 0.03$, $\tilde{u}_s = 1/7$).

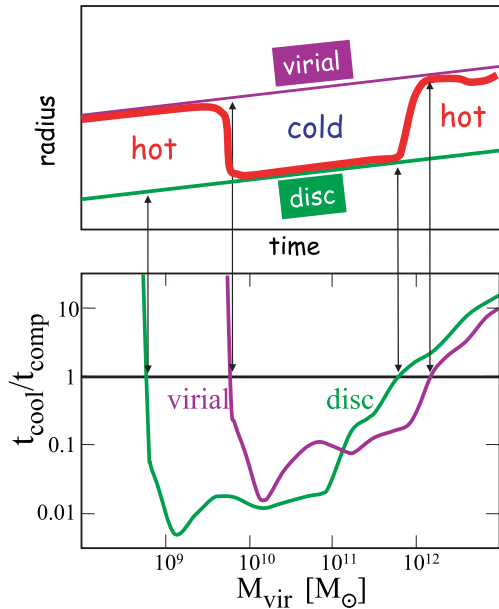


Figure 11. Shock stability as a function of mass at $z = 0$. Bottom: The ratio of rates, $t_{\text{cool}}/t_{\text{comp}}$, as a function of halo mass derived at the disc ($Z = 0.1$, $u_s = 0$) and at the virial radius ($Z = 0.03$, $u_s = 0.15$). The cooling rate is assumed to vanish for $T < 10^4$ K. A stable shock is possible once $t_{\text{cool}}/t_{\text{comp}} > 1$. Top: An illustration of the evolution of shock radius between the disc vicinity and the virial radius as the halo mass grows in time, based on the stability criterion shown in the bottom panel.

The top panel illustrates the deduced evolution of shock radius between the disc and the virial radius as the halo mass grows in time. Haloes below $\sim 10^9 M_\odot$ have stable virial shocks. As the halo grows to $\sim 6 \times 10^8 M_\odot$, the conditions near the disc become unfavourable for a stable shock, but the virial shock persists. Only when the halo becomes $\sim 6 \times 10^9 M_\odot$ shock heating stops completely. Thus, the shock heating prevents star formation in smaller haloes even before the epoch of re-ionization. Once the first stars form in the dense centres of big clouds, their UV radiation re-ionizes the gas, dissociates molecules and helps preventing cooling. This provides a natural explanation for the absence of luminous galaxies with haloes below $\sim 10^9 M_\odot$, predicting a large population of mini haloes with $V_v < 10$ km s $^{-1}$ which are completely dark. There is only a narrow window, $6 \times 10^9 - 6 \times 10^{11} M_\odot$, for haloes that allow cold flows and can form luminous discs at low redshifts. This range is somewhat broader at high redshifts, with the lower bound dropping below $10^9 M_\odot$ and the upper bound rising above $10^{12} M_\odot$ prior to $z \sim 3$. Within this window, it seems that below $V_v \sim 30$ km s $^{-1}$ most of the haloes are totally dark⁶ and the others are populated by gas-poor dwarf spheroidals. This is the scale predicted by photoionization feedback (Section 5).

8 CONCLUSION

8.1 Summary of results

The classic argument of cooling on a dynamical time-scale (Rees & Ostriker 1977; White & Rees 1978), with order-of-magnitude

⁶ This can be deduced from the discrepancy between the flat faint-end luminosity function and the steep halo mass function predicted in Λ CDM, given the Tully–Fisher like velocity–luminosity relation of dwarf galaxies (e.g. Dekel & Woo 2003, fig. 3).

estimates of the time-scales involved, provided an inspiring qualitative upper bound for luminous galaxies, at a halo mass of $M \sim 10^{12-13} M_\odot$. An analytic study of the actual shock-heating process (Birnboim & Dekel 2003, and this paper) now yields a more concrete halo critical scale at $M \simeq 6 \times 10^{11} M_\odot$, somewhat smaller than the original estimate. The criterion for critical shock stability,

