# Galaxies and Halos in the Sloan Digital Sky Survey

Timothy A. McKay[*] and the SDSS Collaboration[†]

[*]*Department of Physics, University of Michigan, 500 East University, Ann Arbor, MI, 48109*
[†]*www.sdss.org*

**Abstract.** Structure formation theory provides very effective predictions of the properties of dark matter halos, including their mass function, clustering, and internal structure. Observations of structure, however, rely on luminous galaxies as tracers. A detailed understanding of the way galaxies occupy dark matter halos is essential for connecting structure formation theory to observation. We describe some of the observables available for contraining the halo occupancy, illustrating each using data from the Sloan Digital Sky Survey. Many of these observables can now be measured with great statistical precision. Comparison of these observables to theory is now limited by systematic uncertainty in the relationship between observable quantities (like velocity dispersion vs. cluster richness) and theoretically favored quantities (like $M_{200}$). We argue for the use of carefully crafted simulations in making this connection, and illustrate their use in some example analyses.

## PRELIMINARIES

We now have a reasonably well established framework for understanding the formation of structure in the universe. At the time of recombination we begin with a nearly uniform fluid composed of weakly interacting cold dark matter (perhaps 90%) and ordinary baryons (perhaps 10%). Very small density fluctuations are imprinted on this fluid, perhaps by quantum fluctuations generated in the Big Bang. These fluctuations are then amplified by the presence of gravitationally interacting matter. This amplification is dominated by dark matter.

As the Universe evolves dark matter halos begin to separate from the overall Hubble flow. These dark matter halos, in one way of accounting, *are* the structure. If we knew the mass, location, and internal structure of all the halos, we would have a complete description of the matter distribution. Identifying these dark halos, and probing their properties, is a major goal.

Within dark matter halos, a radiative differentiator acts. While the dark matter is thought to interact only gravitationally, the minority component of baryonic material can radiate away its energy: cooling and sinking to the center of the potential well until it is rotationally supported. Within this baryonic core stars form, evolve, and feed back material and entropy into the surrounding halo, forming galaxies. It is these galaxies, and not the halos, which we observe. In this way of thinking, they are labels, lovely markers hinting at the presence of much more extensive dark matter halos.

# Connecting Galaxies with Halos

To connect our observable, the locations and motions of luminous galaxies, to theoretical predictions, mostly of the properties of dark matter halos (for example [1]), we need to understand the relationship between the luminous properties of galaxies and their dark matter environments.

To first order, this relationship can be described as bias. In one form the bias is defined as the ratio of the galaxy-galaxy and mass-mass correlation functions. On large scales, greater than a few Mpc, bias is observed to be simple [2], as expected. More complex behaviour is generically expected on smaller scales. A kind of conspiracy of effects converts the complex mass-mass correlation function to the simple unbroken power law observed in the galaxy-galaxy correlation function [3].

This is just the first order picture, and a more complete portrait, including the dependence of this bias on galaxy and halo properties, contains substantial information about galaxy formation, and perhaps about the nature of dark matter.

# The Halo-Occupancy Distribution Function

Both analytic theory and N-body numerical experiments can predict the properties of dark matter halos, including their mass function, their clustering properties, and their internal structure. Comparing this theory to observation requires an understanding of how galaxies 'occupy' these halos.

One useful way to describe the relationship between halos and the galaxies which populate them relies on the 'halo occupancy distribution function' P(N|M). This function describes the probability that a halo of mass M will host a total of N galaxies. This function, together with some information about spatial and velocity bias between the dark matter and galaxies within each halo, provides a useful way of thinking about bias [4]. It provides an essentially complete description of the relationship between the distribution of galaxies we observe, and the distribution of dark matter we are so eager to learn. Given the details of this HOD, it is possible to calculate essentially any observable of large- scale structure.

The approach outlined here is pretty sketchy. For example, if we're going to count galaxies, what will count? Is this a luminosity limited sample? In what band? How should we identify the dark matter halos, how should we define their masses? In the end, this initial HOD description will have to be expanded to accomodate the differing occupancy of different types of galaxies.

