

Reconstruction of the early Universe as a convex optimization problem

Y. Brenier,¹ U. Frisch,^{2,3*} M. Hénon,² G. Loeper,¹ S. Matarrese,⁴ R. Mohayaee²
and A. Sobolevskii^{2,5}

¹CNRS, UMR 6621, Université de Nice-Sophia-Antipolis, Parc Valrose, 06108 Nice Cedex 02, France

²CNRS, UMR 6529, Observatoire de la Côte d'Azur, BP 4229, 06304 Nice Cedex 4, France

³Institute for Advanced Study, Einstein Drive, Princeton, NJ 08540, USA

⁴Dipartimento di Fisica 'G. Galilei' and INFN, Sezione di Padova, via Marzolo 8, 35131-Padova, Italy

⁵Department of Physics, M. V. Lomonossov Moscow University, Leninskie Gory, 119992 Moscow, Russia

Accepted 2003 August 10. Received 2003 August 1; in original form 2003 April 10

ABSTRACT

We show that the deterministic past history of the Universe can be uniquely reconstructed from knowledge of the present mass density field, the latter being inferred from the three-dimensional distribution of luminous matter, assumed to be tracing the distribution of dark matter up to a known bias. Reconstruction ceases to be unique below those scales – a few Mpc – where multistreaming becomes significant. Above $6 h^{-1}$ Mpc we propose and implement an effective Monge–Ampère–Kantorovich method of unique reconstruction. At such scales the Zel'dovich approximation is well satisfied and reconstruction becomes an instance of optimal mass transportation, a problem which goes back to Monge. After discretization into N point masses one obtains an assignment problem that can be handled by effective algorithms with not more than $O(N^3)$ time complexity and reasonable CPU time requirements. Testing against N -body cosmological simulations gives over 60 per cent of exactly reconstructed points.

We apply several interrelated tools from optimization theory that were not used in cosmological reconstruction before, such as the Monge–Ampère equation, its relation to the mass transportation problem, the Kantorovich duality and the auction algorithm for optimal assignment. A self-contained discussion of relevant notions and techniques is provided.

Key words: hydrodynamics – cosmology: theory – early Universe – large-scale structure of Universe.

1 INTRODUCTION

Can one follow back in time to initial locations the highly structured present distribution of mass in the Universe, as mapped by redshift catalogues of galaxies? At first this seems an ill-posed problem since little is known concerning the peculiar velocities of galaxies, so that equations governing the dynamics cannot just be integrated back in time. In fact, it is precisely one of the goals of reconstruction to determine the peculiar velocities. Since the pioneering work of Peebles (1989), a number of reconstruction techniques have been proposed, which frequently provided non-unique answers.¹

Cosmological reconstruction should, however, take advantage of our knowledge that the initial mass distribution was quasi-uniform at baryon–photon decoupling, about 14 billion years ago (see, e.g., Susperregi & Binney 1994). In a recent Letter to Nature (Frisch et al.

2002), four of us have shown that, with suitable assumptions, this a priori knowledge of the initial density field makes reconstruction a well-posed instance of what is called the optimal mass transportation problem.

A well-known fact is that, in an expanding universe with self-gravitating matter, the initial velocity field is ‘slaved’ to the initial gravitational field, which is potential; both fields thus depend on a single scalar function. Hence the number of unknowns matches the number of constraints, namely the single density function characterizing the present distribution of mass.

This observation alone, of course, does not ensure uniqueness of the reconstruction. For this, two restrictions will turn out to be crucial. First, from standard redshift catalogues it is impossible to resolve individual streams of matter with different velocities if they occupy the same space volume. This ‘multistreaming’ is typically confined to relatively small scales of a few megaparsecs (Mpc), below which reconstruction is hardly feasible. Secondly, to reconstruct a given finite patch of the present Universe, we need to know its initial shape at least approximately.

It is our purpose in the present paper to clarify the physical nature of the factors permitting a unique reconstruction and of obstacles

*E-mail: uriel@obs-nice.fr

¹ We put the present work in the context of several important existing techniques in Section 7.

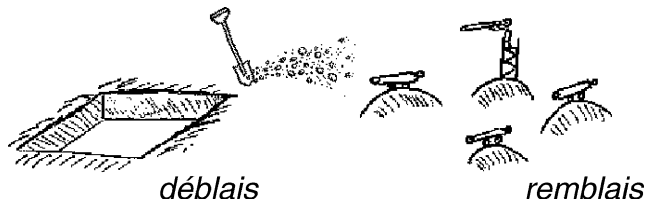


Figure 1. A sketch of Monge’s mass transportation problem in which one searches the optimal way of transporting earth from cuts (*déblais*) to fills (*remblais*), each of prescribed shape; the cost of transporting a molecule of earth is a given function of the distance. The MAK method of reconstructing the early Universe described in this paper corresponds to a quadratic cost.

limiting it, and to give a detailed account of the way some recent developments in the optimal mass transportation theory are applicable. (Fig. 1 may give the reader some feeling of what mass transportation is about.)

The paper is organized as follows. In Section 2 we formulate the reconstruction problem in an expanding universe and state the main result concerning uniqueness of the solution.

In the next three sections we devise and test a reconstruction technique called MAK (for Monge–Ampère–Kantorovich) within a restricted framework where the Lagrangian map from initial to present mass locations is taken as potential. In Section 3 we discuss the validity of the potentiality assumption and its relation to various approximations used in cosmology; then we derive the Monge–Ampère equation, a simple consequence of mass conservation, introduce its modern reformulation as a Monge–Kantorovich problem of optimal mass transportation and finally discuss different limitations on uniqueness of the reconstruction. In Section 4 we show how discretization turns optimization into an instance of the standard assignment problem; we then present effective algorithms for its solution, foremost the ‘auction’ algorithm of D. Bertsekas. Section 5 is devoted to testing the MAK reconstruction against N -body cosmological simulations.

In Section 6, we show how the general case, without the potentiality assumption, can also be recast as an optimization problem with a unique solution and indicate a possible numerical strategy for such a reconstruction. In Section 7 we compare our reconstruction method with other approaches in the literature. In Section 8 we discuss perspectives and open problems.

A number of topics are left for appendices. In Appendix A we derive the Eulerian and Lagrangian equations in the form used throughout the paper (and provide some background for non-cosmologists). Appendix B is devoted to the history of optimal mass transportation theory, a subject more than two centuries old (Monge 1781), which has undergone significant progress within the last two decades. Appendix C is a brief elementary introduction to the technique of duality in optimization, which we use several times throughout the paper. Appendix D gives details of the uniqueness proof that is only outlined in Section 6.

Finally, a word concerning notation (see also Appendix A). We are using comoving coordinates denoted by \mathbf{x} in a frame following the expansion of the Universe. Our time variable is not the cosmic time but the so-called linear growth factor, here denoted by τ , the use of which gives to certain equations the same form as for compressible fluid dynamics in a non-expanding medium. The subscript 0 refers to the present time (redshift $z = 0$), while the quantities evaluated at the initial epoch take the subscript or superscript ‘in’. Following cosmological usage, the Lagrangian coordinate is denoted by \mathbf{q} .

2 RECONSTRUCTION IN AN EXPANDING UNIVERSE

The most widely accepted explanation of the large-scale structure seen in galaxy surveys is that it results from small primordial fluctuations that grew under gravitational self-interaction of collisionless cold dark matter (CDM) particles in an expanding universe (see, e.g., Bernardeau et al. (2002) and references therein). The relevant equations of motion, derived in Appendix A, are the Euler–Poisson equations² written here for a flat, matter-dominated Einstein–de Sitter universe (for a more general case see, e.g., Catelan et al. 1995):

$$\partial_\tau \mathbf{v} + (\mathbf{v} \cdot \nabla_x) \mathbf{v} = -\frac{3}{2\tau} (\mathbf{v} + \nabla_x \varphi_g), \quad (1)$$

$$\partial_\tau \rho + \nabla_x \cdot (\rho \mathbf{v}) = 0, \quad (2)$$

$$\nabla_x^2 \varphi_g = \frac{\rho - 1}{\tau}, \quad (3)$$

where \mathbf{v} denotes the velocity, ρ denotes the density (normalized by the background density $\bar{\rho}$) and φ_g is a rescaled gravitational potential. All quantities are expressed in comoving spatial coordinates \mathbf{x} and linear growth factor τ , which is used as the time variable; in particular, \mathbf{v} is the Lagrangian τ -time derivative of the comoving coordinate of a fluid element.

2.1 Slaving in early-time dynamics and its fossils

The right-hand sides of the momentum and Poisson equations (1) and (3) contain denominators proportional to τ . Hence, a necessary condition for the problem not to be singular as $\tau \rightarrow 0$ is

$$\mathbf{v}_{\text{in}}(\mathbf{x}) + \nabla_x \varphi_g^{\text{in}} = 0, \quad \rho_{\text{in}}(\mathbf{x}) = 1. \quad (4)$$

In other words, (i) the initial velocity must be equal to (minus) the gradient of the initial gravitational potential and (ii) the initial normalized mass distribution is uniform. We shall refer to these conditions as *slaving*. Note that the density contrast $\rho - 1$ vanishes initially, but the rescaled gravitational potential and the velocity, as defined here, stay finite thanks to our choice of the linear growth factor as a time variable. Therefore, we refer to the initial mass distribution as being ‘quasi-uniform’.

In the following, when we mention the Euler–Poisson *initial-value problem*, it is always understood that we start at $\tau = 0$ and assume slaving. Hence we are extending the Newtonian matter-dominated post-decoupling description back to $\tau = 0$. By examination of the Lagrangian equations for $\mathbf{x}(\mathbf{q}, \tau)$ near $\tau = 0$, which can be linearized because the displacement $\mathbf{x} - \mathbf{q}$ is small, it is easily shown that slaving implies the absence of the ‘decaying mode’, which behaves as $\tau^{-3/2}$ in an Einstein–de Sitter universe and is thus singular at $\tau = 0$ (for details see Appendix A).

Slaving is also a sufficient condition for the initial problem to be well posed. It is indeed easily shown recursively that (1)–(3) admit a solution in the form of a formal Taylor series in τ (a related expansion involving only potentials may be found in Catelan et al. 1995):

$$\mathbf{v}(\mathbf{x}, \tau) = \mathbf{v}^{(0)}(\mathbf{x}) + \tau \mathbf{v}^{(1)}(\mathbf{x}) + \tau^2 \mathbf{v}^{(2)}(\mathbf{x}) + \dots, \quad (5)$$

$$\varphi_g(\mathbf{x}, \tau) = \varphi_g^{(0)}(\mathbf{x}) + \tau \varphi_g^{(1)}(\mathbf{x}) + \tau^2 \varphi_g^{(2)}(\mathbf{x}) + \dots, \quad (6)$$

$$\rho(\mathbf{x}, \tau) = 1 + \tau \rho^{(1)}(\mathbf{x}) + \tau^2 \rho^{(2)}(\mathbf{x}) + \dots. \quad (7)$$

Furthermore, $\mathbf{v}^{(n)}(\mathbf{x})$ is easily shown to be curl-free for any n .

² Also often called the Euler equations.

Several important consequences of slaving extend to later times as ‘fossils’ of the earliest dynamics. First, as already stressed in the introduction, the whole dynamics is determined by only one scalar field (e.g. the initial gravitational potential), which we can hope to determine from the knowledge of the present density field.

Secondly, slaving trivially rules out multistreaming up to the time of formation of caustics. Since we are working with collisionless matter, the dynamics should in principle be governed by the Vlasov–Poisson³ kinetic equation, which allows at each (\mathbf{x}, τ) point a non-trivial distribution function $f(\mathbf{x}, \mathbf{v}, \tau)$. Slaving selects a particular class of solutions for which the distribution function is concentrated on a single-speed manifold, thereby justifying the use of the Euler–Poisson equation without having to invoke any hydrodynamical limit (see, e.g., Vergassola et al. 1994; Catelan et al. 1995).

Thirdly, it is easily checked from (1) that the initial slaved velocity, which is obviously curl-free, remains so for all later times (up to formation of caustics). Note that this vanishing of the curl holds in Eulerian coordinates. A similar property in Lagrangian coordinates can only hold approximately but will play an important role in the following (Section 3).

2.2 Formulation of the reconstruction problem

The present Universe is replete with high-density structures: clusters (point-like objects), filaments (line-like objects) and perhaps sheets or walls.⁴

The internal structure of such *mass concentrations* certainly displays multistreaming and cannot be described in terms of a single-speed solution to the Euler–Poisson equations. In N -body simulations, multistream regions are usually found to be of relatively small extension in one or several space directions, typically not more than a few Mpc, and hence have a small volume, although they contain a significant fraction of the total mass (see, e.g., Weinberg & Gunn 1990).

In order not to have to deal with tiny multistream regions, we replace the true mass distribution by a ‘macroscopic’ one, which has a regular part and a singular (collapsed) part, the latter concentrated on objects of dimension less than three, such as points or lines.

The general problem of reconstruction is to find as much information as possible on the history of the evolution that carries the initial uniform density into the present macroscopic mass distribution, including the evolution of the velocities. In principle we would like to find a solution of the Euler–Poisson initial-value problem leading to the present density field $\rho_0(\mathbf{x})$.

A more restricted problem, which we call the ‘displacement reconstruction’, is to find the Lagrangian map $\mathbf{q} \mapsto \mathbf{x}(\mathbf{q})$ and its inverse $\mathbf{x} \mapsto \mathbf{q}(\mathbf{x})$, or, in other words, to answer the question: where does a given ‘Monge molecule’⁵ of matter originate from? Of course, the inverse Lagrangian map will not be single-valued on mass concentrations. Furthermore, for practical cosmological applications, we define a ‘full reconstruction problem’ as (i) displacement reconstruction and (ii) obtaining the initial and present peculiar velocity fields, $\mathbf{v}_{\text{in}}(\mathbf{q})$ and $\mathbf{v}_0(\mathbf{x})$.

We shall show in this paper that the displacement reconstruction problem is uniquely solvable and that the full reconstruction prob-

lem has a unique solution outside of mass concentrations; as to the latter, they are traced back to *collapsed regions* in the Lagrangian space, the shape and positions of which are well defined but the inner structure of density and velocity fluctuations is irretrievably lost.

3 POTENTIAL LAGRANGIAN MAPS: THE MAK RECONSTRUCTION

In this and the next two sections we shall assume that the Lagrangian map from initial positions to present ones is potential

$$\mathbf{x} = \nabla_{\mathbf{q}} \Phi(\mathbf{q}), \quad (8)$$

and furthermore that the potential $\Phi(\mathbf{q})$ is convex, which is, as we shall see, related to the absence of multistreaming.

3.1 Approximations leading to maps with convex potentials

The motivation for the potential assumption, first used by Bertschinger & Dekel (1989),⁶ comes from the Zel’dovich approximation (Zel’dovich 1970), denoted here by ZA, and its refinements. To recall how the ZA comes about, let us start from the equations for the Lagrangian map $\mathbf{x}(\mathbf{q}, \tau)$, written in the Lagrangian coordinate \mathbf{q} (Appendix A)

$$D_{\tau}^2 \mathbf{x} = -\frac{3}{2\tau} (D_{\tau} \mathbf{x} + \nabla_{\mathbf{x}} \phi_{\text{g}}), \quad (9)$$

$$\nabla_{\mathbf{x}}^2 \phi_{\text{g}} = \frac{1}{\tau} [(\det \nabla_{\mathbf{q}} \mathbf{x})^{-1} - 1], \quad (10)$$

where D_{τ} is the Lagrangian time derivative and $\nabla_{x_i} \equiv (\partial q_j / \partial x_i) \nabla_{q_j}$ is the Eulerian gradient rewritten in Lagrangian coordinates. As shown in Appendix A, in one space dimension the Hubble drag term $D_{\tau} \mathbf{x}$ and the gravitational acceleration term $\nabla_{\mathbf{x}} \phi_{\text{g}}$ cancel exactly. Slaving, discussed in Section 2.1, means that the same cancellation holds to leading order in any dimension for small τ . The ZA extends this as an approximation without the restriction of small τ . Within the ZA, the acceleration $D_{\tau}^2 \mathbf{x}$ vanishes. Hence the Lagrangian map has the form

$$\begin{aligned} \mathbf{x}(\mathbf{q}, \tau) &= \mathbf{q} + \tau (D_{\tau} \mathbf{x})_{\text{in}}(\mathbf{q}) = \mathbf{q} - \tau \nabla_{\mathbf{q}} \phi_{\text{g}}^{\text{in}}(\mathbf{q}) \\ &= \nabla_{\mathbf{q}} \Phi(\mathbf{q}, \tau) \end{aligned} \quad (11)$$

with the potential

$$\Phi(\mathbf{q}, \tau) \equiv \frac{|\mathbf{q}|^2}{2} - \tau \phi_{\text{g}}^{\text{in}}(\mathbf{q}). \quad (12)$$

Furthermore, taking the time derivative of (11), we see that the velocity $D_{\tau} \mathbf{x}(\mathbf{q}, \tau)$ is curl-free with respect to the Lagrangian coordinate \mathbf{q} .

Potentiality of the Lagrangian map (and consequently the Lagrangian potentiality of the velocity) is perhaps the most important feature of the ZA. Unlike the vanishing of the acceleration, it does not depend on the choice of the linear growth factor as the time variable. However, *unaccelerated* but *vortical* flow would fail to exhibit the cancellation necessary for the ZA to hold. It is noteworthy that the potentiality is not limited to the ZA: indeed, the latter can be formulated as the first order of a systematic Lagrangian perturbation theory in which, up to second order, the Lagrangian map is still potential under slaving (Moutarde et al. 1991; Buchert 1992; Buchert & Ehlers 1993; Munshi, Sahni & Starobinsky 1994; Catelan 1995).