$$t_{\text{cool}}^{-1} = t_{\text{comp}}^{-1}, \quad (42)$$

is a balance between the cooling rate and the post-shock compression rate, which restores the pressure supporting the shock against gravitational collapse. The compression time is somewhat larger than the crossing time at the shock position. The absolute magnitudes of these time-scales are irrelevant – they could in principle both be longer than the Hubble time, because what matters for shock heating versus cold flows is only the relative rates of the competing processes. The most relevant critical scale is obtained in the inner halo, because as the halo grows, the shock first becomes stable in the inner halo, and it then propagates outwards to the virial radius. Haloes of mass below the threshold mass build discs in their centres by cold flows, while in haloes above the threshold much of the gas is shock heated. These results are confirmed by spherical hydrodynamical simulations. The same phenomenon is seen at a comparable scale in general cosmological hydrodynamical simulations. They reveal that in haloes near the critical scale, and even in larger haloes preferentially at $z \geq 2$ and especially in field galaxies, cold streams along the filaments feeding the galaxy penetrate through the hot medium, and allow further disc growth and star formation.

The interplay between these cold flows and shock heating, the gravitational clustering scale, and the different feedback processes acting below and above a similar mass scale, is proposed to be responsible for the robust bimodality imprinted on the observed galaxy properties. Cold flows in haloes much bigger than the clustering scale allow massive starbursts at $z \geq 2$, while shock heating in comparable haloes at later times shuts off star formation and leads to big red galaxies. While SN and radiative feedbacks regulate star formation below the critical scale, the presence of dilute, shock-heated gas in more massive haloes allows the AGN feedback (or another energetic source capable of affecting big galaxies) to keep the shock-heated gas hot and prevent further disc growth and star formation. The observed bimodality and many of the related phenomena are argued to arise naturally from such a scenario (Section 6). The shock-heating process also plays a role in introducing a lower bound to haloes hosting galaxies, at $\sim 10^9 M_\odot$. The mass range where disc galaxies can form today turns out to be quite narrow, between a few times $10^9 M_\odot$ to slightly below $10^{12} M_\odot$.

8.2 Re-engineering of SAMs

Once the new physical processes are properly incorporated in the detailed models of galaxy formation, they solve many of the apparent conflicts between theory and observation. At a crude level, one might have naively thought that since the cooling time is anyway assumed to be shorter in smaller haloes, the details of the cold flows and shock heating would not matter much to the final result. However, a closer inspection reveals that there are several key features which make a qualitative difference as stated below.

(i) **Star formation.** The near-supersonic cold streams provide a new efficient mechanism for early star formation. This is in contrast to the gradual infall of cooling shock-heated gas assumed in most SAMs, which starts from near rest, accretes smoothly into the disc, and joins the quiescent mode of star formation there. We find that

the expected cold-gas supply is significantly more efficient than assumed in most current models even in small haloes (Cattaneo et al. in preparation).

(ii) **Heating inside out.** The concept of an expanding ‘cooling radius’ used in current SAMs is limited to massive haloes where a virial shock exists. Otherwise, it is the shock causing the heating which is propagating from the inside out.

(iii) **Shutdown of star formation.** The combination of shock heating and AGN feedback provides a mechanism for shutting off disc growth and star formation above a threshold halo mass.

(iv) **Cold streams.** Cold streams that penetrate through the hot media continue to make discs and produce stars in haloes above M_{shock} . This happens mostly at $z \geq 2$, and preferentially in less grouped galaxies, allowing big blue galaxies mostly at high z and some at low-density environments, and enforcing a sharp shutdown of star formation at late times and especially in clustered galaxies.

A practical schematic recipe for the critical halo mass below which cold streams prevail and above which one may assume a shutdown of gas supply and star formation is

$$M_{\text{crit}} = \begin{cases} M_{\text{shock}}, & z < z_{\text{crit}} \\ M_{\text{shock}} \left(\frac{M_{\text{shock}}}{f M_*(z)} \right), & z > z_{\text{crit}} \end{cases} \quad (43)$$

where the critical redshift z_{crit} is defined by $f M_*(z_{\text{crit}}) = M_{\text{shock}}$, the clustering scale $M_*(z)$ is given by equation (A18), and f is a numerical factor of order of a few. Our best estimates for the parameters are $\log M_* \simeq 11.8$ (but possibly another value in the range 11.3–12.3) and $f \simeq 3$ (possible range: 1–10). Using the approximation $\log M_* \simeq 13.1 - 1.3z$ ($z \leq 2$), we obtain $z_{\text{crit}} \simeq 1.4$ for $f = 3$. This recipe should allow big blue systems at $z \geq 2$, eliminate big blue systems and make big red galaxies at $z \leq 1$, and generate a bimodality near M_{shock} . This scheme can be refined to allow for a smooth transition above the critical scale by applying the shutdown to a varying fraction of the gas and by breaking the streams into clumps which will generate high peaks of starbursts.