The tools available for constraining the HOD range from the N-body simulations of dark matter to observations of real galaxies. The simulations can inform our interpretations of observations at an important level, helping us to understand the effects of projection and to determine our selection functions. The observations obviously provide feedback to the simulations, providing direct HOD constraints, measuring various scaling relations on halo scales, and determining the dependence of these things on galaxy properties.

The process of constraining the HOD will have to include an iterative dance among

all these elements, revising models and making new observations until all the pieces fit together.

## IDENTIFYING HALOS AND MEASURING MASSES

To generate constraints on the HOD, we must first identify halos. We describe here measurements based on two approaches. In the first, we identify halos with individual galaxies, labeled with their properties including luminosity, type, and environment. This is surely a very first order method. While at low halo mass there is a roughly one-to-one correspondence between bright galaxies and halos, this is not true at higher mass, where individual halos host groups and clusters of galaxies.

A second way to identify halos is to use group and cluster finders, which label the halos with galaxy content. This method is closer to the proper spirit of the HOD, but has the drawback of being sensitive only to relatively massive halos.

The details of these halo identifications need to be understood. The efficiency and purity of their selection has to be known as a function of redshift, and its dependence on the selection of tracer galaxies must be known. It's also very important to understand how the centers for these halos are identified, and how their spatial extent is constrained. Analysis of simulated universes can provide important insight here.

Once halos are identified, we need to measure their masses. We have two basic probes of mass. In dynamical measurements of mass, the observables are the positions and velocities for a set of luminous test particles. These probe the dynamical effect of gravity on the test particles, which sample the velocity field around the halo. In lensing measurements of mass, the observable is a shear field and the geometry associated with it. Lensing probes the space-time geometry around the objects of interest. It provides a probe of the projected galaxy-mass correlation function around halos.

Note that neither measures mass very directly. Inferring masses from these measurements in either case requires careful consideration of a variety of effects. Of particular importance is our choice of probes. Since the relationship between galaxy properties and mass is strong, our choice of test particles, and of lens and source galaxies affects strongly what we observe. All of that rich behavior needs to be understood. In both cases, measurements of these halo mass probes can now be made with high signal-to-noise. Essentially all the work from now on lies in accurately understanding how these precisely determined observables relate to theoretically well determined quantities, like halo mass function and structure.

## SDSS HALO STUDIES

The data which enable the studies we describe come from the Sloan Digital Sky Survey. The SDSS is a large collaboration, involving perhaps 200 scientists at a number of institutions. It is designed to make comprehensive astronomical observations. Over the coming few years the SDSS will complete an imaging survey of $10^4$ square degrees of the sky. Images are obtained in 5 colors for about $10^8$ galaxies. In addition to imaging,

the SDSS will measure high quality spectra for about $10^6$ galaxies and $10^5$ quasars. This set of observations will support a very broad range of science goals, in much the same sense that data collected by large high energy physics experiments is useful for many purposes. We expect SDSS data to be an important tool for astrophysics for several decades. For the analyses described here, both the imaging and spectroscopic data play an important role.