⁶ In connection with what was called later the Lagrangian POTENT method (Dekel, Bertschinger & Faber 1990).

³ Actually written for the first time by Jeans (1919).

⁴ Whether the Great Wall and the Sculptor Wall are sheet-like or filament-like is a moot point (Sathyaprakash et al. 1998).

⁵ For Monge and his contemporaries, the word ‘molecule’ meant a Leibniz infinitesimal element of mass; see Appendix B.

It is well known that the ZA map defined by (11) ceases in general to be invertible due to the formation of multistream regions bounded by caustics. Since particles move along straight lines in the ZA, the formation of caustics proceeds just as in ordinary optics in a uniform medium in which light rays are also straight.⁷ One of the problems with the ZA is that caustics, which start as localized objects, quickly grow in size and give unrealistically large multistream regions.

A modification of the ZA that has no multistreaming at all, but sharp mass concentrations in the form of shocks and other singularities, has been introduced by Gurbatov & Saichev (1984; see also Gurbatov, Saichev & Shandarin 1989; Shandarin & Zel'dovich 1989). It is known as the *adhesion model*. In Eulerian coordinates it amounts to using a multidimensional Burgers equation (see, e.g., Frisch & Bec 2002)

$$\partial_\tau \mathbf{v} + (\mathbf{v} \cdot \nabla_x) \mathbf{v} = \nu \nabla_x^2 \mathbf{v}, \quad \mathbf{v} = -\nabla_x \phi_v, \quad (13)$$

taken in the limit where the viscosity ν tends to zero. In Lagrangian coordinates, the adhesion model is obtained from the ZA by replacing the velocity potential $\Phi(\mathbf{q}, t)$ given by (12) by its *convex hull* $\Phi_c(\mathbf{q}, t)$ in the \mathbf{q} variable (Vergassola et al. 1994).

Convexity is a concept that plays an important role in this paper, and a few words on it are in order here (see also Appendix C1). A body in the three-dimensional (3D) space is said to be *convex* if, whenever it contains two points, it also contains the whole segment joining them. A function $f(\mathbf{q})$ is said to be convex if the set of all points lying above its graph is convex. The convex hull of the function $\Phi(\mathbf{q})$ is defined as the largest convex function for which the graph lies below that of $\Phi(\mathbf{q})$. In two dimensions it can be visualized by wrapping the graph of $\Phi(\mathbf{q})$ tightly from below with an elastic sheet.

Note that $\Phi(\mathbf{q}, \tau)$ given by (12) is obviously convex for small enough τ since it is then very close to the parabolic function $|\mathbf{q}|^2/2$. After caustics form, convexity is lost in the ZA but recovered with the adhesion model. It may then be shown that those regions in the Lagrangian space where $\Phi(\mathbf{q}, t)$ does not coincide with its convex hull will be mapped in the Eulerian space to sheets, lines and points, each of which contains a finite amount of mass. At these locations the Lagrangian map does not have a uniquely defined Lagrangian antecedent but such points form a set of vanishing volume. Everywhere else, there is a unique antecedent and hence no multistreaming.

Although the adhesion model has a number of known shortcomings, such as non-conservation of momentum in more than one dimension, it has been found to be in better agreement with N -body simulations than the ZA (Weinberg & Gunn 1990). Other single-speed approximations to multistream flow, overcoming difficulties of the adhesion model, are discussed, for example, by Shandarin & Sathyaprakash (1996), Buchert & Domínguez (1998) and Fanelli & Aurell (2002). In such models, multistreaming is completely suppressed by a mechanism of momentum exchange between neighbouring streams with different velocities. This is of course a common phenomenon in ordinary fluids, where it is due to viscous diffusion; dark matter is, however, essentially collisionless and the usual mechanism for generating viscosity does not operate, so that a non-collisional mechanism must be invoked. A qualitative explanation using the modification of the gravitational forces after the formation of caustics has been proposed by Shandarin & Zel'dovich (1989). In our opinion the mechanism limiting multistreaming to rather nar-

row regions is poorly understood and deserves considerable further investigation.

3.2 The Monge–Ampère equation: a consequence of mass conservation and potentiality

We now show that the assumption that the Lagrangian map is derived from a convex potential leads to a pair of non-linear partial differential equations, one for this potential and another for its Legendre transform.

Let us first assume that the present distribution of mass has no singular part, an assumption that we shall relax later. Since in our notation the initial quasi-uniform mass distribution has unit density, mass conservation implies $\rho_0(\mathbf{x}) d^3\mathbf{x} = d^3\mathbf{q}$, which can be rewritten in terms of the Jacobian matrix $\nabla_{\mathbf{q}}\mathbf{x}$ as

$$\det \nabla_{\mathbf{q}} \mathbf{x} = \frac{1}{\rho_0(\mathbf{x}(\mathbf{q}))}. \quad (14)$$

Under the potential assumption (8), this takes the form

$$\det(\nabla_{q_i} \nabla_{q_j} \Phi(\mathbf{q})) = \frac{1}{\rho_0(\nabla_{\mathbf{q}} \Phi(\mathbf{q}))}. \quad (15)$$

A similar equation also follows from equations (1) and (2) of Bertschinger & Dekel (1989).

A simpler equation, in which the unknown appears only in the left-hand side, namely equation (19) below, is obtained for the potential of the *inverse Lagrangian map* $\mathbf{q}(\mathbf{x})$. Key is the observation that the inverse of a map with a convex potential also has a convex potential, and that the two potentials are Legendre transforms of each other.⁸ A purely local proof of this statement is to observe that potentiality of $\mathbf{q}(\mathbf{x})$ is equivalent to the symmetry of the *inverse Jacobian matrix* $\nabla_{\mathbf{x}}\mathbf{q}$, which follows because it is the inverse of the symmetrical matrix $\nabla_{\mathbf{q}}\mathbf{x}$; convexity is equivalent to the positive-definiteness of these matrices. Obviously the function

$$\Theta(\mathbf{x}) \equiv \mathbf{x} \cdot \mathbf{q}(\mathbf{x}) - \Phi(\mathbf{q}(\mathbf{x})), \quad (16)$$

which is the Legendre transform of $\Phi(\mathbf{q})$, is the potential for the inverse Lagrangian map. The modern definition of the Legendre transformation (see Appendix C1), needed for generalization to non-smooth mass distributions, is

$$\Theta(\mathbf{x}) = \max_{\mathbf{q}} \mathbf{x} \cdot \mathbf{q} - \Phi(\mathbf{q}), \quad (17)$$

$$\Phi(\mathbf{q}) = \max_{\mathbf{x}} \mathbf{x} \cdot \mathbf{q} - \Theta(\mathbf{x}). \quad (18)$$

In terms of the potential Θ , mass conservation is immediately written as

$$\det(\nabla_{x_i} \nabla_{x_j} \Theta(\mathbf{x})) = \rho_0(\mathbf{x}). \quad (19)$$

This equation, which has the determinant of the second derivatives of the unknown in the left-hand side and a prescribed (positive) function in the right-hand side, is called the (elliptic) Monge–Ampère equation (see Appendix B for a historical perspective).

Note that our Monge–Ampère equation may be viewed as a non-linear generalization of the Poisson equation (used for reconstruction by Nusser & Dekel 1992; see also Section 7.1), to

⁷ Catastrophe theory has been used to classify the different types of singularities thus obtained (Arnol'd, Shandarin & Zel'dovich 1982).

⁸ Besides our problem, this fact prominently appears in two other fields of physics: in classical mechanics, the Lagrangian and Hamiltonian functions are Legendre transforms of each other – their gradients relate the generalized velocity and momentum – and so are, in thermodynamics, the internal energy and the Gibbs potential, implying the same relation between extensive and intensive parameters of state.

which it reduces if particles have moved very little from their initial positions.

In actual reconstructions we have to deal with mass concentration in the present distribution of matter. Thus the density in the right-hand side of (19) has a singular component (a Dirac distribution concentrated on sets carrying the concentrated mass) and the potential Θ ceases to be smooth. As we now show, a generalized meaning can nevertheless be given to the Monge–Ampère equation by using the key ingredient in its derivation, namely mass conservation, in integrated form.

For a non-smooth convex potential Θ , taking the gradient $\nabla_x \Theta(\mathbf{x})$ still makes sense if one allows it to be multivalued at points where the potential is not differentiable. The gradient at such a point \mathbf{x} is then the set of all possible slopes of planes touching the graph of Θ at $(\mathbf{x}, \Theta(\mathbf{x}))$ (this idea is given a precise mathematical formulation in Appendix C1). As \mathbf{x} varies over an arbitrary domain \mathcal{D}_E in the Eulerian space, its image $\mathbf{q}(\mathbf{x})$ sweeps a domain $\mathbf{q}(\mathcal{D}_E)$ in the Lagrangian space, and mass conservation requires that

$$\int_{\mathcal{D}_E} \rho_0(\mathbf{x}) d^3 \mathbf{x} = \int_{\nabla_x \Theta(\mathcal{D}_E)} d^3 \mathbf{q}, \quad (20)$$

where we take into account that $\mathbf{q}(\mathbf{x}) = \nabla_x \Theta(\mathbf{x})$. Equation (20) must hold for any Eulerian domain \mathcal{D}_E ; this requirement is known as the *weak formulation* of the Monge–Ampère equation (19). A symmetric formulation may be written for (15) in terms of $\mathbf{x}(\mathbf{q}) = \nabla_q \Phi(\mathbf{q})$. For further material on the weak formulation see, e.g., Pogorelov (1978).

Considerable literature has been devoted to the Monge–Ampère equation in recent years (see, e.g., Caffarelli 1999; Caffarelli & Milman 1999). We mention now a few results that are of direct relevance for the reconstruction problem.

In a nutshell, one can prove that when the domains occupied by the mass initially and at present are bounded and convex, the Monge–Ampère equation – in its weak formulation – is guaranteed to have a unique solution, which is smooth unless one or both of the mass distributions is non-smooth. The actual construction of this solution can be done by a variational method discussed in the next section.

A similar result also holds when the present density field is periodic and the same periodicity is assumed for the map.

Also relevant, as we shall see in Section 3.4, is a recent result of Caffarelli & Li (2003): if the Monge–Ampère equation is considered in the whole space, but the present density contrast $\delta = \rho - 1$ vanishes outside of a bounded set, then the solution $\Theta(\mathbf{x})$ is determined uniquely up to prescription of its asymptotic behaviour at infinity, which is specified by a quadratic function of the form

$$\theta(\mathbf{x}) \equiv (\mathbf{x} \cdot \mathbf{A} \mathbf{x}) + (\mathbf{b} \cdot \mathbf{x}) + c, \quad (21)$$

for some positive-definite symmetric matrix \mathbf{A} with unit determinant, vector \mathbf{b} and constant c .

3.3 Optimal mass transportation

As we are going to see now, the Monge–Ampère equation (19) is equivalent to an instance of what is called the ‘optimal mass transportation problem’. Suppose we are given two distributions $\rho_{in}(\mathbf{q})$ and $\rho_0(\mathbf{x})$ of the same amount of mass in two three-dimensional convex bounded domains \mathcal{D}_{in} and \mathcal{D}_0 . The optimal mass transportation problem is then to find the most cost-effective way of rearranging by a suitable map one distribution into the other, the cost of transporting a unit of mass from a position $\mathbf{q} \in \mathcal{D}_{in}$ to $\mathbf{x} \in \mathcal{D}_0$ being a prescribed function $c(\mathbf{q}, \mathbf{x})$.

Denoting the map by $\mathbf{x}(\mathbf{q})$ and its inverse $\mathbf{q}(\mathbf{x})$, we can write the problem as the requirement that the cost

$$I \equiv \int_{\mathcal{D}_{in}} c(\mathbf{q}, \mathbf{x}(\mathbf{q})) \rho_{in}(\mathbf{q}) d^3 \mathbf{q} = \int_{\mathcal{D}_0} c(\mathbf{q}(\mathbf{x}), \mathbf{x}) \rho_0(\mathbf{x}) d^3 \mathbf{x} \quad (22)$$

be minimum, with the constraints of prescribed ‘terminal’ densities ρ_{in} and ρ_0 and of mass conservation $\rho_{in}(\mathbf{q}) d^3 \mathbf{q} = \rho_0(\mathbf{x}) d^3 \mathbf{x}$.⁹

This problem goes back to Monge (1781) who considered the case of a linear cost function $c(\mathbf{q}, \mathbf{x}) = |\mathbf{x} - \mathbf{q}|$ (see Appendix B and Fig. 1).

For our purposes, the central result is that *the problem of finding a potential Lagrangian map with prescribed initial and present mass density fields is equivalent to a mass transportation problem with quadratic cost*. Indeed, it is known (Brenier 1987, 1991) that, when the cost is a quadratic function of the distance, so that

$$I = \int_{\mathcal{D}_{in}} \frac{|\mathbf{x}(\mathbf{q}) - \mathbf{q}|^2}{2} \rho_{in}(\mathbf{q}) d^3 \mathbf{q} = \int_{\mathcal{D}_0} \frac{|\mathbf{x} - \mathbf{q}(\mathbf{x})|^2}{2} \rho_0(\mathbf{x}) d^3 \mathbf{x}, \quad (23)$$

the solution $\mathbf{q}(\mathbf{x})$ to the optimal mass transportation problem is the gradient of a convex function, which then must satisfy the Monge–Ampère equation (19) by mass conservation.

A particularly simple variational proof can be given for the smooth case, when the two mutually inverse maps $\mathbf{x}(\mathbf{q})$ and $\mathbf{q}(\mathbf{x})$ are both well defined.

Performing a variation of the map $\mathbf{x}(\mathbf{q})$, we cause a mass element in the Eulerian space that was located at $\mathbf{x}(\mathbf{q})$ to move to $\mathbf{x}(\mathbf{q}) + \delta \mathbf{x}(\mathbf{q})$. This variation is constrained not to change the density field ρ_0 . To express this constraint it is convenient to rewrite the displacement in Eulerian coordinates $\delta \mathbf{x}_E(\mathbf{x}) \equiv \delta \mathbf{x}(\mathbf{q}(\mathbf{x}))$. Noting that the point \mathbf{x} gets displaced into $\mathbf{y} = \mathbf{x} + \delta \mathbf{x}$, we thus require that $\rho_0(\mathbf{x}) d^3 \mathbf{x} = \rho_0(\mathbf{y}) d^3 \mathbf{y}$ or

$$\rho_0(\mathbf{x}) = \rho_0[\mathbf{x} + \delta \mathbf{x}_E(\mathbf{x})] \det(\nabla_x [\mathbf{x} + \delta \mathbf{x}_E(\mathbf{x})]). \quad (24)$$

Expanding this equation, we find that, to leading order,

$$\nabla_x \cdot [\rho_0(\mathbf{x}) \delta \mathbf{x}_E(\mathbf{x})] = 0, \quad (25)$$

an equation which just expresses the physically obvious fact that the mass flux $\rho_0(\mathbf{x}) \delta \mathbf{x}_E(\mathbf{x})$ should have zero divergence. Performing the variation on the functional I given by (23), we obtain

$$\begin{aligned} \delta I &= \int_{\mathcal{D}_{in}} [\mathbf{x}(\mathbf{q}) - \mathbf{q}] \cdot \delta \mathbf{x}(\mathbf{q}) \rho_{in}(\mathbf{q}) d^3 \mathbf{q} \\ &= \int_{\mathcal{D}_0} [\mathbf{x} - \mathbf{q}(\mathbf{x})] \cdot [\rho_0(\mathbf{x}) \delta \mathbf{x}_E(\mathbf{x})] d^3 \mathbf{x} = 0, \end{aligned} \quad (26)$$

which has to hold under the constraint (25). In other words, the displacement $\mathbf{x} - \mathbf{q}(\mathbf{x})$ has to be orthogonal (in the L_2 functional sense) to all divergenceless vector fields and, thus, must be a gradient. Since \mathbf{x} is obviously a gradient, it follows that $\mathbf{q}(\mathbf{x}) = \nabla_x \Theta(\mathbf{x})$ for a suitable potential Θ .

It remains to prove the convexity of Θ . First, we prove that the map $\mathbf{x} \mapsto \mathbf{q}(\mathbf{x}) = \nabla_x \Theta(\mathbf{x})$ is *monotonic*, i.e. by definition, that for any \mathbf{x}_1 and \mathbf{x}_2

$$(\mathbf{x}_2 - \mathbf{x}_1) \cdot [\mathbf{q}(\mathbf{x}_2) - \mathbf{q}(\mathbf{x}_1)] \geq 0. \quad (27)$$

Indeed, should this inequality be violated for some $\bar{\mathbf{x}}_1, \bar{\mathbf{x}}_2$, the continuity of $\mathbf{q}(\mathbf{x})$ would imply that for all $\mathbf{x}_1, \mathbf{x}_2$ close enough to $\bar{\mathbf{x}}_1, \bar{\mathbf{x}}_2$

$$|\mathbf{q}(\mathbf{x}_1) - \mathbf{x}_1|^2 + |\mathbf{q}(\mathbf{x}_2) - \mathbf{x}_2|^2 > |\mathbf{q}(\mathbf{x}_2) - \mathbf{x}_1|^2 + |\mathbf{q}(\mathbf{x}_1) - \mathbf{x}_2|^2. \quad (28)$$

⁹ Note that $\mathbf{x}(\mathbf{q}) = \mathbf{q}$ does not solve the above problem as it violates the latter constraint unless the terminal densities are identical.