In addition, one may wish to have an effective minimum requirement for the central black hole mass in order to ensure enough feedback energy for maintaining the gas hot. This may emphasize the bimodality gap in colour and bulge-to-disc ratio. However, the proposed shutdown by halo mass may be enough for ensuring sufficient bulge mass and black hole mass.

The SAMs should be re-engineered to incorporate these processes and thus help working out the detailed implications of the proposed scenario. Preliminary attempts to do that, using two different SAMs, indicate that the incorporation of the new proposed processes outlined above indeed leads to significantly better fits with the observed bimodality features along the lines proposed in Section 6 (C05).

8.3 Open issues

In parallel, the physics of the involved ingredients should be properly worked out in more detail, starting with the following two hypotheses that were laid out in Section 6.

(i) **Fate of cold streams.** A detailed investigation is required of the way by which the cold streams evolve and eventually merge with the central disc, the associated star formation and the resulting feedback process. While progress can be made using toy models and simplified simulations, a proper analysis will require simulations of higher resolution than are currently available. In particular, whether or not the predicted starbursts could be associated with the big dusty sources indicating massive star formation at high red-

shifts, such as the SCUBA sources (Chapman et al. 2003), remains to be determined once the theory is worked out and the observed characteristics of these sources are clarified.

(ii) **AGN feedback.** The physics of AGN feedback is another unknown. One wishes to understand how the available energy originating near the central black hole is transferred to the hot gas spread over the halo. The physics of how thermal conductivity may heat the gas is also to be investigated. The increased efficiency of these feedback mechanisms in the presence of a hot medium as opposed to their effect on cold flows and clumps are to be quantified.

Parallel attempts to work out the details of the physical input and to incorporate it in quantitative models of galaxy formation will lead to progress in our understanding of the galaxy bimodality and the associated features.

ACKNOWLEDGMENTS

We thank our collaborators A. Cattaneo, S. M. Faber, A. Kravtsov, E. Neistein, J. Primack, P. Seleson, R. Somerville and E. Zinger. We thank R. Dave, N. Katz, D. Keres and D. Weinberg for sharing with us the results of their simulations. We acknowledge stimulating discussions with J. Binney, D. Lin, G. Kauffmann, G. Mamon, J. P. Ostriker and D. Weinberg. This research has been supported by ISF 213/02 and NASA ATP NAG5-8218. AD acknowledges support from a Miller Visiting Professorship at UC Berkeley, a Visiting Professorship at UC Santa Cruz and a Blaise Pascal International Chair by Ecole Normale Supérieure at the Institut d’Astrophysique, Paris.