## Halos identified by individual galaxies

We begin with the lensing studies, measuring the correlation between locations of foreground galaxies and distortion in the shapes of distant sources. The SDSS data used for this study are drawn from the Early Data Release [5]. They include imaging and spectroscopic data for about 4% of the SDSS survey region. From these data, we select a sample of 34,693 foreground 'lens' objects. Every one of these objects has a spectroscopic redshift and highly accurate 5 color photometry. For this purpose they are drawn entirely from the SDSS 'main' galaxy sample. We also select a fainter background sample of 3,615,718 'source' objects.
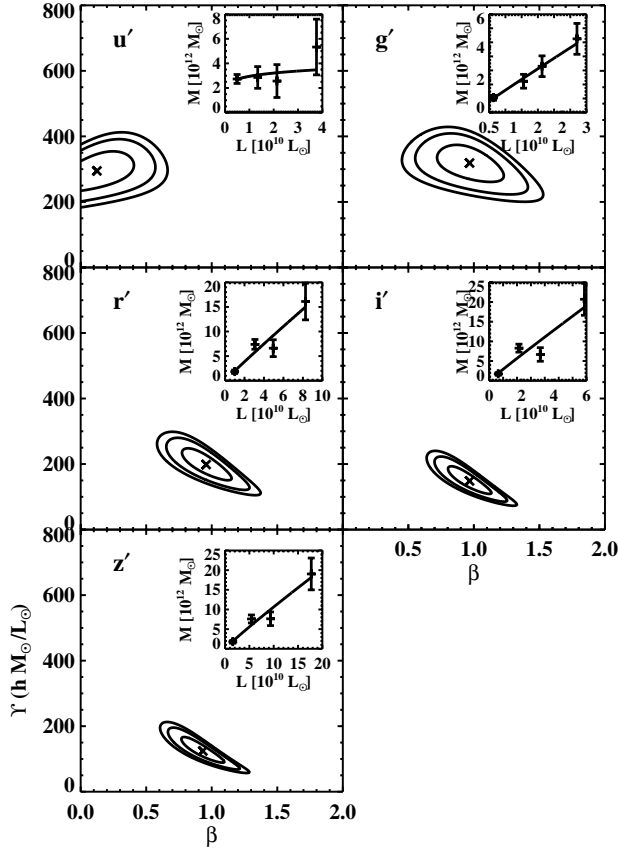
Details of the lens and source sample selection, and the subsequent analysis, can be found in McKay et al. [6]. While the foreground redshift distribution is accurately measured, with a median redshift of 0.1, the background source galaxy redshift distribution is estimated from the magnitude distribution. Using these samples we measure the galaxy-mass correlation function $\xi_{GM}(r)$ around our lens galaxies. It is important to note that the signal we are measuring is extremely small. The peak distortion is only about 0.5%. Despite this tiny signal, the $\xi_{GM}(r)$ is detected at $>13$ $\sigma$ in the g, r, and i bands. The observed $\xi_{GM}(r)$ is well fit by a power law of the form:

$$\xi_{GM}(r) = (2.5 \pm 0.7 h M_\odot pc^{-2}) \times (r/1 Mpc)^{-0.8 \pm 0.2} \tag{1}$$

Now that we have detected $\xi_{GM}(r)$, we can begin to study how it varies with the luminous properties of the lens galaxies. As a first check, we divide all lens galaxies into four luminosity bins in each of the five SDSS colors. We then compare $\xi_{GM}$ in each luminosity bin to probe mass-to-light scalings.