This in turn means that if we interchange the destinations of small patches around \bar{x}_1 and \bar{x}_2 , sending them not to the corresponding patches around $q(\bar{x}_1)$ and $q(\bar{x}_2)$ but vice versa, then the value of the functional I will decrease by a small yet positive quantity, and therefore it cannot be minimum for the original map.¹⁰

To complete the argument, observe that convexity of a smooth function $\Theta(x)$ follows if the matrix of its second derivatives $\nabla_{x_i} \nabla_{x_j} \Theta(x)$ is positive definite for all x . Substituting $q(x) = \nabla_x \Theta(x)$ into (27), assuming that x_2 is close to x_1 and Taylor expanding, we find that

$$(x_2 - x_1) \cdot [\nabla_{x_i} \nabla_{x_j} \Theta(x_1)(x_2 - x_1)] \geq 0. \quad (29)$$

As x_2 is arbitrary, this proves the desired positive definiteness and thus establishes the equivalence of the Monge–Ampère equation (19) and of the mass transportation problem with quadratic cost.

This equivalence is actually proved under much weaker conditions, not requiring any smoothness (Brenier 1987, 1991). The proof makes use of the ‘relaxed’ reformulation of the mass transportation problem due to Kantorovich (1942). Instead of solving the highly non-linear problem of finding a map $q(x)$ minimizing the cost (22) with prescribed terminal densities, Kantorovich considered the *linear programming* problem of minimizing

$$\bar{I} \equiv \int_{\mathcal{D}_{\text{in}}} \int_{\mathcal{D}_0} c(q, x) \rho(q, x) d^3 q d^3 x, \quad (30)$$

under the constraint that the joint distribution $\rho(q, x)$ is non-negative and has marginals $\rho_{\text{in}}(q)$ and $\rho_0(x)$, the latter being equivalent to

$$\int_{\mathcal{D}_0} \rho(q, x) d^3 x = \rho_{\text{in}}(q), \quad \int_{\mathcal{D}_{\text{in}}} \rho(q, x) d^3 q = \rho_0(x). \quad (31)$$

Note that if we assume any of the two following forms for the joint distribution:

$$\begin{aligned} \rho(q, x) &= \rho_0(x) \delta(q - q(x)) \\ \rho(q, x) &= \rho_{\text{in}}(q) \delta(x - x(q)), \end{aligned} \quad (32)$$

we find that \bar{I} reduces to the cost I as defined in (22). This relaxed formulation allowed Kantorovich to establish the existence of a minimizing joint distribution.

The relaxed formulation can be used to show that the minimizing solution actually defines a map, which need not be smooth if one or both of the terminal distribution have a singular component (in our case, when mass concentrations are present). The derivation (Brenier 1987, 1991) makes use of the technique of duality (Appendix C2), which will also appear in discussing algorithms (Section 4.2) and reconstruction beyond the potential hypothesis (Section 6).

We have thus shown that the Monge–Kantorovich optimal mass transportation problem can be applied to solving the Monge–Ampère equation. The actual implementation (Section 4), done for a suitable discretization, will be henceforth called Monge–Ampère–Kantorovich.

3.4 Sources of uncertainty in reconstruction

In this section we discuss various sources of non-uniqueness of the MAK reconstruction: multistreaming, collapsed regions, reconstruction from a finite patch of the Universe.

¹⁰ As we shall see in Section 4.1, the converse is not true: monotonicity alone does not imply that the integral I is a minimum; the minimizing map must also be potential.

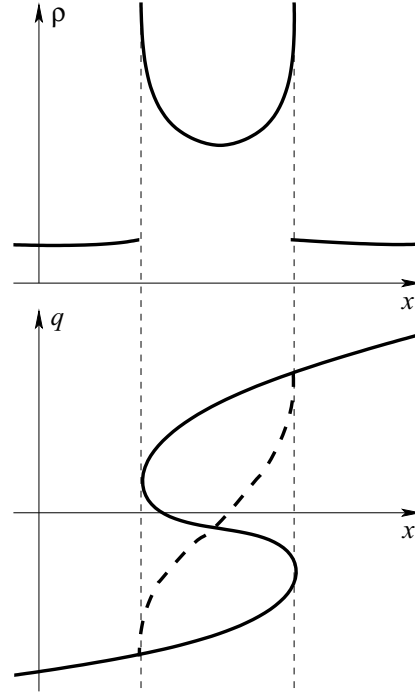


Figure 2. A one-dimensional example of non-unique reconstruction of the Lagrangian map in the presence of multistreaming. The density distribution (upper graph) is generated by a multistreaming Lagrangian map (thick line of lower graph) but may also be generated by a spurious single-stream Lagrangian map (dashed line).

We have stated before that our uniqueness result applies only in so far as we can treat present-epoch high-density multistream regions as if they were truly collapsed, ignoring their width. We now give a simple one-dimensional example of non-uniqueness in which a thick region of multistreaming is present. Fig. 2 shows a multistream Lagrangian map $x(q)$ and the associated density distribution; the inverse map $q(x)$ is clearly multivalued. The same density distribution may, however, be generated by a spurious single-stream Lagrangian map shown in the same figure. There is no way to distinguish between the two inverse Lagrangian maps if the various streams cannot be disentangled.

Suppose now that the present density has a singular part, i.e. there are mass concentrations present that have vanishing (Eulerian) volumes but possess finite masses. Obviously any such object originates from a domain in the Lagrangian space which occupies a finite volume. A one-dimensional example is again helpful. Fig. 3 shows a Lagrangian map in which a whole Lagrangian shock interval $[q_1, q_2]$ has collapsed into a single point on the x -axis. Outside of this point the Lagrangian map is uniquely invertible but the point itself has many antecedents. Note that the graph of the Lagrangian map may be inverted by just interchanging the q and x axes, but its inverse contains a piece of vertical line. The position of the Lagrangian shock interval which has collapsed by the present epoch is uniquely defined by the present mass field but the initial velocity fluctuations in this interval cannot be uniquely reconstructed. In particular, there is no way to know whether collapse has started before the present epoch. We can of course arbitrarily assume that collapse has just happened at the present epoch; if we also suppose that particles have travelled with a constant speed, i.e. use the Zel’dovich/adhesion approximation, then the initial velocity profile within the Lagrangian shock interval will be linear (Fig. 3). Any

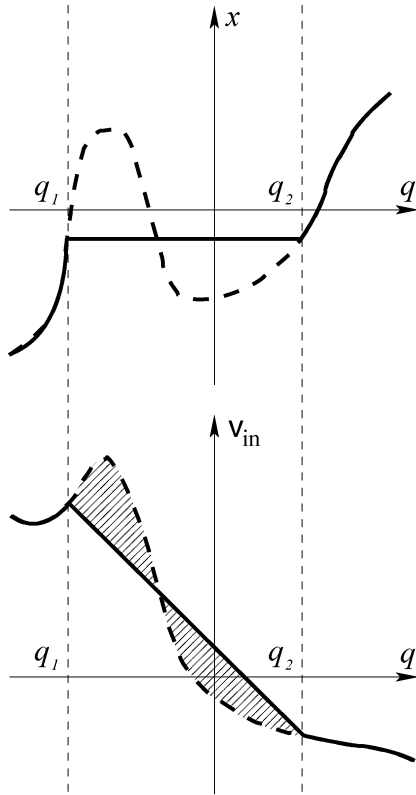


Figure 3. Two initial velocity profiles $v_{in}(q)$ (bottom, solid and dashed lines) leading to the same Lagrangian map $x = q + \tau v_{in}(q)$ (top, solid line) in the adhesion approximation. The Zel'dovich approximation would give multistreaming (top, dashed line). Hatched areas (bottom) are equal in the adhesion dynamics.

other smooth velocity profile joining the same end points would have points where its slope (velocity gradient) is more negative than that of the linear profile (Fig. 3) and thus would have started to collapse before the present epoch (in one dimension caustics appear at the time which is minus the inverse of the most negative initial velocity gradient).

All of this carries over to more than one dimension. The MAK reconstruction gives a unique antecedent for any Eulerian position outside mass concentrations. Each mass concentration in the Eulerian space, taken globally, has a uniquely defined Lagrangian antecedent region but the initial velocity field inside the latter is unknown. In other words, displacement reconstruction is well defined but full reconstruction, based on the Zel'dovich/adhesion approximation for velocities, is possible only outside of mass concentrations (note, however, that velocities in the Eulerian space are still reconstructed at almost all points). We call the corresponding initial Lagrangian domains *collapsed regions*.

Finally, we consider a uniqueness problem arising from knowing the present mass distribution only truncated over a finite Eulerian domain \mathcal{D}_0 , as is necessarily the case when working with a real catalogue. If we also know the corresponding Lagrangian domain \mathcal{D}_{in} and both domains are bounded and convex, then uniqueness is guaranteed (see Section 3.2). What we know for sure about \mathcal{D}_{in} is its volume, which (in our units) is equal to the total mass contained in \mathcal{D}_0 . Its shape and position may, however, be constrained by further information. For example, if we know that the typical displacement of mass elements since decoupling is about 10 Mpc in comoving coordinates (see Section 5) and our data extend over a patch of

typical size 100 Mpc, then there is not more than a 10 per cent uncertainty on the shape of \mathcal{D}_{in} . Additional information on peculiar velocities may also be used to constrain \mathcal{D}_{in} .

Note also that a finite-size patch \mathcal{D}_0 with unknown antecedent \mathcal{D}_{in} will give rise to a unique reconstruction (up to a translation) if we assume that it is surrounded by a uniform background extending to infinity. This is a consequence of the result of Caffarelli & Li mentioned at the end of Section 3.2. The arbitrary linear term in (21) corresponds to a translation; as to the quadratic term, it is constrained by the cosmological principle of isotropy to be exactly $|q|^2/2$.

4 THE MAK METHOD: DISCRETIZATION AND ALGORITHMS

In this section we show how to compute the solution to the Monge–Ampère–Kantorovich problem from the known present density field. First, the problem is discretized into an assignment problem (Section 4.1), then we present some general tools that make the assignment problem computationally tractable (Section 4.2) and finally we present, to the best of our knowledge, the most effective method for solving our particular assignment problem, based on the auction algorithm of D. Bertsekas (Section 4.3), and details of its implementation for the MAK reconstruction (Section 4.4).

4.1 Reduction to an assignment problem

Perhaps the most natural way of discretizing a spatial mass distribution is to approximate it by a finite system of identical Dirac point masses, with possibly more than one mass at a given location. This is compatible both with N -body simulations and with the intrinsically discrete nature of observed luminous matter. Assuming that we have N unit masses both in the Lagrangian and the Eulerian space, we may write

$$\rho_0(\mathbf{x}) = \sum_{i=1}^N \delta(\mathbf{x} - \mathbf{x}_i), \quad \rho_{in}(\mathbf{q}) = \sum_{j=1}^N \delta(\mathbf{q} - \mathbf{q}_j). \quad (33)$$

For discrete densities of this form, the mass conservation constraint in the optimal mass transportation problem (Section 3.3) requires that the map $\mathbf{q}(\mathbf{x})$ induce a one-to-one pairing between positions of the unit masses in the \mathbf{x} and \mathbf{q} spaces, which may be written as a permutation of indices that sends \mathbf{x}_i to $\mathbf{q}_{j(i)}$. Substituting this into the quadratic cost functional (23), we obtain

$$I = \sum_{i=1}^N \frac{|\mathbf{x}_i - \mathbf{q}_{j(i)}|^2}{2}. \quad (34)$$

We thus reduced the problem to the purely combinatorial one of finding a permutation $j(i)$ [or its inverse $i(j)$] that minimizes the quadratic cost function (34).

This problem is an instance of the general *assignment problem* in combinatorial optimization: for a cost matrix c_{ij} , find a permutation $j(i)$ that minimizes the cost function

$$I = \sum_{i=1}^N c_{ij(i)}. \quad (35)$$

As we shall see in the following sections, there exist effective algorithms for finding minimizing permutations.

Before proceeding with the assignment problem, we should mention an alternative approach in which discretization is performed only in the Eulerian space and the initial mass distribution is kept continuous and uniform. Minimization of the quadratic cost function will then give rise to a tessellation of the Lagrangian space

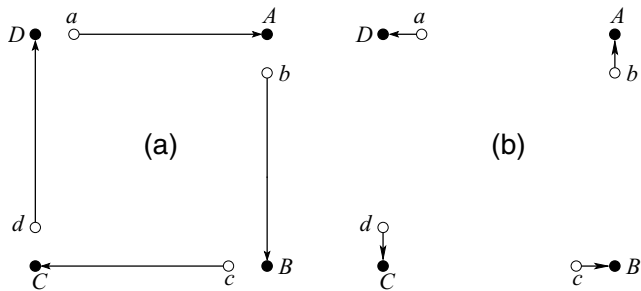


Figure 4. Two monotonic assignments sending white points to black ones: (a) an assignment that is vastly non-optimal in terms of quadratic cost but cannot be improved by any pair interchange; (b) the optimal assignment, shown for comparison.

into polyhedral regions which end up collapsed into the discrete Eulerian Dirac masses. Basically, the reason why these regions are polyhedra is that the convex potential $\Phi(\mathbf{q})$ of the Lagrangian map has a gradient that takes only finitely many values. This problem, which has been studied by Aleksandrov and Pogorelov (see, e.g., Pogorelov 1978), is closely related to Minkowski’s (1897) famous problem of constructing a convex polyhedron with prescribed areas and orientations of its faces (in our setting, areas and orientations correspond to masses and values of the gradient). Uniqueness in the Minkowski problem is guaranteed up to a translation. Starting with Minkowski’s own very elegant solution, various methods of constructing solutions to such geometrical questions have been devised. So far, we have not been able to make use of such ideas in a way truly competitive with discretization in both spaces and then solving the assignment problem.

The solution to our assignment problem (with quadratic cost) has the important property that it is monotonic: for any two Lagrangian positions \mathbf{q}_1 and \mathbf{q}_2 , the corresponding Eulerian positions \mathbf{x}_1 and \mathbf{x}_2 are such that

$$(\mathbf{x}_1 - \mathbf{x}_2) \cdot (\mathbf{q}_1 - \mathbf{q}_2) \geq 0. \quad (36)$$

This is of course the discrete counterpart of (27). In one dimension, when all the Dirac masses are on the same line, monotonicity implies that the leftmost Lagrangian position goes to the leftmost Eulerian position, the second leftmost Lagrangian position to the second leftmost Eulerian position, etc. It is easily checked that this correspondence minimizes the cost (34).

In more than one dimension, a correspondence between Lagrangian and Eulerian positions that is just monotonic will usually not minimize the cost (a simple two-dimensional counterexample is given in Fig. 4).¹¹ Actually, a much stronger condition, called *cyclic monotonicity*, is needed in order to minimize the cost. It requires k -monotonicity for any k between 2 and N ; the latter is defined by taking any k Eulerian positions with their corresponding Lagrangian antecedents and requiring that the cost (34) should not decrease under an arbitrary reassignment of the Lagrangian positions within the set of Eulerian positions taken. Note that the usual monotonicity corresponds to 2-monotonicity (stability with respect to pair exchanges).

A strategy called the path-interchange Zel’dovich approximation (PIZA) for constructing monotonic correspondences between

Lagrangian and Eulerian positions has been proposed by Croft & Gaztañaga (1997). In PIZA, a randomly chosen tentative correspondence between initial and final positions is successively improved by swapping randomly selected *pairs* of initial particles whenever (36) is not satisfied. After the cost (34) ceases to decrease between iterations, an approximation to a monotonic correspondence is established, which is generally neither unique, as already observed by Valentine, Saunders & Taylor (2000) in testing PIZA reconstruction, nor optimal. We shall return to this in Sections 5 and 7.3.

4.2 Nuts and bolts of solving the assignment problem

For a general set of N unit masses, the assignment problem with the cost function (34) has a single solution which can obviously be found by examining all $N!$ permutations. However, unlike computationally hard problems, such as the travelling salesman problem, the assignment problem can be handled in ‘polynomial time’ – actually in not more than $O(N^3)$ operations. All methods achieving this use a so-called dual formulation of the problem, based on a relaxation similar to that applied by Kantorovich to the optimal mass transportation (Section 3.3; a brief introduction to duality is given in Appendix C2). In this section we explain the basics of this technique, using a variant of a simple mechanical model introduced in a more general setting by Hénon (1995, 2002).

Consider the general assignment problem of minimizing the cost (35) over all permutations $j(i)$. We replace it by a ‘relaxed’, linear programming problem of minimizing

$$\tilde{I} = \sum_{i,j=1}^N c_{ij} f_{ij}, \quad (37)$$

where auxiliary variables f_{ij} satisfy

$$f_{ij} \geq 0, \quad \sum_{k=1}^N f_{kj} = \sum_{k=1}^N f_{ik} = 1 \quad (38)$$

for all i, j , an obvious discrete analogue of (31). We show now that it is possible to build a simple mechanical device (Fig. 5), which solves

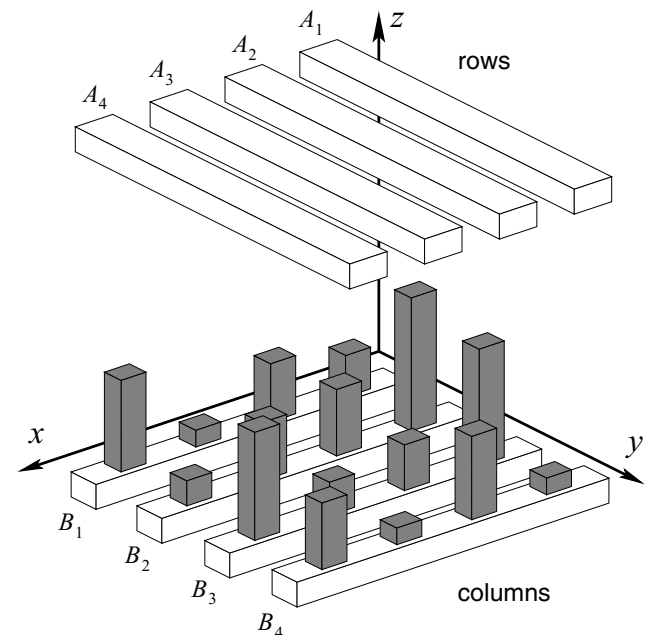


Figure 5. An analogue computer solving the assignment problem for $N = 4$.