REFERENCES

- Abadi M. G., Navarro J. F., Steinmetz M., Eke V. R., 2003, *ApJ*, 591, 499
 Abazajian K. et al., 2005, *ApJ*, 625, 613
 Baldry I. K., Glazebrook K., Brinkmann J., Ivezić Ž., Lupton R. H., Nichol R. C., Szalay A. S., 2004, *ApJ*, 600, 681
 Balogh M. et al., 2004, *MNRAS*, 348, 1355
 Bardeen J. M., Bond J. R., Kaiser N., Szalay A. S., 1986, *ApJ*, 304, 15
 Bell E. F., Baugh C. M., Cole S., Frenk C. S., Lacey C. G., 2003a, *MNRAS*, 343, 343
 Bell E. F., McIntosh D. H., Katz N., Weinberg M. D., 2003b, *ApJ*, 585, L117
 Bell E. F., McIntosh D. H., Katz N., Weinberg M. D., 2003c, *ApJS*, 149, 289
 Bell E. F. et al., 2004, *ApJ*, 608, 752
 Benson A. J., Bower R. G., Frenk C. S., White S. D. M., 2000, *MNRAS*, 314, 557
 Benson A. J., Bower R. G., Frenk C. S., Lacey C. G., Baugh C. M., Cole S., 2003, *ApJ*, 599, 38
 Berlind A. A., Blanton M. R., Hogg D. W., Weinberg D. H., Dav R., Eisenstein D. J., Katz N., 2005, *ApJ*, 629, 625
 Binney J., 1977, *ApJ*, 215, 483
 Binney J., 2004, *MNRAS*, 347, 1093
 Birnboim Y., Dekel A., 2003, *MNRAS*, 345, 349 (BD03)
 Blanton M. R., Eisenstein D., Hogg D. W., Zehavi I., 2004, preprint (astro-ph/0411037)
 Blanton M. R., Eisenstein D., Hogg D. W., Schlegel D. J., Brinkmann J., 2005, *ApJ*, 629, 143
 Blumenthal G. R., Faber S. M., Primack J. R., Rees M. J., 1984, *Nat*, 311, 517
 Bryan G. L., Norman M. L., 1998, *ApJ*, 495, 80
 Bullock J. S., Kravtsov A. V., Weinberg D. H., 2000, *ApJ*, 539, 517
 Bullock J. S., Kolatt T. S., Sigad Y., Somerville R. S., Kravtsov A. V., Klypin A. A., Primack J. R., Dekel A., 2001, *MNRAS*, 321, 559
 Carroll S. M., Press W. H., Turner E. L., 1992, *ARA&A*, 30, 499
 Cattaneo A., Dekel A., Devriendt J., Guiderdoni B., Blaizot J., 2006, *MNRAS*, submitted (astro-ph/0601295) (C05)