To characterize the variation of $\xi_{GM}(r)$ with lens luminosity, we fit the measured $\xi_{GM}(r)$ from 20-260 $h^{-1}$ kpc with a singular isothermal sphere model. For this best fit model we integrate the associated mass out to 260 $h^{-1}$ kpc and call this $M_{260}$. This outer radius is chosen because contributions to $\xi_{GM}(r)$ due to neighboring galaxies are estimated to be less than 10% at this radius. We then examine how this parameter varies with luminosity. Figure 1 shows for each color the actual $M_{260}$ to light scaling, and then $\chi^2$ contours for the best fit normalization and power law index in each color.
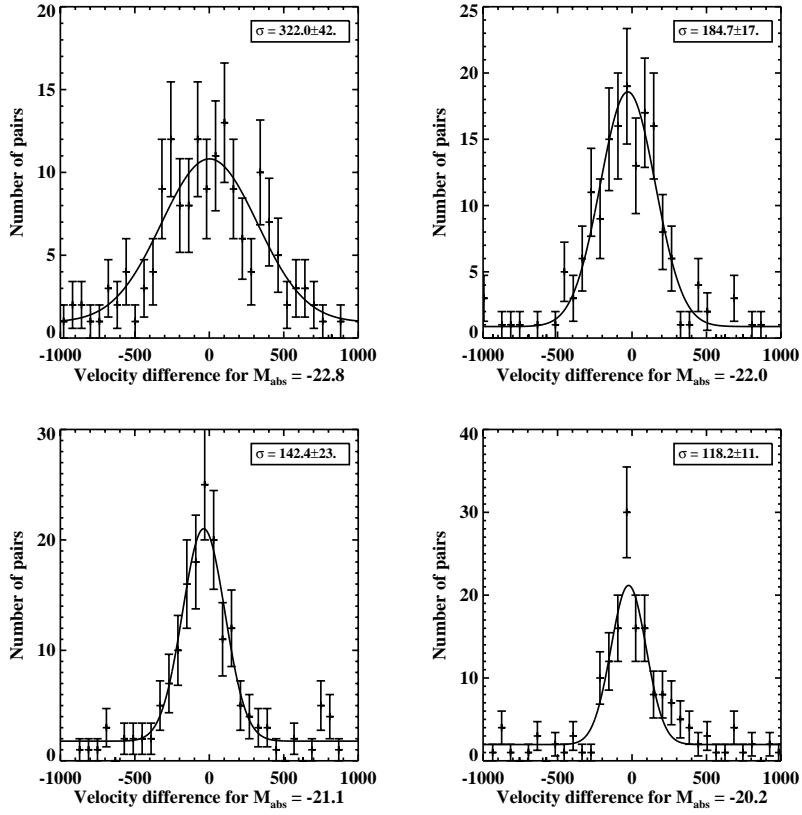
There is little relationship between $M_{260}$ and luminosity in u. This is not surprising as the u luminosity of a galay is often dominated by recent, short-lived, bursts of star formation, and hence does not well reflect the galaxy's mass. But the relationship between $M_{260}$ and luminosity in the other bands is strong, and in every case consistent with linear. In this case, the normalization can be described as a mass-to-light ratio. For the i band the best fit value for this is $M_{260}/L_i = 124 \pm 15 M_\odot/L_\odot$

**FIGURE 1.** The five panels in this figure summarize the relation between $M_{260}$ and luminosity in each of the five SDSS bands. For each band the small inset figure shows this directly. Points in these inset figures are the measured $M_{260}$ and mean luminosity of galaxies in four luminosity bins. The line in these inset figures shows the best fit to a power law relation between $M_{260}$ and luminosity of the form: $M_{260} = \Upsilon \times \left(L_{central}/10^{10}L_{\odot}\right)^{\beta}$. The larger figure shows 68%, 95%, and 99% confidence contours for the fit parameters $\Upsilon$ and $\beta$.

It is important to be cautious in interpreting any such measurement of a mass-to-light ratio. Mass is not the observable quantity. Masses are derived only under rather naive assumptions (a singular isothermal sphere mass model) which while consistent with the data, are not well constrained by it. What we really have is is a measurement of the scaling between luminosity and a fit parameter of a model. If we fit to different models, we may find different scalings. But it is clear at least that the mass of halos on these large scales varies with luminosity.

It is useful to test the conclusions of these lensing measurements using dynamical mass probes [7]. We begin by identifying a set of luminous test particles in orbit around galaxies. In order to identify simple systems, we look for a set of relatively isolated host

**FIGURE 2.** This figure shows velocity difference histograms for faint satellites of four groups of relatively isolated host galaxies. The hosts are grouped by absolute magnitude, ranging from rather $M_{abs}$ = -20.2 at the lower right to $M_{abs}$ = -22.8 at the upper left. The increase in satellite dispersion with host luminosity is clear.

galaxies surrounded by fainter, less luminous, satellites. Each host galaxy has only a few satellites, so we can only measure the average dynamical response of a class of satellites to their hosts.

By constructing a velocity difference histogram for a class of galaxies, we probe the velocity structure of the galaxy-galaxy correlation function $\xi_{GG}(r, \Delta v)$. This velocity structure represents the average dynamical effect of the hosts in the same way that $\xi_{GM}(r)$ represents their average projected surface mass density. We find that the velocity difference width increases significantly with host galaxy luminosity (see Figure 2). While deriving masses from these velocity distributions is quite model dependent, a simple virial mass estimate yields mass-to-light scalings consistent with those derived from lensing [7].