¹¹ Note that in one dimension, in the continuous case, any map is a gradient and we have already observed in Section 3.3 that if a gradient map is monotonic it is the gradient of a convex function.

this relaxed problem and that the solution will in fact determine a minimizing permutation in the original assignment problem (i.e. for any i or j fixed, only one f_{ij} will be unity and all others zero). The device acts as an *analogue computer*: the numbers involved in the problem are represented by physical quantities, and the equations are replaced by physical laws.

Define coordinate axes x, y, z in space, with the z -axis vertical. We take two systems of N horizontal rods, parallel to the x and y axes, respectively, and call them *columns* and *rows*, referring to columns and rows of the cost matrix. Each rod is constrained to move in a corresponding vertical plane while preserving the horizontal orientation in space. For a row rod A_i , we denote the z coordinate of its bottom face by α_i and for a column rod B_j , we denote the z coordinate of its top face β_j . Row rods are placed above column rods, therefore $\alpha_i \geq \beta_j$ for all i, j (see Fig. 5).

Upper (row) rods are assumed to have unit weight, and lower (column) rods to have negative unit weight, or unit ‘buoyancy’. Therefore, both groups of rods are subject to gravitational forces pulling them together. However, this movement is obstructed by N^2 small vertical studs of negligible weight put on column rods just below row rods. A stud placed at projected intersection of column B_j and row A_i has length $C - c_{ij}$ with a suitably large positive constant C and thus constrains the quantities α_i and β_j to satisfy the stronger inequality

$$\alpha_i - \beta_j \geq C - c_{ij}. \quad (39)$$

The potential energy of the system is, up to a constant,

$$U = \sum_{i=1}^N \alpha_i - \sum_{j=1}^N \beta_j. \quad (40)$$

In linear programming, the problem of minimizing (40) under the set of constraints given by (39) is called the *dual problem* to the ‘relaxed’ one (37)–(38) (see Appendix C2); the α and β variables are called the *dual variables*.

The analogue computer does in fact solve the dual problem. Indeed, first hold the two groups of rods separated from each other and then release them, so that the system starts to evolve. Rows will go down, columns will come up, and contacts will be made with the studs. Aggregates of rows and columns will be progressively formed and modified as new contacts are made, giving rise to a complex evolution. Eventually the system reaches an equilibrium, in which its potential energy (40) is minimum and all constraints (39) are satisfied (Hénon 2002). Moreover, it may be shown that the solution to the original problem (37)–(38) is expressible in terms of the forces exerted by the rods on each other at equilibrium and is typically a one-to-one correspondence between the A_i s and the B_j s (for details, see Appendix C3).

The common feature of many existing algorithms for solving the assignment problem, which makes them more effective computationally than the simple enumeration of all $N!$ permutations, is the use of the intrinsically continuous, geometric formulation in terms of the pair of linear programming problems (37)–(38) and (40)–(39). The mechanical device provides a concrete model for this formulation; in fact, assignment algorithms can be regarded as descriptions of specific procedures to make the machine reach its equilibrium state.¹² An introduction to algorithmic aspects of solving the assignment problem, including a proof of the $O(N^3)$ theoretical bound

on the number of operations, based on the Hungarian method of Kuhn (1955), may be found in Papadimitriou & Steiglitz (1982).

In spite of the general $O(N^3)$ theoretical bound, various algorithms may show very different performance when applied to a specific optimization problem. During the preparation of the earlier publication (Frisch et al. 2002) the dual simplex method of Balinski (1986) was used, with some modifications inspired by algorithm B of Hénon (2002). Several other algorithms were tried subsequently, including an adaptation of algorithm A of the latter reference and the algorithm of Burkard & Derigs (1980), itself based on the earlier work of Tomizawa (1971). For the time being, the fastest running code by far is based on the auction algorithm of Bertsekas (1992, 2001), arguably the most effective of existing ones, which is discussed in the next section. Needless to say, all of these algorithms arrive at the same solution to the assignment problem with given data but can differ by several orders of magnitude in the time it takes to complete the computation.

4.3 The auction algorithm

We explain here the essence of the auction algorithm in terms of our mechanical device.¹³ Note that the original presentation of this algorithm (Bertsekas 1981, 1992, 2001) is based on a different perspective, that of an *auction*, in which the optimal assignment appears as an economic rather than a mechanical equilibrium; the interested reader will benefit much from reading these papers.

Put initially the column rods at zero height and all row rods well above them, so that no contacts are made and constraints (39) are satisfied. To decrease the potential energy, now let the row rods descend while keeping the column rods fixed. Eventually all row rods will meet studs placed on column rods and stop. Some column rods may then come in contact with multiple row rods. Such rods are overloaded: if they were not prevented from moving they would descend.

Note that at this stage any column rod A_i has established a contact with a row rod B_j for which the stud length $C - c_{ij}$ is the maximum and the cost c_{ij} the minimum among other B s; for $c_{ij} = |\mathbf{x}_i - \mathbf{q}_j|^2/2$, this means that any Eulerian position \mathbf{x}_i is coupled to its nearest Lagrangian neighbour \mathbf{q}_j . This coupling is a reasonable guess for the optimal assignment; should it happen to be one-to-one, then the equilibrium, and with it the optimal assignment, would be reached. It is usually not, so there are overloaded B rods and the following procedure is applied to find a compromise between minimization of the total cost and the requirement of one-to-one correspondence.

Take any overloaded rod B_j and let it descend while keeping other column rods fixed. As B_j descends, row rods touching it will follow its motion until they meet studs of other column rods and stay behind. The downward motion of B_j is stopped only when the last row rod touching B_j is about to lose its contact. We then turn to any other overloaded column rod and repeat the procedure as often as needed.

This general step can be viewed as an *auction* in which row rods bid for the descending column rod, offering prices equal to decreases in their potential energy as they follow its way down. As the column rod descends, thereby increasing its price, the auction is won by the row rod able to offer the largest *bidding increment*, i.e. to decrease its potential energy by the largest amount while not violating the constraints posed by studs of the rest of column rods. For computational purposes it suffices to compute bidding increments for all

¹² This applies to algorithms that never violate constraints (39) represented by studs; all practical assignment algorithms known to us fall within this category.

¹³ A movie illustrating the subsequent discussion may be found at <http://www.obs-nice.fr/etc7/movie.html> (requires fast Internet access).

competing row rods from the dual α and β variables and assign the descending column rod B_j to the highest bidder A_i , decreasing their heights β_j and α_i correspondingly.

Observe that, at each step, the total potential energy U defined by (40) decreases by the largest amount that can be achieved by moving the descending column rod without violating the constraints.¹⁴ Since (40) is obviously non-negative, the descent cannot proceed indefinitely, and the process may be expected to converge quite fast to a one-to-one pairing that solves the assignment problem.

However, as observed by Bertsekas (1981, 1992, 2001), this ‘naive’ auction algorithm may end up in an infinite cycle if several row rods bid for a few equally favourable column rods, thus having zero bidding increments. To break such cycles and also to accelerate convergence, a perturbation mechanism is introduced in the algorithm. Namely, the constraints (39) are replaced by weaker ones

$$\alpha_i - \beta_j \geq C - c_{ij} - \epsilon \quad (41)$$

for a small positive quantity ϵ , and in each auction the descending column rod is pushed down by ϵ in addition to decreasing its height by the bidding increment. It can be shown that this reformulated process terminates in a finite number of rounds; moreover, if all stud lengths are integer and ϵ is smaller than $1/N$, then the algorithm terminates at an assignment that is optimal in the unperturbed problem (Bertsekas 1992).

The third ingredient in the Bertsekas algorithm is the idea of ϵ -scaling. When the values of dual variables are already close to the solution of the dual problem, it usually takes relatively few rounds of the auction to converge to a solution. Thus one can start with large ϵ to compute a rough approximation for dual variables quickly, without worrying about the quality of the assignment, and then proceed reducing ϵ in geometric progression until it passes the $1/N$ threshold, ensuring that the assignment thus achieved solves the initial problem.

Bertsekas’ algorithm is especially fast for *sparse assignment problems*, in which rods A_i and B_j can be matched only if the pair (i, j) belongs to a given subset \mathcal{A} of the set of N^2 possible pairs. We call such pairs *valid* and define the *filling factor* to be the proportion of valid pairs $f = |\mathcal{A}|/N^2$. When this factor is small, computation can be considerably faster: to find the bidding increment for a rod A_i , we only need to run over the list of rods B_j such that (i, j) is a valid pair.

Note also that the decentralized structure of the algorithm facilitates its parallelization (see references in Bertsekas 1992, 2001).

4.4 The auction algorithm for the MAK reconstruction

We now describe the adaptation of the auction algorithm to the MAK reconstruction. Experiments with various programs contained in Bertsekas’ publicly available package (<http://web.mit.edu/dimitrib/www/auction.txt>) showed that the most effective for our problem is AUCTION_FLP. It assumes integer costs c_{ij} , which in our case requires

¹⁴ This idea of moving a rod, or adjusting a dual variable, up to the last point compatible with all the constraints, may be actually implemented in a number of ways, giving rise to several possible flavours of the auction algorithm. For example, the above procedure in its most effective implementation requires a parallel computer so that groups of several rods can be tracked simultaneously. On sequential computers another, less intuitive procedure, in which upper rods are dropped one at a time, proves more effective (Bertsekas 1992).

proper scaling of the cost matrix. To achieve this, the unit of length is adjusted so that the size of the reconstruction patch equals 100, and then the square of the distance between an initial and a final position is rounded off to an integer. In our application, row and column rods correspond to Eulerian and Lagrangian positions, respectively. As the MAK reconstruction is planned for application to catalogues of 10^5 and more galaxies, we do not store the cost matrix, which would require an $O(N^2)$ storage space, but rather compute its elements on demand from the coordinates, which requires only $O(N)$ space.

Our problem is naturally adapted for a sparse description if galaxies travel only a short distance compared with the dimensions of the reconstruction patch. For instance, in the simulation discussed in Section 5, the rms distance travelled is only about $10 h^{-1}$ Mpc, or 5 per cent of the size of the simulation box, and the largest distance travelled is about 15 per cent of this size. So we may assume that in the optimal assignment distances between paired positions will be limited. We define then a critical distance d_{crit} and specify that a final position x_i and an initial position q_j form a valid pair only if they are within less than d_{crit} from each other. This critical distance must be adjusted carefully: if it is too small, we risk excluding the optimal assignment; if it is taken too large, the benefit of the sparse description is lost.

However, the saving in computing time achieved by sparse description has to be paid for in storage space: to store the set \mathcal{A} of valid pairs, storage of size $|\mathcal{A}| = fN^2$ is needed, which takes us back to the $O(N^2)$ storage requirement. We have explored two solutions to this problem.

(1) Use a *dense* description nevertheless, i.e. the one where all pairs (i, j) are valid and there is no need to store the set \mathcal{A} . The auction program is easily adapted to this case (in fact this simplifies the code). However, we forfeit the saving in time provided by the sparse structure.

(2) The sparse description can be preserved if the set of valid pairs is computed on demand rather than stored. This is easy if initial positions fill a uniform cubic grid, the simplest discrete approximation to the initial quasi-uniform distribution of matter in the reconstruction problem. Thus, for a given final position x_i , the valid pairs correspond to points of the cubic lattice that lie inside a sphere of radius d_{crit} centred at x_i , so their list can be generated at run time.

Fig. 6 gives the computing time as a function of the number of points N used in the assignment problem. Shown are the dense and sparse versions of the auction algorithm (in the latter, the critical distance squared was taken equal to 200) and the Burkard & Derigs (1980) algorithm, which ranked the next fastest in our experiments. The N initial and final positions are chosen from the file generated by an N -body simulation described in Section 5; the choice is random except for the sparse algorithm, in which the initial positions are required to fill a cubic lattice. Hence, the performance of the sparse auction algorithm shown in the figure is not completely comparable to that of the two other algorithms.

It is evident that the difference in computing time between the dense auction and the Burkard & Derigs algorithms steadily increases. In the vicinity of $N = 10^5$, the dense auction algorithm is about 10 times faster than the other one. For the sparse version, the decrease in computing time is spectacular: as could be expected, the ratio of computing times for the two versions of the auction algorithm is of the order of f . For large N , the $O(N^3)$ asymptotic of the computing time is quite clear for the sparse auction algorithm. For two other algorithms, a similar asymptotic was found for larger N in other experiments (not shown).

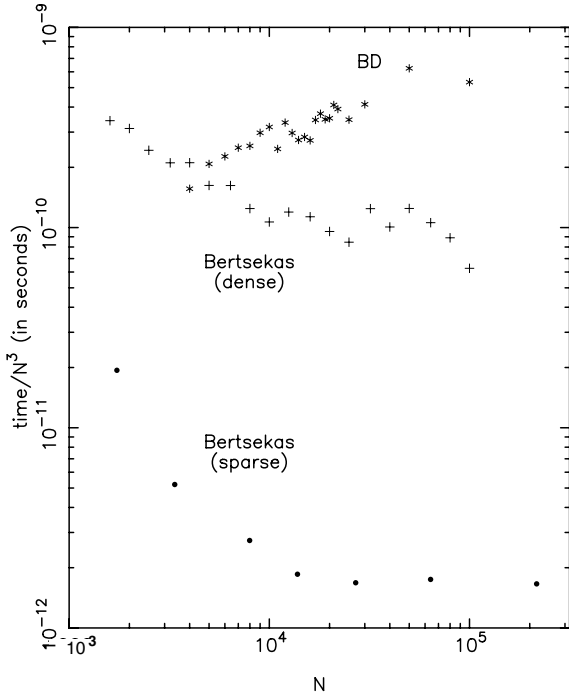


Figure 6. Computing time for different algorithms as a function of the number N of points (divided by N^3 for normalization). Asterisks, the Burkard & Derigs (1980) algorithm (BD); crosses and points, the dense and sparse versions of the auction algorithm (described in the text).

In all three cases shown, the initial positions fill a constant volume while N is varied. This is what we call *constant-volume computations*. In the sparse case, this results in a constant filling factor, equal to the ratio of the volume of the sphere with radius d_{crit} to the volume occupied by the initial positions. Here this filling factor is about $f = 0.019$. Another choice, not shown in the figure, is that of *constant-density computations*, when the initial positions are taken from a volume for which the size increases with N . In this case the time dependence of algorithms for large N is of the order of $N^{3/2}$.

We finally observe that the sparse auction algorithm applied to the MAK reconstruction requires 5 h of single-processor CPU time on a 667-MHz Compaq/DEC Alpha machine for 216 000 points.

5 TESTING THE MAK RECONSTRUCTION

In this section we present results of our testing the MAK reconstruction against data of cosmological N -body simulations. In a typical simulation of this kind, the dark matter distribution is approximated by N particles of identical mass. Initially the particles are put on a uniform cubic grid and given velocities that form a realization of the primordial velocity field for which the statistics is prescribed by a certain cosmological model. Trajectories of particles are then computed according to the Newtonian dynamics in a comoving frame, using periodic boundary conditions. The reconstruction problem is therefore to recover the pairing between the initial (Lagrangian) positions of the particles and their present (Eulerian) positions in the N -body simulation, knowing only the set of computed Eulerian positions in the physical space.

We test our reconstruction against a simulation of 128^3 particles in a box of $200 h^{-1}$ Mpc size (where h is the Hubble parameter in units of $100 \text{ km s}^{-1} \text{ Mpc}^{-1}$) performed using the adaptive P^3M code

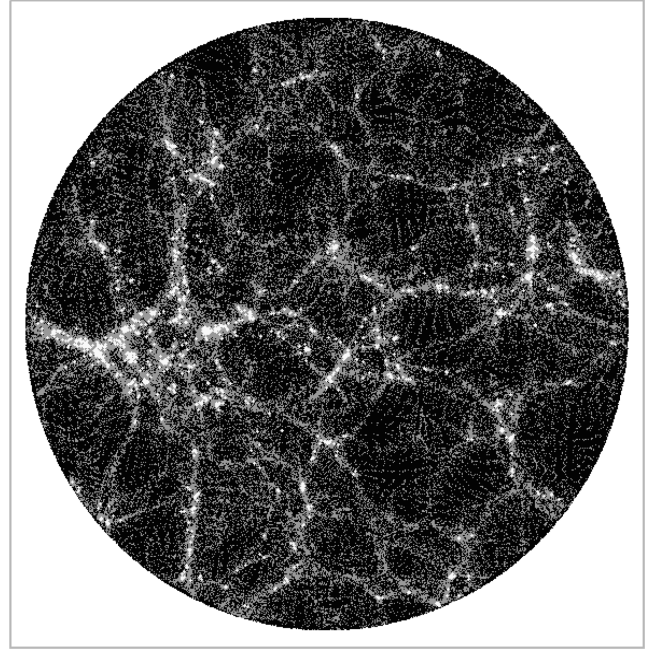


Figure 7. N -body simulation output in the Eulerian space used for testing our reconstruction method (shown is a projection on to the x - y plane of a 10 per cent slice of the simulation box of size $200 h^{-1}$ Mpc). Points are highlighted in yellow when reconstruction fails by more than $6.25 h^{-1}$ Mpc, which happens mostly in high-density regions. This figure is available in colour in the on-line version of the journal on *Synergy*.