- Chapman S. C., Blain A. W., Ivison R. J., Smail I. R., 2003, *Nat*, 422, 695
Chapman S. C., Smail I., Blain A. W., Ivison R. J., 2004, *ApJ*, 614, 671
Ciotti L., Pellegrini S., Renzini A., D'Ercole A., 1991, *ApJ*, 376, 380
Coil A. L. et al., 2004, *ApJ*, 609, 525
Courteau S., Rix H., 1999, *ApJ*, 513, 561
Dalcanton J. J., Yoachim P., Bernstein R. A., 2004, *ApJ*, 608, 189
De Lucia G., Kauffmann G., White S. D. M., 2004, *MNRAS*, 349, 1101
Dekel A., 1981, *A&A*, 101, 79
Dekel A., Silk J., 1986, *ApJ*, 303, 39
Dekel A., Woo J., 2003, *MNRAS*, 344, 1131
Dekel A., Arad I., Devor J., Birnboim Y., 2003a, *ApJ*, 588, 680
Dekel A., Devor J., Metzroni G., 2003b, *MNRAS*, 341, 326
Dickinson M., Papovich C., Ferguson H. C., Budavári T., 2003, *ApJ*, 587, 25
Ding J., Charlton J. C., Bond N. A., Zonak S. G., Churchill C. W., 2003, *ApJ*, 587, 551
Dressler A., 1980, *ApJ*, 236, 351
Dutton A., van den Bosch F. C., Courteau S., Dekel A., 2005, preprint (astro-ph/0501256)
El-Zant A. A., Hoffman Y., Primack J., Combes F., Shlosman I., 2004a, *ApJ*, 607, L75
El-Zant A. A., Kim W., Kamionkowski M., 2004b, *MNRAS*, 354, 169
Fall S. M., Efstathiou G., 1980, *MNRAS*, 193, 189
Fall S. M., Rees M. J., 1985, *ApJ*, 298, 18
Fardal M. A., Katz N., Gardner J. P., Hernquist L., Weinberg D. H., Davé R., 2001, *ApJ*, 562, 605
Ferguson H. C., Babul A., 1998, *MNRAS*, 296, 585
Field G. B., 1965, *ApJ*, 142, 531
Fioc M., Rocca-Volmerange B., 1999, *A&A*, 344, 393
Furlanetto S. R., Schaye J., Springel V., Hernquist L., 2003, *ApJ*, 599, L1
Giavalisco M. et al., 2004, *ApJ*, 600, L103
Gnedin N. Y., 2000, *ApJ*, 542, 535
Haiman Z., Rees M. J., Loeb A., 1996, *ApJ*, 467, 522
Hammer F., Flores H., Elbaz D., Zheng X. Z., Liang Y. C., Cesarsky C., 2005, *A&A*, 430, 115
Hartwick F. D. A., 2004, *ApJ*, 603, 108
Hatton S., Devriendt J. E. G., Ninin S., Bouchet F. R., Guiderdoni B., Vibert D., 2003, *MNRAS*, 343, 75
Heavens A., Panter B., Jimenez R., Dunlop J., 2004, *Nat*, 428, 625
Helsdon S. F., Ponman T. J., 2003, *MNRAS*, 340, 485
Hogg D. W. et al., 2003, *ApJ*, 585, L5
Kannappan S. J., 2004, *ApJ*, 611, L89
Kauffmann G. et al., 2003a, *MNRAS*, 346, 1055
Kauffmann G. et al., 2003b, *MNRAS*, 341, 54
Kauffmann G., White S. D. M., Heckman T. M., Menard B., Brinchmann J., Charlot S., Tremonti C., Brinkmann J., 2004, *MNRAS*, 353, 713
Kereš D., Katz N., Weinberg D. H., Dave R., 2005, *MNRAS*, 363, 2
Klypin A., Zhao H., Somerville R. S., 2002, *ApJ*, 573, 597
Koushiappas S. M., Bullock J. S., Dekel A., 2004, *MNRAS*, 354, 292
Kravtsov A. V., 2003, *ApJ*, 590, L1
Kravtsov A. V., Berlind A. A., Wechsler R. H., Klypin A. A., Gottloeber S., Allgood B., Primack J. R., 2004, *ApJ*, 609, 35
Kravtsov A. V., Gnedin O. Y., 2005, *ApJ*, 623, 650
Kurk J., Röttgering H., Pentericci L., Miley G., Overzier R., 2003, *New Astron. Rev.*, 47, 339
Lahav O., Rees M. J., Lilje P. B., Primack J. R., 1991, *MNRAS*, 251, 128
Loeb A., Barkana R., 2001, *ARA&A*, 39, 19
Ma C., Boylan-Kolchin M., 2004, *Phys. Rev. Lett.*, 93, 021301
Madau P., Ferguson H. C., Dickinson M. E., Giavalisco M., Steidel C. C., Fruchter A., 1996, *MNRAS*, 283, 1388
Madgwick D. S. et al., 2003a, *ApJ*, 599, 997
Madgwick D. S., Somerville R., Lahav O., Ellis R., 2003b, *MNRAS*, 343, 871
Maller A. H., Bullock J. S., 2004, *MNRAS*, 355, 694
Maller A. H., Dekel A., 2002, *MNRAS*, 335, 487
Maller A. H., Dekel A., Somerville R., 2002, *MNRAS*, 329, 423
Marinoni C., Hudson M. J., 2002, *ApJ*, 569, 101
Mathews W. G., Brighenti F., 2003, *ARA&A*, 41, 191
Mo H. J., White S. D. M., 2002, *MNRAS*, 336, 112
Mo H. J., Mao S., White S. D. M., 1998, *MNRAS*, 295, 319
Moustakas L. A. et al., 2004, *ApJ*, 600, L131
Murray N., Quataert E., Thompson T. A., 2005, *ApJ*, 618, 569
Nagai D., Kravtsov A. V., 2003, *ApJ*, 587, 514
Navarro J. F., Steinmetz M., 2000, *ApJ*, 538, 477
Navarro J. F., Frenk C. S., White S. D. M., 1997, *ApJ*, 490, 493
Nilsson K., Fynbo J. P. U., Moller P., Sommer-Larsen J., Ledoux C., 2005, *A&A Lett.*, submitted (astro-ph/0512396)
Osmond J. P. F., Ponman T. J., 2004, *MNRAS*, 350, 1511
Ostriker E. C., 1999, *ApJ*, 513, 252
Pen U., 1999, *ApJ*, 510, L1
Prochaska J. X., Gawiser E., Wolfe A. M., Castro S., Djorgovski S. G., 2003, *ApJ*, 595, L9
Rees M. J., Ostriker J. P., 1977, *MNRAS*, 179, 541
Ruszkowski M., Bruggen M., Begelman M. C., 2004, *ApJ*, 615, 675
Scannapieco E., Oh S. P., 2004, *ApJ*, 608, 62
Schaye J., Aguirre A., Kim T., Theuns T., Rauch M., Sargent W. L. W., 2003, *ApJ*, 596, 768
Shapley A. E., Erb D. K., Pettini M., Steidel C. C., Adelberger K. L., 2004, *ApJ*, 612, 108
Shaviv N. J., Dekel A., 2003, preprint (astro-ph/0305527)
Sheth R. K., Tormen G., 2002, *MNRAS*, 329, 61
Silk J., 1977, *ApJ*, 211, 638
Slyz A., Devriendt J., Bryan G., Silk J., 2005, *MNRAS*, 356, 737
Smail I., Ivison R. J., Blain A. W., Kneib J.-P., 2002, *MNRAS*, 331, 495
Sutherland R. S., Dopita M. A., 1993, *ApJS*, 88, 253
Thomas D., Maraston C., Bender R., de Oliveira C. M., 2005, *ApJ*, 621, 673
Thoul A. A., Weinberg D. H., 1995, *ApJ*, 442, 480
Tremaine S. et al., 2002, *ApJ*, 574, 740
Tremonti C. A. et al., 2004, *ApJ*, 613, 898
Voigt L. M., Fabian A. C., 2004, *MNRAS*, 347, 1130
White S. D. M., Frenk C. S., 1991, *ApJ*, 379, 52
White S. D. M., Rees M. J., 1978, *MNRAS*, 183, 341
Yan R., Madgwick D. S., White M., 2003, *ApJ*, 598, 848
Yang X., Mo H. J., van den Bosch F. C., 2003, *MNRAS*, 339, 1057