It is reassuring that these two completely different ways of probing mass reveal comparable relationships between mass and light. They are subject to totally different systematic errors, and to quite different problems in interpretation. While this comparison is only a first step, it shows that these combined methods hold great promise for quantifying the relationship between galaxies and their dark matter environments. The signals are there. We have only to understand the systematics.

## Halos identified by groups and clusters

These studies can be expanded to halos identified by groups and clusters in a straightforward way. To do this requires a catalog of groups and clusters, something which is not a standard output of the SDSS analysis system.

For this purpose we have been using group and cluster catalogs derived using the maxBCG method [8, 9]. This method takes advantage of the very uniform colors (the E/S0 ridgeline) of galaxies found in many groups and clusters. Since this color shifts with redshift, groups and clusters appear as overdensities of objects in position color space in a way which is remarkably insensitive to projection effects. This tight correlation in color allows cluster members to be indentified and counted. This estimated number of cluster members, $N_{gal}$, provides an estimate of cluster richness. The color itself provides accurate estimates of redshift.
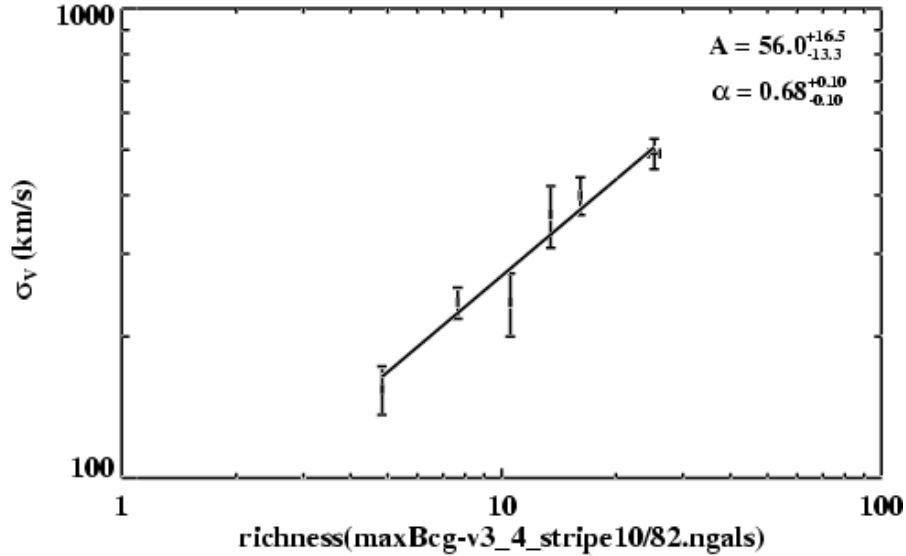
This method is observationally very robust. It's redshift estimates are very good, better than 0.02, and this $N_{gal}$ richness measure is a very well defined count of easily identifiable objects. While it is not yet clear what fraction, or what subclass, of groups this method finds, we have confirmed, by comparison to x-ray selected and other optically selected cluster lists, that it is quite complete in its identification of clusters.

This maxBCG method, fits neatly into the halo occupancy distribution picture of structure formation. In this picture, N-body simulations are used to determine the halo mass function. Then galaxies are included via a halo occupancy distribution function P(N|M), which describes the probability of having N galaxies in a halo of mass M. The maxBCG catalog provides a very clean measure of N, nearly unaffected by projection. It's also complete to roughly z=0.4.

We can now use the same lensing and dynamical measures to probe the relationship between $M_{avg}$ and N. We begin with measurements of the cluster-mass correlation function. Although clusters are rare and our samples are small, they are also massive. So the S/N in cluster shear measurements is very similar to the S/N in galaxy-galaxy lensing measurements [10]. This is generically true across the mass spectrum, so we can use lensing to probe halos of a wide range of masses.