HYDRA (Couchman, Thomas & Pearce 1995).¹⁵ A Λ CDM cosmological model is used with parameters $\Omega_m = 0.3$, $\Omega_\Lambda = 0.7$, $h = 0.65$, $\sigma_8 = 0.9$.¹⁶ The values of these parameters within the model are determined by fitting the observed cosmic microwave background (CMB) spectrum.¹⁷ The output of the N -body simulation is illustrated in Fig. 7 by a projection on to the x - y plane of a 10 per cent slice of the simulation box.

Since the simulation assumes periodic boundary conditions, some Eulerian positions situated near boundaries may have their Lagrangian antecedents at the opposite side of the simulation box. Suppressing the resulting spurious large displacements is crucial for successful reconstruction. Indeed, for a typical particle displacement of $1/20$ of the box size, spurious box-wide leaps of 1 per cent of the particles will generate a contribution to the quadratic cost (34) four times larger than that of the rest. To suppress such leaps, for each Eulerian position that has its antecedent Lagrangian position at the other side of the simulation box, we add or subtract the box size from coordinates of the latter (in other words, we are considering

¹⁵ In a flavour of N -body codes called particle-mesh (PM) codes, Newtonian forces acting on particles are interpolated from the gravitational field computed on a uniform mesh. In very dense regions, precision is increased by adaptively refining the mesh and by direct calculation of local particle-particle (PP) interactions; codes of this type are correspondingly called *adaptive P^3M* .

¹⁶ The use of a Λ CDM model instead of the model without a cosmological constant (Appendix A) leads to some modifications in the basic equations but does not change formulae used for the MAK reconstruction.

¹⁷ Data of the first year *Wilkinson Microwave Anisotropy Probe* (Spergel et al. 2003; see also Bridle et al. 2003) suggest a value $\sigma_8 = 0.84 \pm 0.04$, marginally smaller than that used here. This may slightly extend the range of scales favourable for the MAK reconstruction.

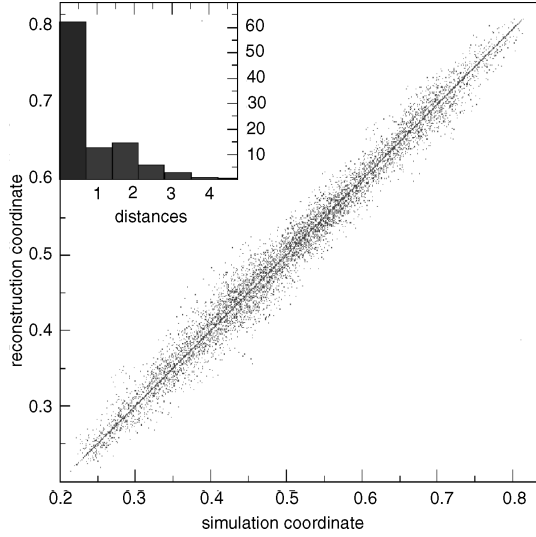


Figure 8. Test of the MAK reconstruction for a sample of $N' = 17\,178$ points initially situated on a cubic grid with mesh $\Delta x = 6.25 h^{-1}$ Mpc. The scatter diagram plots true versus reconstructed initial positions using a quasi-periodic projection which ensures one-to-one correspondence with points on the cubic grid. The histogram inset gives the distribution (in percentages) of distances between true and reconstructed initial positions; the horizontal unit is the sample mesh. The width of the first bin is less than unity to ensure that only exactly reconstructed points fall in it. Note that more than 60 per cent of the points are exactly reconstructed.

the distance on a torus). In what follows we refer to this procedure as the *periodicity correction*.

We first present reconstructions for three samples of particles initially situated on Lagrangian subgrids with meshes given by $\Delta x = 6.25 h^{-1}$ Mpc, $\Delta x/2$ and $\Delta x/4$. To further reduce possible effects of the unphysical periodic boundary condition, we truncate the data by discarding those points for which Eulerian positions are not within the sphere of radius $16\Delta x$ placed at the centre of the simulation box (for the largest Δx its diameter coincides with the box size). The problem is then confined to finding the pairing between the remaining Eulerian positions and the set of

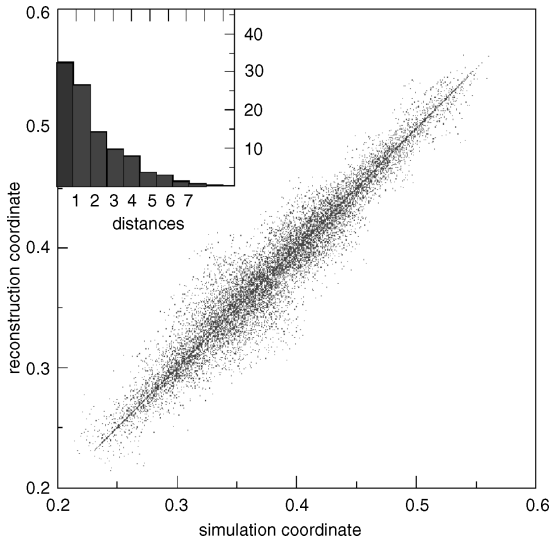


Figure 9. Same as in Fig. 8 but with $N' = 19\,187$ and a sample mesh of $\Delta x/2 = 3.125 h^{-1}$ Mpc. Exact reconstruction is down to 35 per cent.

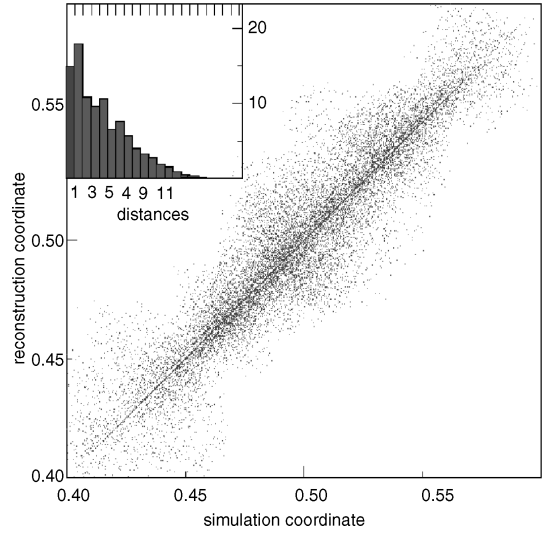


Figure 10. Same as in Fig. 8 but with $N' = 23\,111$ and a sample mesh of $\Delta x/4 = 1.56 h^{-1}$ Mpc. Exact reconstruction is down to 14 per cent.

their periodicity-corrected Lagrangian antecedents in the N -body simulation.

The results are shown in Figs 8–11. The main plots show the scatter of reconstructed versus simulation Lagrangian positions for the same Eulerian positions. For these diagrams we introduce a ‘quasi-periodic projection’

$$\tilde{q} \equiv (q_1 + \sqrt{2}q_2 + \sqrt{3}q_3)/(1 + \sqrt{2} + \sqrt{3}) \quad (42)$$

of the vector \mathbf{q} , which ensures a one-to-one correspondence between \tilde{q} -values and points on the regular Lagrangian grid. The insets are histograms (by percentage) of distances, in reconstruction mesh units, between the reconstructed and simulation Lagrangian positions; the first darker bin, slightly less than one mesh in width, corresponds to perfect reconstruction (thereby allowing a good determination of the peculiar velocities of galaxies).

With the mesh size Δx , Lagrangian positions of 62 per cent of the sample of 17 178 points are reconstructed perfectly and about

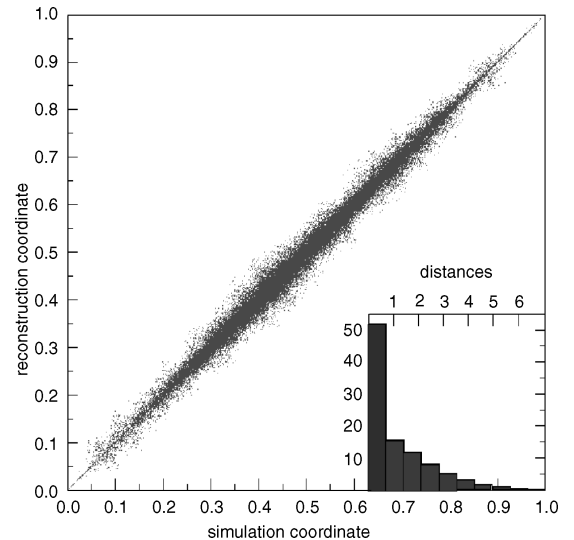


Figure 11. Same as in Fig. 8 with $N' = 10^5$ points selected at random, neighbouring points being typically $4.35 h^{-1}$ Mpc apart. Exact reconstruction is in excess of 50 per cent.

75 per cent are placed within not more than one mesh. With the $\Delta x/2$ grid, we still have 35 per cent exact reconstruction out of 19 187 points, but only 14 per cent for the $\Delta x/4$ grid with 23 111 points.

We also performed a reconstruction on a random sample of 100 000 Eulerian positions taken with their periodicity-corrected Lagrangian antecedents out of the whole set of 128^3 particles, without any restrictions. This reconstruction, with the effective mesh size (average distance between neighbouring points) of $4.35 h^{-1}$ Mpc, gives 51 per cent perfect reconstruction (Fig. 11).

We compared these results with those of the PIZA reconstruction method (see Section 4.1 and Croft & Gaztañaga 1997), which gives a 2-monotone but not necessarily optimal pairing between Lagrangian and Eulerian positions. We applied the PIZA method on the Δx grid and obtained typically 30–40 per cent exactly reconstructed positions, but severe non-uniqueness: for two different seeds of the random generator used to set up the initial tentative assignment, only about half of the exactly reconstructed positions were the same (see figs 3 and 7 of Mohayaee et al. 2003 for an illustration). We also implemented a modification of the PIZA method establishing 3-monotonicity (monotonicity with respect to interchanges of three points instead of pairs) and checked that it does not give a significant improvement over the original PIZA.

In comoving coordinates, the typical displacement of a mass element is about $1/20$ of the box size, that is about $10 h^{-1}$ Mpc. This is not much larger than the coarsest grid of $6.25 h^{-1}$ Mpc used in testing MAK which gave 62 per cent exact reconstruction. Nevertheless, there are 18 other grid points within $10 h^{-1}$ Mpc of any given grid point, so that this high percentage cannot be trivially explained by the smallness of the displacement. Note that without the periodicity correction, the percentage of exact reconstruction for the coarsest grid degraded significantly (from 62 to 45 per cent) and the resulting cost was far from the true minimum.

For real catalogues, reconstruction has to be performed for galaxies with positions that are specified in the *redshift space*, where they appear to be displaced radially (along the line of sight) by an amount proportional to the radial component of the peculiar velocity. Thus, at the present epoch, the redshift position s of a mass element situated at the point x in the physical space is given by

$$s = x + \hat{x}\beta(v \cdot \hat{x}), \quad (43)$$

where v is the peculiar velocity in the comoving coordinates x and the linear growth factor time τ , \hat{x} denotes the unit normal in the direction of x , and the parameter β equals 0.486 in our Λ CDM model.

Following Valentine et al. (2000, see also Monaco & Efstathiou 1999), we use the Zel'dovich approximation to render our MAK quadratic cost function in the s variable. As follows from (11), in this approximation the peculiar velocity is given by

$$v = \frac{1}{\tau}(x - q). \quad (44)$$

At the present time, since $\tau_0 = 1$, this together with (43) gives

$$(s - q) \cdot \hat{x} = (1 + \beta)(x - q) \cdot \hat{x}, \quad (45)$$

$$|s - q|^2 = |x - q|^2 + \beta(\beta + 2)[(x - q) \cdot \hat{x}]^2. \quad (46)$$

Now combining these two equations and using the fact that, by (43), the vectors x and s are collinear and therefore $\hat{x} = \pm \hat{s}$, we may write the quadratic cost function as

$$\frac{1}{2}|x - q|^2 = \frac{1}{2}|s - q|^2 - \frac{\beta(\beta + 2)}{2(\beta + 1)^2}[(s - q) \cdot \hat{s}]^2. \quad (47)$$

The redshift-space reconstruction is then in principle reduced to the physical-space reconstruction. Note, however, that the redshift

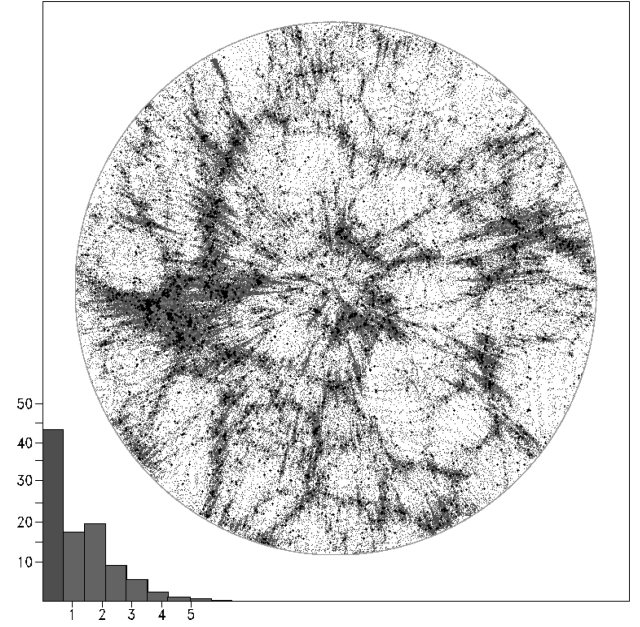


Figure 12. Test of the redshift-space variant of the MAK reconstruction based on the same data as Fig. 8. The circular redshift map (violet points) corresponds to the same physical-space slice as displayed in Fig. 7 (the observer is taken at the centre of the simulation box). Points are highlighted in red when reconstruction fails by more than one mesh. This figure is available in colour in the on-line version of the journal on *Synergy*.

transformation of Eulerian positions may fail to be one-to-one if the peculiar component of velocity field in the proper space coordinates exceeds the Hubble expansion component. This undermines the simple reduction outlined above for catalogues confined to small distances.

We have performed a MAK reconstruction with the redshift-modified cost function (47). The redshift positions were computed for the simulation data with peculiar velocities smoothed over a sphere with radius of $1/100$ of the box size ($2 h^{-1}$ Mpc). This reconstruction led to 43 per cent of exactly reconstructed positions and 60 per cent which are within not more than one Δx mesh from their correct positions (see Fig. 12; a scatter diagram is omitted because it is quite similar to that in Fig. 8). A comparison of the redshift-space MAK reconstruction with the physical-space MAK reconstruction shows that almost 50 per cent of exactly reconstructed positions correspond to the same points. This test shows that the MAK method is robust with respect to systematic errors introduced by the redshift transformation.

Our results demonstrate the essentially potential character of the Lagrangian map above $\sim 6 h^{-1}$ Mpc (within the Λ CDM model) and perhaps at somewhat smaller scales.

Although it is not our intention in this paper to actually implement the MAK reconstruction on real catalogues, a few remarks are in order. The effect of the catalogue selection function can be handled by standard techniques; for instance, one can assign each galaxy a ‘mass’ inversely proportional to the catalogue selection function (Nusser & Branchini 2000; Valentine et al. 2000; Branchini, Eldar & Nusser 2002). Biasing can be taken into account in a similar manner (Nusser & Branchini 2000). These modifications as well as the natural scatter of masses in the observational catalogues require that massive objects be represented by clusters of multiple Eulerian points of unit mass (with the correspondingly increased number of points on a finer grid in the Lagrangian space), which reduces the

problem to a variant of the usual assignment. We also observe that real catalogues involve truncation, that is data available only over a finite region. As already discussed in Section 3.4, this is not a serious problem provided a sufficiently large patch is available. Actually, as noted earlier in this section, the data used in testing have been truncated spherically, without significantly affecting the quality of the reconstruction.

In the redshift-space modification, more accurate determination of peculiar velocities can be done using second-order Lagrangian perturbation theory. Note also that, for the observational catalogues, the motion of the local group itself should also be accounted for (Taylor & Valentine 1999).

6 RECONSTRUCTION OF THE FULL SELF-GRAVITATING DYNAMICS

The MAK reconstruction discussed in Sections 3 and 4 was performed under the assumption of a potential Lagrangian map and of the absence of multistreaming. The tests presented in Section 5 indicate that potentiality works well at scales above $6 h^{-1}$ Mpc, whereas multistreaming is mostly believed to be unimportant above a few Mpc. There could thus remain a substantial range of scales over which the quality of the reconstruction can be improved by relaxing the potentiality assumption and using the full self-gravitating dynamics. Here we show that, as long as the dynamics can be described by a solution to the Euler–Poisson equations, the prescription of the present density field still determines a unique solution to the full reconstruction problem. We give only the main ideas, technical details being left for Appendix D (a mathematically rigorous proof may be found in Loeper 2003). In order to make the exposition self-contained, we also give in Appendix C an elementary introduction to convexity and duality, which are used for the derivation (and also elsewhere in this paper).