APPENDIX A: USEFUL RELATIONS

We summarize here the cosmological relations used in the analysis of Section 3. This is rather basic material, based for example on Lahav et al. (1991); Carroll, Press & Turner (1992) and Mo & White (2002). By specifying it here in a concise and convenient form, we hope to allow the reader to reproduce our results and use them in future work. Additional relations associated with the spherical top-hat collapse model are brought in the appendix of BD03.

A1 Cosmology

The basic parameters characterizing a flat cosmological model in the matter era are the current values of the mean mass density parameter Ω_m and the Hubble constant H_0 . At the time associated with expansion factor $a = 1/(1+z)$, the vacuum-energy density parameter is $\Omega_\Lambda(a) = 1 - \Omega_m(a)$ and

$$\Omega_m(a) = \frac{\Omega_m a^{-3}}{\Omega_\Lambda + \Omega_m a^{-3}}. \quad (\text{A1})$$

The Hubble constant is

$$H(a) = H_0 (\Omega_\Lambda + \Omega_m a^{-3})^{1/2}, \quad (\text{A2})$$

and the age of the universe is

$$t_{\text{univ}}(a) = \frac{2}{3} H(a)^{-1} \frac{\sinh^{-1}(|1 - \Omega_m(a)|/\Omega_m(a))^{1/2}}{(|1 - \Omega_m(a)|)^{1/2}}. \quad (\text{A3})$$

The mean mass density is

$$\rho_u \simeq 1.88 \times 10^{-29} \Omega h^2 a^{-3} \simeq 2.76 \times 10^{-30} \Omega_{m0.3} h_{0.7}^2 a^{-3}, \quad (\text{A4})$$

where $\Omega_{m0.3} \equiv \Omega_m/0.3$, $h \equiv H_0/100 \text{ km s}^{-1} \text{ Mpc}^{-1}$, and $h_{0.7} \equiv h/0.7$.

A2 Virial relations

The virial relations between halo mass, velocity and radius,

$$V_v^2 = \frac{GM_v}{R_v}, \quad \frac{M_v}{(4\pi/3)R_v^3} = \Delta\rho_u \quad (\text{A5})$$

become

$$M_{11} \simeq 6.06 V_{100}^3 A^{3/2} \simeq 342 R_1^3 A^{-3}, \quad (\text{A6})$$

where $M_{11} \equiv M_v/10^{11} M_\odot$, $V_{100} \equiv V_v/100 \text{ km s}^{-1}$, $R_1 \equiv R_v/1 \text{ Mpc}$, and

$$A \equiv (\Delta_{200} \Omega_{m0.3} h_{0.7}^2)^{-1/3} a. \quad (\text{A7})$$

An approximation for $\Delta(a)$ in a flat universe (Bryan & Norman 1998) is

$$\Delta(a) \simeq (18\pi^2 - 82\Omega_\Lambda(a) - 39\Omega_\Lambda(a)^2)/\Omega_m(a). \quad (\text{A8})$$