Figure 3 shows the dependence of the cluster-mass correlation function on cluster richness. To determine this, we fix the r dependence by fitting $\xi_{cluster-mass}(r)$ to an SIS model, and examine changes in the normalization with cluster richness. This amounts to measuring an effective velocity dispersion for the clusters. As with the galaxies, the signal is clearly there to trace the variation in the cluster-mass correlation function with richness $N_{gal}$.

To supplement this with dynamical measurements we determine the cluster-galaxy

**FIGURE 3.** This figure shows the variation of an effective velocity dispersion derived from lensing estimates on cluster richness. The effective $\sigma_v$ is derived by fitting $\xi_{cluster-mass}(r)$ to an SIS model out to a radius of 500 $h^{-1}$ kpc. The parameters A and $\alpha$ represent best fits for the normalization and slope of a power law relation between $\sigma_v$ and $N_{gal}$.

correlation function by finding maxBCG objects for which a spectrum was taken of the central galaxies. We then search around these BCGs for spectroscopic neighbors. Then, in a manner analogous to what we did for the galaxies, we make velocity difference histograms in narrow ranges of richness.

A clear variation of the width of this velocity difference histogram with richness is observed. Both of the cluster mass probes I described measure some kind of effective velocity dispersion. Both show a smooth transition from groups with ~200 km/s velocity dispersions to clusters with velocity dispersions of 900 km/s and more. The estimates from velocity dispersion and lensing are consistent at the 20% level. Both are detected at high S/N, so in principle they can yield precise constraints on richness-mass calibrations. In practice, details of the changing size of objects, their velocity structure, and possible velocity biases need to be understood to take full advantage of this.

## Utilizing observable simulated universes

We clearly have a set of experimentally accessible observables which probe the relationship between halos and mass. Imagine that we have a simulation of the universe made with rich enough physics input to reasonably represent observable data. That is, imagine a simulated universe which contains not only dark matter, but luminous galaxies in something like their full variety.

Given such simulations, we can repeat the same 'observations' done in the real universe in an environment which contains all the physics we believe is relevant. We can then compare various predictions to reality at the level of the observables, rather than interposing models which involve untested or patently incorrect assumptions. The remainder of this proceeding describes some early examples of the kinds of comparisons between observations and simulations which we advocate.

The first example involves the GIF simulations [11]. These simulations are built on top of the VIRGO consortium N- body simulations. They add to the N-body outputs by identifying galaxies with the most massive subhalos. Semianalytic prescriptions are used to provide luminosities, colors, and stellar masses for all of these galaxies.

Most important for this study, each of these galaxies has velocity information derived from the full N-body simulation. This allows us to conduct the same dynamical analysis in the simulated data used in the real data. We can compare 'predictions' from the simulations to observations at the observable level (variation of velocity dispersion with luminosity) rather than at the level of model fits. Furthermore, we can use the simulations to tell us how the observables relate to the 'real' masses of the systems.

There are some important limitations, mostly that these simulations include only the most massive and luminous galaxies. So overlap with the observations is not as complete as we would like. Nevertheless, the variation in the width of the velocity difference histogram with host luminosity seen in the GIF simulation is consistent with that seen in SDSS data [7].

What is most useful about doing this in simulations is that we can directly probe the way in which an observable (like this $M_{260}$) relates to a theoretically interesting quantity, like $M_{200}$, the mass measured out to the point where the overdensity is 200 times the mean density. The GIF analysis suggests that these satellite dynamics are indeed probing masses on halo scales, at least up to a scale factor (about 0.7 in this case). There are many reasons to be cautious in asserting this, but especially because the GIF host luminosity range is quite narrow, only about a factor of three.

We are also developing simulated universes aimed at our measurements of halos identified by groups and clusters. For this work, Risa Wechsler has built simulations on the Hubble Volume simulations which are designed to very specifically match observed SDSS galaxies. That is, rather than use semianalytic prescriptions to 'grow' galaxies, she inserts the galaxies in ways which are constrained by data.