We shall start from an Eulerian variational formulation of the Euler–Poisson equations in an Einstein–de Sitter universe, which is an adaptation of a variational principle given by Giavalisco et al. (1993). We minimize the action

$$I = \frac{1}{2} \int_0^{\tau_0} d\tau \int d^3x \tau^{3/2} \left(\rho |v|^2 + \frac{3}{2} |\nabla_x \varphi_g|^2 \right), \quad (48)$$

under the following four constraints: the Poisson equation (3), the mass conservation equation (2) and the boundary conditions that the density field be unity at $\tau = 0$ and prescribed at the present time $\tau = \tau_0$. The constraints can be handled by the standard method of Lagrange multipliers (here functions of space and time), which allows one to vary independently the fields ρ , φ_g and v . The vanishing of the variation in v gives $v = \tau^{-3/2} \nabla_x \theta$, where $\theta(x, \tau)$ is the Lagrange multiplier for the mass conservation constraint. Hence, the velocity is curl-free. The vanishing of the variation in ρ then gives

$$\partial_\tau \theta + \frac{1}{2\tau^{3/2}} |\nabla_x \theta|^2 + \frac{3}{2\tau} \psi = 0. \quad (49)$$

By taking the gradient, this equation goes over into the momentum equation (1), repeated here for convenience:

$$\partial_\tau v + (v \cdot \nabla_x) v = -\frac{3}{2\tau} (v + \nabla_x \varphi_g). \quad (50)$$

It is noteworthy that, if in the action we replace $3/2$ both in the exponent of τ and in the gravitational energy term by $3\alpha/2$, we obtain (50) but also with a $3\alpha/(2\tau)$ factor in the right-hand side. The Zel’dovich approximation and the associated MAK reconstruction

amount clearly to setting $\alpha = 0$, so as to recover the ‘free-streaming action’

$$I = \frac{1}{2} \int_0^{\tau_0} d\tau \int d^3x \rho |v|^2, \quad (51)$$

the minimization of which is easily shown to be equivalent to that of the quadratic cost function (23).

Assuming the action (48) to be finite, existence of a minimum is mostly a consequence of the action being manifestly non-negative. Here it is interesting to observe that the Lagrangian, which is the *difference* between the kinetic energy and the potential energy, is positive, whereas the Hamiltonian which is their *sum* does not have a definite sign. Consequently, our two-point boundary problem is, as we shall see, well posed but the initial-value problem for the Euler–Poisson system is not well posed since formation of caustics after a finite time cannot be ruled out.¹⁸

Does the variational formulation imply uniqueness of the solution? This would be the case if the action were a strictly convex functional (see Appendix C1), which is guaranteed to have one and only one minimum. The action as written in (48) is not convex in the ρ and v variables, but can be rendered so by introducing the mass flux $J = \rho v$; the kinetic energy term then becomes $|J|^2/(2\rho)$, which is convex in the J and ρ variables.

Strict convexity is particularly cumbersome to establish, but there is an alternative way, known as duality: by a Legendre-like transformation the variational problem is carried into a dual problem written in terms of dual variables; the minimum value for the original problem is the maximum for the dual problem. It turns out that the difference of these equal values can be rewritten as a sum of non-negative terms, each of which must thus vanish. This is then used to prove (i) that the difference between any two solutions to the variational problem vanishes and (ii) that any curl-free solution to the Euler–Poisson equations with the prescribed boundary conditions for the density also minimizes the action. All of this together establishes uniqueness. For details see Appendix D.

Several of the issues raised in connection with the MAK reconstruction appear in almost the same form for the Euler–Poisson reconstruction. First, we are faced again with the problem that, when reconstructing from a finite patch of the present Universe, we need either to know the shape of the initial domain or to make some hypothesis as to the present distribution of matter outside this patch. Secondly, just as for the MAK reconstruction, the proof of uniqueness still holds when the present density $\rho_0(x)$ has a singular part, i.e. when some matter is concentrated. Again, we shall have full information on the initial shape of collapsed regions but not on the initial fluctuations inside them. The particular solution obtained from the variational formulation is the only solution that stays smooth for all times prior to τ_0 .

We also note that, at this moment and probably for quite some time, 3D catalogues sufficiently dense to allow reconstruction will be limited to fairly small redshifts. Eventually, it will, however, become of interest to perform reconstruction ‘along our past light-cone’ with data not all at τ_0 . The variational approach can in principle be adapted to handle such reconstruction.

In previous sections we have seen how to implement reconstruction using MAK, which is equivalent to using the simplified action (51). Implementation using the full Euler–Poisson action (48) is mostly beyond the scope of this paper, but we shall indicate some

¹⁸ If we had considered electrostatic repulsive interactions the conclusions would be reversed.

possible directions. In principle it should be possible to adapt to the Euler–Poisson reconstruction the method of the augmented Lagrangian which has been applied to the two-dimensional Monge–Ampère equation (Benamou & Brenier 2000). An alternative strategy, which allows reduction to MAK-type problems, uses the idea of ‘kicked burgulence’ (Bec, Frisch & Khanin 2000) in which, in order to solve the one- or multidimensional Burgers equation

$$\partial_\tau v + (v \cdot \nabla_x)v = f(x, \tau), \quad v = -\nabla_x \varphi_v, \quad (52)$$

one approximates the force by a sum of delta-functions in time:

$$f(x, \tau) \approx \sum_i \delta(\tau - \tau_i) g_i(x). \quad (53)$$

In the present case, $g_i(x)$ are proportional to the right-hand side of (50) evaluated at the kicking times τ_i . The action then becomes a sum of free-streaming Zel’dovich-type actions plus discrete gravitational contributions stemming from the kicking times. Between kicks one can use our MAK solution. At kicking times the velocity undergoes a discontinuous change that is related to the gravitational potential (and thus to the density) at those times. The densities at kicking times can be determined by an iterative procedure. The kicking strategy also allows one to perform redshift-space reconstruction by applying the redshift-space modified cost (Section 5) at the last kick.

7 COMPARISON WITH OTHER RECONSTRUCTION METHODS

Reconstruction started with Peebles’ (1989) work, in which he compared reconstructed and measured peculiar velocities for a small number of Local Group galaxies, situated within a few Mpc. The focus of reconstruction work has now moved to tackling the rapidly growing large 3D surveys (see, e.g., Frieman & Szalay 2000). It is not our intention here to review all the work on reconstruction;¹⁹ rather we shall discuss how some of the previously used methods can be reinterpreted in the light of the optimization approach to reconstruction. For convenience we shall divide methods into perturbative (Section 7.1), probabilistic (Section 7.2) and variational (Section 7.3). Methods such as POTENT (Dekel et al. 1990), the purpose of which is to obtain the full peculiar velocity field from its radial components using the (Eulerian) curl-free property, are not directly within our scope. Note that in its original Lagrangian form (Bertschinger & Dekel 1989; Dekel et al. 1990) POTENT was assuming a curl-free velocity in Lagrangian coordinates, an assumption closely related to the potential assumption made for MAK, as already pointed out in Section 3.1. Even closer is the relation between MAK and the PIZA method of Croft & Gaztañaga (1997), discussed in Section 7.3, which is also based on minimization of quadratic action.

7.1 Perturbative methods

Nusser & Dekel (1992) have proposed using the Zel’dovich approximation backwards in time to obtain the initial velocity fluctuations and thus (by slaving) the density fluctuations. Schematically, their procedure involves two steps: (i) obtaining the present potential velocity field and (ii) integrating the Zel’dovich–Bernoulli equation back in time. Using the equality (in our notation) of the velocity

and gravitational potentials, they point out that the velocity potential can be computed from the present density fluctuation field by solving the Poisson equation. This is a perturbative approximation to reconstruction in so far as it replaces the Monge–Ampère equation (19) by a linearized form. Indeed, when using the Zel’dovich approximation we have $q = x - \tau v = x + \tau \nabla_x \varphi_v(x)$. We know that $q = \nabla_x \Theta(x)$ with Θ satisfying the Monge–Ampère equation. The latter can thus be rewritten as

$$\det[\delta_{ij} + \tau \nabla_{x_i} \nabla_{x_j} \varphi_v(x)] = \rho(x), \quad (54)$$

where δ_{ij} denotes the identity matrix. If we now use the relation $\det(\delta_{ij} + \epsilon A_{ij}) = 1 + \epsilon \sum_i A_{ii} + O(\epsilon^2)$ and truncate the expansion at order ϵ , we obtain the Poisson equation

$$\tau \nabla_x^2 \varphi_v(x) = \rho(x) - 1 = \delta(x). \quad (55)$$

Of course, in one dimension no approximation is needed. From a physical point of view, equating the velocity and gravitational potentials at the present epoch amounts to using the Zel’dovich approximation in reverse and is actually inconsistent with the forward Zel’dovich approximation: the slaving which makes the two potentials equal initially does not hold in this approximation at later epochs. Replacing the Monge–Ampère equation by the Poisson equation is not consistent with a uniform initial distribution of matter and will in general lead to spurious multistreaming in the initial distribution. Of course, if the present-epoch velocity field happens to be known one can try applying the Zel’dovich approximation in reverse. Nusser and Dekel observe that calculating the inverse Lagrangian map by $q = x - \tau v$ does not work well (spurious multistreaming appears) and instead integrate back in time the Zel’dovich–Bernoulli equation²⁰

$$\partial_t \varphi_v = \frac{1}{2} (\nabla_x \varphi_v)^2, \quad (56)$$

which is obviously equivalent to the Burgers equation (13) with the viscosity $\nu = 0$. One way of performing this reverse integration, which guarantees the absence of multistreaming, is to use the Legendre transformation (18) to calculate $\Phi(q)$ from $\Theta(x) = |x|^2/2 - \tau \varphi_v(x)$ and then obtain the reconstructed initial velocity field as

$$v_{\text{in}}(q) = v_0(\nabla_q \Phi(q)). \quad (57)$$

This procedure can, however, lead to spurious shocks in the reconstructed initial conditions, due to inaccuracies in the present-epoch velocity data, unless the data are suitably smoothed. Finally, the improved reconstruction method of Gramann (1993) can be viewed as an approximation to the Monge–Ampère equation beyond the Poisson equation, which captures part of the non-linearity.

7.2 Probabilistic methods

Weinberg (1992) presents an original approach to reconstruction, which turns out to have hidden connections to optimal mass transportation. The key observations in his ‘Gaussianization’ technique are the following: (i) the initial density fluctuations are assumed to be Gaussian, (ii) the rank order of density values is hardly changed between initial and present states, (iii) the bulk displacement of large-scale features during dynamical evolution can be neglected.

²⁰ In the non-cosmological literature this equation is usually called Hamilton–Jacobi in the context of analytical mechanics (Landau & Lifshitz 1960) and Kardar–Parisi–Zhang (Kardar, Parisi & Zhang 1986) in condensed matter physics.

¹⁹ For a comparison of six different techniques, see Narayanan & Croft (1999).

Assumption (i) is part of the standard cosmological paradigm. Assumption (iii) can of course be tested in N -body simulations. As we have seen in Section 5, a displacement of $10 h^{-1}$ Mpc is typical and can indeed be considered small compared with the size of the simulation boxes ($64 h^{-1}$ Mpc in Weinberg's simulations and $200 h^{-1}$ Mpc in ours). Assumption (ii) means that the correspondence between initial and present values of the density ρ (or of the contrast $\delta = \rho - 1$) is monotonic. This map, which can be determined from the empirical present data, can then be applied to all the data to produce a reconstructed initial density field. Finally, by running an N -body simulation initialized on the reconstructed field one can test the validity of the procedure, which turns out to be quite good and can be improved further by hybrid methods (Narayanan & Weinberg 1998; Kolatt et al. 1996) combining Gaussianization with the perturbative approaches of Nusser & Dekel (1992) or Gramann (1993).

This technique is actually connected with mass transportation: starting with the work of Fréchet (1957a,b, see also Rachev 1984), probabilists have been asking the following question: given two random variables m_1 and m_2 with two laws, say probability density functions (PDFs) p_1 and p_2 , can one find a joint distribution of (m_1, m_2) with PDF $p_{12}(m_1, m_2)$ having the following properties: (i) p_1 and p_2 are the marginals, i.e. when p_{12} is integrated over m_2 (respectively, m_1) one recovers p_1 (respectively, p_2), (ii) the correlation $\langle m_1 m_2 \rangle$ is maximum? Since $\langle m_1^2 \rangle$ and $\langle m_2^2 \rangle$ are obviously prescribed by the constraint that we know p_1 and p_2 , maximizing the correlation is the same as minimizing the quadratic distance $\langle (m_1 - m_2)^2 \rangle$. This is precisely an instance of the mass transportation problem with quadratic cost, as we defined it in Section 3.3. As we know, the optimal solution is obtained by a map from the space of m_1 values to that of m_2 values which is the gradient of a convex function. If m_1 and m_2 are scalar variables, the map is just monotonic, as in the Gaussianization method (in the discrete setting this was already observed in Section 4.1). Hence Weinberg's method may be viewed as requiring maximum correlation (or minimum quadratic distance in the above sense) between initial and present distributions of density fluctuations.

In principle the Gaussianization method can be extended to multipoint distributions, leading to a difficult multidimensional mass transportation problem, which can be discretized into an assignment problem just as in Section 4.1. The contact of the maximum correlation assumption to the true dynamics is probably too flimsy to justify using such heavy machinery.

7.3 Variational methods

All variational approaches to reconstruction, starting with that of Peebles (1989), have common features: one uses a suitable Lagrangian and poses a two-point variational problem with boundary conditions prescribed at the present epoch by the observed density field, and at early times by requiring a quasi-uniform distribution of matter (more precisely, as we have seen in Section 2.1, by requiring that the solutions not be singular as $\tau \rightarrow 0$).

The path-interchange Zel'dovich approximation method of Croft & Gaztañaga (1997) and our MAK reconstruction techniques use a free-streaming Lagrangian in linear growth rate time. As we have seen in Section 3.1, this amounts to assuming adhesion dynamics. Once discretized for numerical purposes, the variational problem becomes an instance of the assignment problem. Croft & Gaztañaga (1997) have proposed a restricted procedure for solving it, which does not account for the Lagrangian potentiality and yields non-unique approximate solutions. As we have seen in Sections 4 and

5, the exact and unique solution can be found with reasonable CPU resources.

Turning now to the Peebles least-action method, let us first describe it schematically, using our notation. In its original formulation it is applied to a discrete set of galaxies (assumed of course to trace mass) in an Einstein–de Sitter universe. The action, in our notation, can be written as

$$I = \int_0^{\tau_0} d\tau \frac{3}{2\tau^{1/2}} \left(\sum_i \frac{m_i \tau^2}{3} \left| \frac{d\mathbf{x}_i}{d\tau} \right|^2 + \frac{3G}{2} \sum_{i \neq j} \frac{m_i m_j}{|\mathbf{x}_i - \mathbf{x}_j|} + \pi G \bar{\rho}_0 \sum_i m_i |\mathbf{x}_i|^2 \right), \quad (58)$$

where m_i is the mass and \mathbf{x}_i is the comoving coordinate of the i th galaxy (see also Nusser & Branchini 2000). This is supplemented by the boundary condition that the present positions of the galaxies are known and that the early-time velocities satisfy²¹

$$\tau^{3/2} \frac{d\mathbf{x}_i}{d\tau} \rightarrow 0 \quad \text{for } \tau \rightarrow 0. \quad (59)$$

This particle approach was extended by Gialvalisco et al. (1993) to a continuous distribution in Eulerian coordinates and leads then to the action analogous to (48), which we have used in Section 6. The procedure also involves a ‘Galerkin truncation’ of the particle trajectories to finite sums of trial functions of the form

$$x_i^\mu(\tau) = x_i^\mu(\tau_0) + \sum_{n=0}^{N-1} C_{i,n}^\mu f_n(\tau), \quad (60)$$

$$f_n(\tau) = \tau^n (\tau_0 - \tau), \quad n = 0, 1, \dots, N-1. \quad (61)$$

The reconstructed peculiar velocities for the Local Group were used by Peebles to calibrate the Hubble and density parameters, which turned out to differ from the previously assumed values. However, the peculiar velocity of one dwarf galaxy, N6822, failed to match the observed value (see Fig. 13). This led Peebles (1990) to partially relax the assumption of *minimum* action, also allowing for saddle points in the action. Somewhat better agreement with observations is then obtained, but at the expense of lack of uniqueness.

In the context of the present approach, various remarks can be made. The boundary condition (59) is trivially satisfied if the velocities $d\mathbf{x}/d\tau$ remain bounded. Actually, we have seen in Section 2.1 that, as a consequence of slaving, the velocity has a regular expansion in powers of τ , which implies its boundedness as $\tau \rightarrow 0$. The important point is that the function $f_n(\tau)$ appearing in (60) should be expandable in powers of τ , as is the case with the Ansatz (61).

In Section 6 we have established uniqueness of the reconstruction with a prescribed present density and under the assumption of absence of multistreaming (but we allow for mass concentrations). This restriction is meaningful only in the continuous case: in the discrete case, unless the particles are rather closely packed, the concept of multistreaming is not clear but there have been attempts to relate uniqueness to an absence of ‘orbit crossing’ (see, e.g., Gialvalisco et al. 1993; Whiting 2000). Of course, at the level of the underlying dark matter, multistreaming is certainly not ruled out at sufficiently small scales; at such scales unique reconstruction is not possible.

²¹ This condition, which is written $a^2 d\mathbf{x}_i/dt \rightarrow 0$ in Peebles' notation, ensures the vanishing of the corresponding boundary term after an integration by parts in the time variable.