The virial temperature can be defined by

$$\frac{kT_v}{m} = \frac{1}{2} V_v^2. \quad (\text{A9})$$

For an isotropic, isothermal sphere, this equals σ^2 , where σ is the one-dimensional velocity dispersion and the internal energy per unit mass is $e = (3/2) \sigma^2$. Thus

$$V_{100}^2 \simeq 2.79 T_6 \quad M_{11} \simeq 28.2 T_6^{3/2} A^{3/2}, \quad (\text{A10})$$

where $T_6 \equiv T_v/10^6 \text{ K}$.

A3 Press Schechter

Linear fluctuation growth is given by (Lahav et al. 1991; Carroll et al. 1992; Mo & White 2002)

$$D(a) = \frac{g(a)}{g(1)} a, \quad (\text{A11})$$

where

$$g(a) \simeq \frac{5}{2} \Omega_m(a) \times \left[\Omega_m(a)^{4/7} - \Omega_\Lambda(a) + \frac{1 + \Omega_m(a)/2}{1 + \Omega_\Lambda(a)/70} \right]^{-1}. \quad (\text{A12})$$

The CDM power spectrum is approximated by (Bardeen et al. 1986)

$$P(k) \propto k T^2(k), \quad (\text{A13})$$

with

$$T(k) = \frac{\ln(1 + 2.34q)}{2.34q} \times [1 + 3.89q + (16.1q)^2 + (5.46q)^3 + (6.71q)^4]^{-1/4}, \quad (\text{A14})$$

where

$$q = \frac{k}{(\Omega_m h^2 \text{ Mpc}^{-1})}. \quad (\text{A15})$$

It is normalized by σ_8 at $R = 8 h^{-1} \text{ Mpc}$, where

$$\sigma^2(R) = \frac{1}{2\pi} \int_0^\infty dk k^2 P(k) \tilde{W}^2(kR), \quad (\text{A16})$$

and with the Fourier transform of the top-hat window function

$$\tilde{W}(x) = 3(\sin x - x \cos x)/x^3. \quad (\text{A17})$$

In the Press Schechter (PS) approximation, the characteristic halo mass $M_*(a)$ is defined as the mass of the 1σ fluctuation,

$$1 = \nu(M, a) = \frac{\delta_c}{D(a) \sigma(M)}, \quad \delta_c \simeq 1.69, \quad (\text{A18})$$

where M and the comoving radius R are related via the universal density today: $M = (4\pi/3) \bar{\rho}_0 R^3$. The mass of 2σ fluctuations is obtained by setting $\nu(M, a) = 2$, etc. Based on the improved formalism of Sheth & Tormen (2002), the fraction of total mass in haloes of masses exceeding M is

$$F(> M, a) \simeq 0.4 \left(1 + \frac{0.4}{\nu^{0.4}} \right) \text{erfc} \left(\frac{0.85 \nu}{\sqrt{2}} \right). \quad (\text{A19})$$

This fraction for 1 , 2 and 3σ fluctuations is 22, 4.7 and 0.54 per cent, respectively.

Fig. 2 shows the PS mass M_* as a function of redshift. For the standard Λ CDM with $\sigma_8 = 0.9$, its value at $z = 0$ is $M_{*0} = 1.36 \times 10^{13} M_\odot$. One can see that an excellent practical fit in the range $0 \leq z \leq 2$ is provided by a power law in this semilog plot: $\log M_* \approx 13.134 - 1.3z$. At larger redshifts this gradually becomes an underestimate. Trying to provide crude power-law approximations, we find that $M_* \propto a^{4.2} \propto t^{3.5}$ are crude approximations in the range $0 \leq z \leq 1$, and that $M_* \propto a^5 \propto t^4$ are good to within a factor of 2 in the range $0 \leq z \leq 2$. These power laws become overestimates at higher redshifts.

This paper has been typeset from a $\text{\TeX}/\text{\LaTeX}$ file prepared by the author.