The basic algorithm has several steps: Choose an appropriate number of galaxies, with r-band luminosities drawn from the observed SDSS luminosity function. For each galaxy, choose a mass particle in the simulation in a way which matches the observed luminosity dependent clustering seen in the SDSS data. Add passive luminosity evolution, and assign galaxy colors by selecting a real SDSS galaxy with similar luminosity and local density, then translating the SED of this galaxy to the simulated galaxy redshift.

This method produces simulated clusters with properties (like the E/S0 ridgeline) which are remarkably similar to real clusters. By running the maxBCG algorithm on these simulated universes, we can in principle calibrate the relation between richness $N_{gal}$ and mass. But to apply this, we must be certain that the $N_{gal}$ counted in the simulation really reflects what we count in the SDSS data. To check this we are conducting a number of tests. First, we can try matching the halo number distribution and space

density to translate $N_{gal}^{sim}$ to $N_{gal}^{data}$. This is somewhat problematic though, as we'd like to use the cluster number density to constrain cosmology.

As a result it's better to use the cluster-galaxy and cluster-mass correlation functions as constraints. We note in passing that a number of these comparisons are complex. Our simulations do not currently include any special status for a brightest cluster galaxy, despite evidence that they have very unusual dynamics and locations. The simulations also lack higher order galaxy property correlations, velocity bias and so on. Much work remains, but progress is being made pretty rapidly.

## CONCLUSIONS

If we are to determine the distribution of matter in the universe, we must understand in some detail the relationship between the dark matter halos which dominate the mass budget and the luminous galaxies which illuminate them. We must use the locations of galaxies to identify halos, after which we can study their properties by both lensing and dynamical means.

Relating the observables, $\xi_{GM}(r)$ and $\xi_{GG}(r, \Delta v)$, to halo mass requires modeling. Even with this small subset of SDSS data, uncertainty in this modeling already dominates our ability to interpret these results. As we make ever more precise, systematic uncertainty in this modeling will continue to limit our ability to interpret the results.

The most straightforward predictions of structure formation theory describe the halo mass function, clustering, and structure. Since these are not observables, comparison of observations to theory relies on ill-determined intermediate modeling steps. To avoid this step, we must propagate theoretical predictions forward to the observable level. This can be done by including galaxies, with all their observable properties, in simulations of structure formation. Initial attempts at this kind of prediction are provided by, for example, the GIF simulations [11].

We describe some example comparisons between SDSS observations and predictions of observables made using such simulations. While important details in the simulations remain to be checked, initial results are very encouraging. It seems likely that this method, comparing observations to theory at the observable level through structure formation simulations which include galaxies, will play an important role in the interpretation of precise new observations of the distribution of matter in the universe.

## ACKNOWLEDGMENTS

## REFERENCES

1.  Jenkins, A. et al. 1998, Astrophysical Journal, 499, 20
2.  Verde, L. et al. 2002, Monthly Notices of the Royal Astronomical Society, 335, 432
3.  Benson, A. J., Cole, S., Frenk, C. S., Baugh, C. M., & Lacey, C. G. 2000, Monthly Notices of the Royal Astronomical Society, 311, 793
4.  Berlind, A. A. & Weinberg, D. H. 2002, Astrophysical Journal, 575, 587
5.  Stoughton, C. et al. 2002, Astronomical Journal, 123, 485
6.  McKay, T., et al. 2001, astro-ph/0108013
7.  McKay, T. A. et al. 2002, Astrophysical Journal Letters, 571, L85
8.  Gladders, M. D. & Yee, H. K. C. 2000, Astronomical Journal, 120, 2148
9.  Annis, J., et al. 2003, in preparation
10.  Sheldon, E. S. et al. 2001, Astrophysical Journal, 554, 881
11.  Kauffmann, G., Colberg, J. M., Diaferio, A., & White, S. D. M. 1999, Monthly Notices of the Royal Astronomical Society, 303, 188