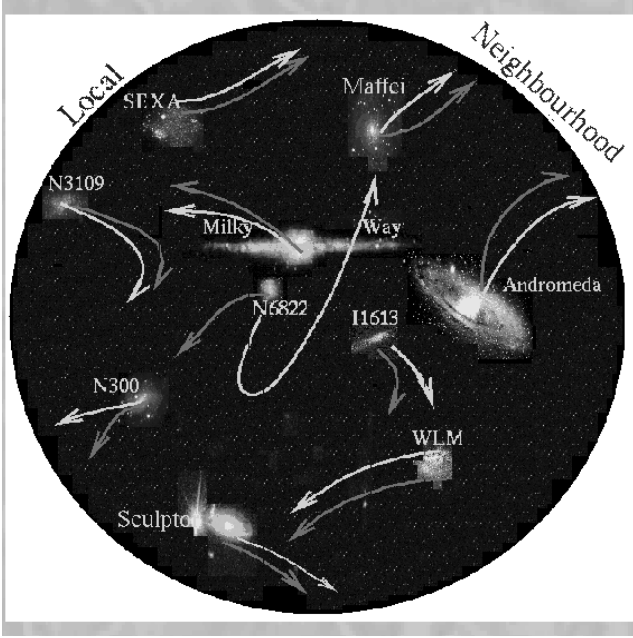


Figure 13. A schematic demonstration of Peebles’ reconstruction of the trajectories of the members of the local neighbourhood using a variational approach based on the minimization of Euler–Lagrange action. The arrows go back in time, starting from the present and pointing towards the initial positions of the sources. In most cases there is more than one allowed trajectory due to orbit crossing (closely related to the multistreaming of the underlying dark matter fluid). The darker (pink in the on-line version) orbits correspond to taking the minimum of the action, whereas the brighter (yellow) orbits were obtained by taking the saddle-point solution. Of particular interest is the orbit of N6822, which in the former solution is on its first approach towards us and in the second solution is in its passing orbit. A better agreement between the evaluated and observed velocities was shown to correspond to the saddle-point solution. This figure is available in colour in the on-line version of the journal on *Synergy*.

In the truly discrete case, e.g. when considering a dwarf galaxy, there is no reason to prefer the true minimum action solution over any other stationary action solution.

8 CONCLUSION

The main theoretical result of this paper is that reconstruction of the past dynamical history of the Universe, knowing only the present spatial distribution of mass, is a well-posed problem with a unique solution. More precisely, reconstruction is uniquely defined down to those scales, a few megaparsecs, where multistreaming becomes important. The presence of concentrated mass in the form of clusters, filaments, etc. is not an obstacle to a unique displacement reconstruction; the mass within each such structure originates from a collapsed region of known shape but with unknown initial density and velocity fluctuations inside. There are of course practical limitations to reconstruction stemming from the knowledge of the present mass distribution over only a limited patch of the Universe; these were discussed in Section 3.4.

In this paper we have also presented in detail and tested a reconstruction method called MAK which reduces reconstruction to an assignment problem with quadratic cost, for which effective algorithms are available. MAK, which is exact for dynamics governed by the adhesion model, works very well above $6 h^{-1}$ Mpc and can in principle be adapted to full Euler–Poisson reconstruction.

We note that a very common method for testing ideas concerning the early Universe is to take some model of early density fluctuations and then run N -body simulations with assumed cosmological parameters until the present epoch. Confrontation with the observed *statistical properties* of the present Universe then helps in selecting plausible models and in narrowing the choice of cosmological parameters. This *forward method* is conceptually very different from reconstruction; the latter not only works backward but, more importantly, it is a *deterministic* method that gives us a detailed map of the early Universe and how it relates to the present one. Reconstruction thus allows us to obtain the peculiar velocities of galaxies and is probably the only method that can hope to do this for a large number of galaxies. In those instances where we have partial information on peculiar velocities (from independent distance measurements), e.g. for the NearBy Galaxies (NBG) catalogue of Tully (1988), such information can be used to calibrate cosmological parameters or to provide additional constraints, which are in principle redundant but can improve the quality.

The detailed reconstruction of early density fluctuations, which will become possible using large 3D surveys such as 2dF and SDSS (see, e.g., Frieman & Szalay 2000), will allow us to test assumptions such as the Gaussianity of density fluctuations at decoupling. Note, however, that such a reconstruction gives us full access only to the complement of collapsed regions; any statistical information thus obtained will be biased, roughly by overemphasizing underdense regions.

Finally, we have no reason to hide the pleasure we experience in seeing this heavenly problem bring together and indeed depend crucially on so many different areas of mathematics and physics, from fluid dynamics to Monge–Ampère equations, mass transportation, convex geometry and combinatorial optimization. Probably this is the first time that the three-dimensional Monge–Ampère equation has been tackled numerically for practical purposes. As usual, we can expect that the techniques, here applied to cosmic reconstruction, will find many applications, for example to the optimal matching of two holographic or tomographic images or to the correction of images in multidimensional colour space.

ACKNOWLEDGMENTS

Special thanks are due to E. Branchini (observational and conceptual aspects) and to D. Bertsekas (algorithmic aspects). We also thank J. S. Bagla, J. Bec, E. Bertschinger, T. Buchert, A. Domínguez, H. Frisch, J. Gaite, C. l’Hostis, L. Moscardini, A. Noullez, M. Rees, V. Sahni, S. Shandarin, A. Shnirelman, E. Slezak, E. Spiegel, A. Starobinsky, P. Thomas and B. Villone for comments and useful information. This work was supported by the European Union under contract HPRN-CT-2000-00162, by the BQR programme of Observatoire de la Côte d’Azur, by the TMR programme of the European Union (UF), by MIUR (SM), by the French Ministry of Education, the McDonnell Foundation, and the Russian Foundation for Basic Research under grant RFBR 02-01-1062 (AS). RM was supported by a Marie Curie Fellowship HPMF-CT-2002-01532.

REFERENCES

- Ambrosio L., 2003, in Tatsien L., ed., Proc. ICM Vol. 3. World Scientific, Singapore, p. 131
- Ampère A.-M., 1820, J. Ecole R. Polytech., 11, 1
- Arnol’d V.I., Shandarin S.F., Zel’dovich Y.B., 1982, Geophys. Astrophys. Fluid Dynam., 20, 111
- Balinski M.L., 1986, Math. Prog., 34, 125
- Bec J., Frisch U., Khanin K., 2000, J. Fluid Mech., 416, 239

- Benamou J.-D., Brenier Y., 2000, *Numer. Math.*, 84, 375
- Bernardeau F., Colombi S., Gaztañaga E., Scoccimarro R., 2002, *Phys. Rep.*, 367, 1
- Bertschinger E., Dekel A., 1989, *ApJ*, 336, L5
- Bertsekas D.P., 1981, *Math. Prog.*, 21, 152
- Bertsekas D.P., 1992, *Comput. Optim. Appl.*, 1, 7
- Bertsekas D., 2001, *Encyclopedia of Optimization*, Vol. I. Kluwer, Dordrecht
- Bour É., 1862, *J. École Polytechn.*, 22, Cahier 39, 149
- Branchini E., Eldar A., Nusser A., 2002, *MNRAS*, 335, 53
- Brenier Y., 1987, *C.R. Acad. Sci. Paris I*, 305, 805
- Brenier Y., 1991, *Comm. Pure Appl. Math.*, 44, 375
- Bridle S.L., Lahav O., Ostriker J.P., Steinhardt P.J., 2003, *Sci*, 299, 1532
- Buchert T., 1992, *MNRAS*, 254, 729
- Buchert T., Domínguez A., 1998, *A&A*, 335, 395
- Buchert T., Ehlers J., 1993, *MNRAS*, 264, 375
- Burkard R.E., Derigs U., 1980, *Assignment and matching problems: solution methods with FORTRAN programs*, *Lecture Notes in Economics and Mathematical Systems*, Vol. 184. Springer-Verlag, Berlin
- Caffarelli L., Li Y.Y., 2003, *Comm. Pure Appl. Math.*, 56, 549
- Caffarelli L.A., 1999, in Christ M., Kenig C.E., Sadosky C., eds, *Chicago Lectures in Mathematics, Harmonic Analysis and Partial Differential Equations*. Univ. Chicago Press, Chicago, pp. 117–126
- Caffarelli L.A., Milman M., eds, 1999, *Monge Ampère Equation: Applications to Geometry and Optimization*. *Contemporary Mathematics*, Vol. 226. American Mathematical Society, Providence
- Catelan P., 1995, *MNRAS*, 276, 115
- Catelan P., Lucchin F., Matarrese S., Moscardini L., 1995, *MNRAS*, 276, 39
- Coles P., Lucchin F., 2002, *Cosmology: the Origin and Evolution of Cosmic Structure*. Wiley, Chichester
- Couchman H.M.P., Thomas P.A., Pearce F.R., 1995, *ApJ*, 452, 797
- Croft R.A.C., Gaztañaga E., 1997, *MNRAS*, 285, 793
- Dekel A., Bertschinger E., Faber S.M., 1990, *ApJ*, 364, 349
- Fanelli D., Aurell E., 2002, *A&A*, 395, 399
- Fenchel W., 1949, *Can. J. Math.*, 1, 73
- Fréchet M., 1957a, *C.R. Acad. Sci. Paris I*, 244, 689
- Fréchet M., 1957b, *Publ. Inst. Statist. Univ. Paris*, 6, 183
- Frieman J.A., Szalay A.S., 2000, *Phys. Rep.*, 333, 215
- Frisch U., Bec J., 2002, in Lesieur M., Yaglom A., David F., eds, *École de Physique des Houches, session LXXIV, New Trends in Turbulence*. EDP Sciences/Springer-Verlag, Berlin, pp. 341–384
- Frisch U., Matarrese S., Mohayaee R., Sobolevski A., 2002, *Nat*, 417, 260
- Gangbo W., McCann R.J., 1996, *Acta Math.*, 177, 113
- Giavalisco M., Mancinelli B., Mancinelli P.J., Yahil A., 1993, *ApJ*, 411, 9
- Goursat É., 1896, *Leçons sur l'Intégration des Équations aux Dérivées Partielles du Second Ordre*, Vol. I. Imprimerie Guthier-Villars, Paris
- Gramann M., 1993, *ApJ*, 405, 449
- Gurbatov S.N., Saichev A.I., 1984, *Izv. Vysš. Učebn. Zaved. Radiofiz.*, 27, 456
- Gurbatov S.N., Saichev A.I., Shandarin S.F., 1989, *MNRAS*, 236, 385
- Hénon M., 1995, *C.R. Acad. Sci. Paris I*, 321, 741
- Hénon M., 2002, *A Mechanical Model for the Transportation Problem*. Preprint (math.OC/0209047)
- Jeans J.H., 1919, *Problems of Cosmogony and Stellar Dynamics*. Cambridge Univ. Press, Cambridge
- Kantorovich L.V., 1942, *C.R. (Dokl.) Acad. Sci. USSR*, 321, 199
- Kardar M., Parisi G., Zhang Y., 1986, *Phys. Rev. Lett.*, 56, 889
- Kolatt T., Dekel A., Ganon G., Willick J.A., 1996, *ApJ*, 458, 419
- Kuhn H.W., 1955, *Naval Res. Logist. Q.*, 2, 83
- Landau L.D., Lifshitz E.M., 1960, *Mechanics. Course of Theoretical Physics*, Vol. 1. Pergamon, Oxford
- Loeper G., 2003, *The Inverse Problem for the Euler–Poisson System in Cosmology*. *Arch. Ration. Mech. Anal.*, submitted (math.AP/0306430)
- Mandelbrojt S., 1939, *C.R. Acad. Sci. Paris*, 209, 977
- Minkowski H., 1897, *Nachr. Ges. Wiss. Göttingen (Math. Phys. Klasse)*, pp. 198–219
- Mohayaee R., Frisch U., Matarrese S., Sobolevskii A., 2003, *A&A*, 406, 393
- Monaco P., Efstathiou G., 1999, *MNRAS*, 308, 763
- Monge G., 1781, *Histoire de l'Académie Royale des Sciences*, p. 666
- Monge G., 1784, *Histoire de l'Académie Royale des Sciences*, p. 118
- Moutarde F., Alimi J.-M., Bouchet F.R., Pellat R., Ramani A., 1991, *ApJ*, 382, 377
- Munshi D., Sahni V., Starobinsky A.A., 1994, *ApJ*, 436, 517
- Narayanan V.K., Croft R.A.C., 1999, *ApJ*, 515, 471
- Narayanan V.K., Weinberg D.H., 1998, *ApJ*, 508, 440
- Nusser A., Branchini E., 2000, *MNRAS*, 313, 587
- Nusser A., Dekel A., 1992, *ApJ*, 391, 443
- Papadimitriou C.H., Steiglitz K., 1982, *Combinatorial Optimization: Algorithms and Complexity*. Prentice-Hall, Englewood Cliffs
- Peebles P.J.E., 1980, *The Large-scale Structure of the Universe*. Princeton Univ. Press, Princeton
- Peebles P.J.E., 1989, *ApJ*, 344, L53
- Peebles P.J.E., 1990, *ApJ*, 362, 1
- Pogorelov A.V., 1978, *The Minkowski Multidimensional Problem*. Winston, Washington
- Rachev S.T., 1984, *Teor. Veroyatnost. Primenen.*, 29, 625
- Rockafellar R.T., 1970, *Convex Analysis*. Princeton Univ. Press, Princeton
- Sathyaprakash B.S., Sahni V., Shandarin S., Fisher K.B., 1998, *ApJ*, 507, L109
- Shandarin S.F., Sathyaprakash B.S., 1996, *ApJ*, 467, L25
- Shandarin S.F., Zel'dovich Y.B., 1989, *Rev. Mod. Phys.*, 61, 185
- Spergel D.N. et al., 2003, *ApJS*, 148, 175
- Susperregi M., Binney J., 1994, *MNRAS*, 271, 719
- Taylor A., Valentine H., 1999, *MNRAS*, 306, 491
- Tomizawa N., 1971, *Networks*, 1, 173
- Tully R.B., 1988, *Nearby Galaxies Catalog*. Cambridge Univ. Press, Cambridge
- Valentine H., Saunders W., Taylor A., 2000, *MNRAS*, 319, L13
- Vergassola M., Dubrulle B., Frisch U., Noullez A., 1994, *A&A*, 289, 325
- Weber E. v., 1900, in Burkhardt H., ed., *Encyklopädie der Mathematischen Wissenschaften*, Vol. 2, Analysis. Teubner, Leipzig, p. 294
- Weinberg D.H., 1992, *MNRAS*, 254, 315
- Weinberg D.H., Gunn J.E., 1990, *MNRAS*, 247, 260
- Whiting A.B., 2000, *ApJ*, 533, 50
- Zel'dovich Y.B., 1970, *A&A*, 5, 84

APPENDIX A: EQUATIONS OF MOTION IN AN EXPANDING UNIVERSE

On distances covered by present and forthcoming redshift galaxy catalogues, the Newtonian description constitutes a realistic approximation to the dynamics of self-gravitating cold dark matter filling the Universe (Peebles 1980; Coles & Lucchin 2002). This description gives, in proper space coordinates denoted here by \mathbf{r} and cosmic time t , the familiar Euler–Poisson system for the density $\varrho(\mathbf{r}, t)$, velocity $\mathbf{U}(\mathbf{r}, t)$ and the gravitational potential $\phi(\mathbf{r}, t)$:

$$\partial_t \mathbf{U} + (\mathbf{U} \cdot \nabla_r) \mathbf{U} = -\nabla_r \phi_g, \quad (\text{A1})$$

$$\partial_t \varrho + \nabla_r \cdot (\varrho \mathbf{U}) = 0, \quad (\text{A2})$$

$$\nabla_r^2 \phi_g = 4\pi G \varrho, \quad (\text{A3})$$

where G is the gravitational constant.

In a homogeneous isotropic universe, the density and velocity fields take the form

$$\varrho(\mathbf{r}, t) = \bar{\varrho}(t), \quad \mathbf{U}(\mathbf{r}, t) = H(t) \mathbf{r} = \frac{\dot{a}(t)}{a(t)} \mathbf{r}, \quad (\text{A4})$$

where the coefficient $H(t)$ is the *Hubble parameter* and $a(t)$ is the *expansion scalefactor* defined so that integration of the velocity field $\dot{\mathbf{r}} = \mathbf{U}(\mathbf{r}, t) = H(t) \mathbf{r}$ yields $\mathbf{r} = a(t) \mathbf{x}$, where \mathbf{x} is called the *comoving coordinate*.

The *background density* $\bar{\varrho}(t)$ gives rise to the background gravitational potential $\bar{\phi}_g$, which by (A1) and (A4) satisfies

$$-\nabla_r \bar{\phi}_g = \frac{\ddot{a}}{a} \mathbf{r}. \quad (\text{A5})$$

For the background density, mass conservation (A2) then gives

$$\bar{\varrho} a^3 = \bar{\varrho}_0, \quad (\text{A6})$$

where $\bar{\varrho}_0 = \bar{\varrho}(t_0)$ with t_0 the present epoch and $a(t_0)$ is normalized to unity. Equations (A5), (A6) and (A3) imply the *Friedmann equation* for $a(t)$:

$$\ddot{a} = -\frac{4}{3}\pi G \bar{\varrho}_0 \frac{1}{a^2} \quad (\text{A7})$$

with conditions posed at $t = t_0$:

$$a(t_0) = 1, \quad \dot{a}(t_0) = H_0 > 0, \quad (\text{A8})$$

where H_0 is the present value of the Hubble parameter, which is positive for an expanding universe.

For simplicity we restrict ourselves to the case of the *critical density*, corresponding to the flat, matter-dominated Einstein–de Sitter universe (without a cosmological constant):

$$\bar{\varrho}_0 = \frac{3H_0^2}{8\pi G} \quad (\text{A9})$$

and adjust the origin of the time axis such that the solution takes the form of a power law

$$a(t) = \left(\frac{t}{t_0}\right)^{2/3} \quad (\text{A10})$$

with $H_0 = 2/(3t_0)$ and $\bar{\varrho}_0 = 1/(6\pi G t_0^2)$.

The observed Hubble expansion of the Universe suggests that the density, velocity and gravitational fields may be decomposed into a sum of terms describing the uniform expansion and fluctuations against the background:

$$\varrho = \bar{\varrho}(t)\rho, \quad \mathbf{U} = \frac{\dot{a}(t)}{a(t)}\mathbf{r} + a(t)\mathbf{u}, \quad \phi_g = \bar{\phi}_g + \tilde{\phi}_g. \quad (\text{A11})$$

The term $a(t)\mathbf{u}$ is called the *peculiar velocity*. In cosmology, one also often employs the *density contrast* defined as $\delta = \rho - 1$, which gives the fluctuation against the normalized background density. Taking ρ , \mathbf{u} and $\tilde{\phi}_g$ as functions of the comoving coordinate $\mathbf{x} = \mathbf{r}/a(t)$ and using (A5)–(A7), we rewrite the Euler–Poisson system in the form

$$\partial_t \mathbf{u} + (\mathbf{u} \cdot \nabla_x) \mathbf{u} = -2\frac{\dot{a}}{a} \mathbf{u} - \frac{1}{a} \nabla_x \tilde{\phi}_g, \quad (\text{A12})$$

$$\partial_t \rho + \nabla_x \cdot (\rho \mathbf{u}) = 0, \quad (\text{A13})$$

$$\nabla_x^2 \tilde{\phi}_g = \frac{4\pi G \bar{\varrho}_0}{a} (\rho - 1). \quad (\text{A14})$$

Note the *Hubble drag* term $-2(\dot{a}/a)\mathbf{u}$ in the right-hand side of (A12) representing the relative slowdown of peculiar velocities due to the uniform expansion.

Formally linearizing (A12)–(A14) around the trivial zero solution, one obtains the following ordinary differential equation for the *linear growth factor* $\tau(t)$ of density fluctuations:

$$\frac{d}{dt}(a^2 \dot{\tau}) = 4\pi G \bar{\varrho}_0 \frac{\tau}{a}. \quad (\text{A15})$$

The only solution of this equation that stays bounded (indeed, vanishes) at small times is usually referred to as the *growing mode*. As we shall shortly see, it is convenient to choose the amplitude factor τ of the growing mode to be a new ‘time variable’, which in an

Einstein–de Sitter universe is proportional to $t^{2/3}$. It is normalized such that $\tau_0 = \tau(t_0) = 1$. Rescaling the peculiar velocity and the gravitational potential according to

$$\mathbf{u} = \dot{\tau} \mathbf{v}, \quad \tilde{\phi}_g = \frac{4\pi G \bar{\varrho}_0 \tau}{a} \varphi_g \quad (\text{A16})$$

and using the fact that in an Einstein–de Sitter universe $d \ln(a^2 \dot{\tau})/d\tau = 3/(2\tau)$, we arrive at the following form of the *Euler–Poisson system*, which we use throughout this paper:

$$\partial_\tau \mathbf{v} + (\mathbf{v} \cdot \nabla_x) \mathbf{v} = -\frac{3}{2\tau} (\mathbf{v} + \nabla_x \varphi_g), \quad (\text{A17})$$

$$\partial_\tau \rho + \nabla_x \cdot (\rho \mathbf{v}) = 0, \quad (\text{A18})$$

$$\nabla_x^2 \varphi_g = \frac{\rho - 1}{\tau}. \quad (\text{A19})$$

Suppose initially, i.e. at $\tau = 0$, a mass element is located at a point with the comoving coordinate \mathbf{q} . Transported by the peculiar velocity field in the comoving coordinates, this element describes a trajectory $\mathbf{x}(\mathbf{q}, \tau)$. Using the *Lagrangian coordinate* \mathbf{q} to parametrize the whole continuum of mass elements, we recast (A17) and (A19) in the form

$$D_\tau^2 \mathbf{x} = -\frac{3}{2\tau} (D_\tau \mathbf{x} + \nabla_x \varphi_g), \quad (\text{A20})$$

$$\nabla_x^2 \varphi_g = \frac{1}{\tau} [(\det \nabla_q \mathbf{x})^{-1} - 1]. \quad (\text{A21})$$

The density and peculiar velocity in Lagrangian variables are given by

$$\rho(\mathbf{x}(\mathbf{q}, \tau), \tau) = (\det \nabla_q \mathbf{x})^{-1}, \quad (\text{A22})$$

$$\mathbf{v}(\mathbf{x}(\mathbf{q}, \tau), \tau) = D_\tau \mathbf{x}(\mathbf{q}, \tau),$$

which automatically satisfy the mass conservation law (A18). Here D_τ is the operator of Lagrangian time derivative, which in Lagrangian variables is the usual partial time derivative at constant \mathbf{q} and in Eulerian variables coincides with the material derivative $\partial_\tau + \mathbf{v} \cdot \nabla_x$. The notation ∇_x in Lagrangian variables stands for the $\mathbf{x}(\mathbf{q}, \tau)$ -dependent differential operator with components $\nabla_{x_i} \equiv (\partial q_j / \partial x_i) \nabla_{q_j}$, which expresses the Eulerian gradient rewritten in Lagrangian coordinates, using the inverse Jacobian matrix. Note that ∇_x and D_τ do not commute and that terms with ∇_x in the Lagrangian equations are implicitly non-linear.

In one dimension, equation (A21) has an interesting consequence:

$$\nabla_x \varphi_g = -\frac{x - q}{\tau}. \quad (\text{A23})$$

Indeed, in one dimension (A21) takes the form

$$\nabla_x^2 \varphi_g = \frac{1}{\tau} [(\nabla_q x)^{-1} - 1]. \quad (\text{A24})$$

Multiplying this equation by $\nabla_q x$ and expressing the first of the two x -derivatives acting on φ_g as a q -derivative, we obtain

$$\nabla_q (\nabla_x \varphi_g) = \nabla_q \frac{q - x}{\tau}. \quad (\text{A25})$$

Equation (A23) is obtained from (A25) by integrating in q . The absence of an arbitrary τ -dependent constant is established either by assuming vanishing at large distances of both φ_g and of the displacement $x - q$ or, in the space-periodic case, by assuming the vanishing of period averages.

Using (A23) to eliminate the φ_g term in (A20) and introducing the notation ξ for the displacement $x - q$, we obtain

$$D_\tau^2 \xi = -\frac{3}{2\tau} \left(D_\tau \xi - \frac{\xi}{\tau} \right). \quad (\text{A26})$$

The only solution to this equation that remains well behaved for $\tau \rightarrow 0$ is the linear one $\xi \propto \tau$. This solution has the two terms on the right-hand side of the one-dimensional version of (A20) cancelling each other and hence gives a vanishing ‘acceleration’ $D_\tau^2 x$.

An approximate vanishing of acceleration takes place in higher dimensions as well. For early times, the *Lagrangian map* $\mathbf{x}(\mathbf{q}, \tau)$ stays close to the identity, with displacements $\xi(\mathbf{q}, \tau) = \mathbf{x}(\mathbf{q}, \tau) - \mathbf{q}$ small. Linearizing (A20) and (A21) around zero displacement, we obtain the system

$$D_\tau^2 \xi = -\frac{3}{2\tau}(D_\tau \xi + \nabla_q \varphi_g), \quad (\text{A27})$$

$$\nabla_q^2 \varphi_g = -\frac{1}{\tau} \nabla_q \cdot \xi, \quad (\text{A28})$$

where we use the fact that $\nabla_x \simeq \nabla_q$ and $\det \nabla_q \mathbf{x} \simeq 1 + \nabla_q \cdot \xi$. Using (A28) to eliminate φ_g in (A27), we obtain for $\theta \equiv \nabla_q \cdot \xi$ an equation that coincides with (A26) up to the change of variable $\xi \mapsto \theta$. Choosing the well-behaved linear solution for θ , solving for ξ and using the above argument to eliminate a τ -dependent constant, we see that, in the linearized equations, terms in the right-hand side of (A27) cancel each other and the acceleration vanishes. This simplification justifies using the linear growth factor τ as a time variable.

APPENDIX B: HISTORY OF MASS TRANSPORTATION

The subject of mass transportation was started by Gaspard Monge (1781) in a paper²² entitled *Théorie des déblais et des remblais* (Theory of cuts and fills) the preamble of which is worth quoting entirely (our translation):

‘When earth is to be moved from one place to another, the usage is to call *cuts* the volumes of earth to be transported and *fills* the space to be occupied after transportation.

The cost of transporting one molecule being, all things otherwise equal, proportional to its weight and to the distance [*espace*] travelled and consequently the total cost being proportional to the sum of products of molecules each multiplied by the distance travelled, it follows that for given shapes and positions of the cuts and fills, it is not indifferent that any given molecule of the cuts be transported to this or that place in the fills, but there ought to be a certain distribution of molecules of the former into the latter, according to which the sum of these products will be the least possible, and the cost of transportation will be a *minimum*.’

Although clearly posed, the ‘mass transportation problem’ was not solved, in more than one dimension, until Leonid Kantorovich (1942) formulated a ‘relaxed’ version, now called the Monge–Kantorovich problem: instead of a ‘distribution of molecules of the former *into* the latter’, he allowed a distribution in the product space where more than one position in the fills could be associated with a position in the cuts and where the initial and final distributions are prescribed marginals (see Section 3.3). In *cosmospeak*, he allowed multistreaming with given initial and final mass distributions. Using the techniques of duality and of linear programming that he had invented (see Appendix C2), Kantorovich was then able to solve the mass transportation problem in this relaxed formulation. The techniques developed by Kantorovich found many applications, notably in economics, which in fact was his original motivation (he

was awarded, together with T.C. Koopmans, the 1975 Nobel prize in this field).

Before turning to more recent developments we must say a few words concerning the history of the Monge–Ampère equation. It was considered for the first time by Ampère (1820) for an unknown function $z(x, y)$ of two scalar variables. The equation is to be found on p. 65 of Ampère’s huge (188 pages) mathematical memoir in the form

$$Hr + 2Ks + Lt + M + N(rt - s^2) = 0, \quad (\text{B1})$$

where in modern notation $r = \partial^2 z / \partial x^2$, $s = \partial^2 z / (\partial x \partial y)$, $t = \partial^2 z / \partial y^2$, and H, K, L, M, N are functions of x, y, z and the two first-order derivatives $p = \partial z / \partial x$ and $q = \partial z / \partial y$. This extends the earlier work by Monge (1784, p. 126) concerning the equation without the Hessian term ($N = 0$). Both Ampère and Monge were interested in methods of explicit integration of these equations. Ampère also pointed out the way the equation changes under Legendre transformations but there is no physical interpretation in terms of Lagrangian coordinates.²³ There is evidence that until the beginning of the 20th century the scientific community attributed the equation with the Hessian solely to Ampère (see e.g. Bour 1862, p. 186 and Weber 1900, p. 367). However, the joint attribution of (B1) to ‘Monge and Ampère’ is already found in Goursat (1896).

The subjects of mass transportation and of the Monge–Ampère equation came together when one of us (YB) showed the equivalence of the elliptic Monge–Ampère equation and of the mass transportation problem with quadratic cost: when initial and final distributions are non-singular, the optimal solution is actually one-to-one, so that nothing is lost by the Kantorovich relaxation trick (Brenier 1987, 1991). For an extension of this result to general costs see Gangbo & McCann (1996); a review of the many recent papers on the subject is given by Ambrosio (2003).

APPENDIX C: BASICS OF CONVEXITY AND DUALITY

C1 Convexity and the Legendre transformation

A convex body may be defined by the condition that it coincides with the intersection of all half-spaces containing it. Obviously, it is sufficient to take only those half-spaces limited by planes that touch the body; such planes are called *supporting*.

Now take a convex function $f(\mathbf{q})$, so that the set of points in the $(3 + 1)$ -dimensional (\mathbf{q}, f) space lying above its graph is convex. It follows that we can write

$$f(\mathbf{q}) = \max_x \mathbf{x} \cdot \mathbf{q} - f^*(\mathbf{x}), \quad (\text{C1})$$

where the expression $\mathbf{x} \cdot \mathbf{q} - f^*(\mathbf{x})$ specifies a supporting plane with the slope \mathbf{x} for the set of points lying above the graph of f (see Fig. C1 for the one-dimensional case).

The function $f^*(\mathbf{x})$, which specifies how high one should place a supporting plane to touch the graph, is called the *Legendre transform* of $f(\mathbf{q})$.²⁴

²³ According to the biography of Ampère by L. Pearce Williams in the *Dictionary of Scientific Biography*, Ampère’s paper was written – after he had switched from mathematics to chemistry and physics – with the purpose of facilitating his election to the Paris Academy of Science; one can then speculate that his mention of the Legendre transformation was influenced by Legendre’s presence in this academy.

²⁴ It was introduced in the one-dimensional case by Mandelbrojt (1939) and then generalized by Fenchel (1949).

²² The author’s name appears in this paper as ‘M. Monge’, where the ‘M.’ stands for ‘Monsieur.’

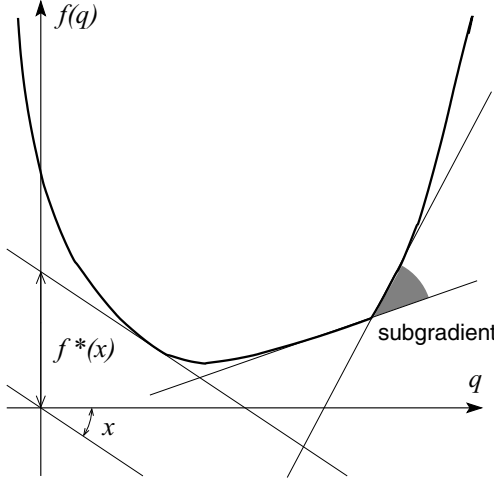


Figure C1. A convex function $f(q)$ and the geometrical construction of its Legendre transform $f^*(x)$. Also illustrated is the subgradient of $f(q)$ at a non-smooth point.

From equation (C1) follows the inequality (known as the *Young inequality*)

$$f(q) + f^*(x) \geq x \cdot q \quad \text{for all } x, q, \quad (\text{C2})$$

where both sides coincide if and only if the supporting plane with the slope x touches the graph of f at q . This fact, together with the obvious symmetry of this inequality, implies that

$$f^*(x) = \max_q x \cdot q - f(q). \quad (\text{C3})$$

Thus, the Legendre transform of a convex function is itself convex and the Legendre transform of the Legendre transform recovers the initial convex function.

If, however, we apply (C1) to a *non-convex* function f , we obtain a convex function f^* , where the Legendre transform will give the *convex hull* of f , the largest convex function for which the graph lies below that of f .

When f is both convex and differentiable, (C2) becomes an equality for $x = \nabla_q f(q)$. If f^* is also differentiable, then one has also $q = \nabla_x f^*(x)$. This is actually Legendre's original definition of the transformation, which is thus limited to smooth functions. Furthermore, if the original function is not convex and thus has the same gradient at separate locations, Legendre's purely local definition will give a multivalued Legendre transform. (In the context of the present paper this corresponds to multistreaming.)

Not all convex functions are differentiable [e.g. $f(q) = |q|$]. However, the Young inequality can be employed to define a useful generalization of the gradient: the *subgradient* of f at q is the set of all x for which the equality in (C2) holds (see Fig. C1). If f is smooth at q , then $\nabla_q f(q)$ will be the only such point; otherwise, there will be a (convex) set of them.

If a convex function has the same subgradient at more than one point, the function is said to lack *strict convexity*. In fact, strict convexity and smoothness are complementary: lack of one in a convex function implies lack of the other in the Legendre transform.

For further background on convex analysis and geometry, see Rockafellar (1970).

C2 Duality in optimization

Suppose we want to minimize a convex function $\Phi(q)$ subject to a set of linear constraints that may be written in matrix notation as

$Aq = b$ (vectors q satisfying this constraint are called *admissible* in optimization parlance). We now observe that

$$\inf_{Aq=b} \Phi(q) = \inf_q \sup_x \Phi(q) - x \cdot (Aq - b). \quad (\text{C4})$$

Indeed, should Aq not equal b , the sup operation in x will give infinity, so such q will not contribute to minimization. Here we use the inf/sup notation instead of min/max because the extremal values may not be reached, e.g. when they are infinite.

Using (C1), we rewrite this in the form

$$\begin{aligned} \inf_q \sup_{x,y} y \cdot q - \Phi^*(y) - x \cdot (Aq - b) \\ = \inf_q \sup_{x,y} (y - A^T x) \cdot q - \Phi^*(y) + x \cdot b, \end{aligned} \quad (\text{C5})$$

where $\Phi^*(y)$ is the Legendre transform of $\Phi(q)$ and A^T is the transpose of A . Taking inf in q first, we see that the expression in the right-hand side will be infinite unless $y = A^T x$. We then obtain the optimization problem of finding

$$\sup_x x \cdot b - \Phi^*(A^T x), \quad (\text{C6})$$

which is called *dual* to the original one. Note that there are no constraints on the dual variable x : any value is admissible.

Denoting solutions of problems (C4) and (C6) by q^* and x^* , we see that

$$\Phi(q^*) + \Phi^*(A^T x^*) - x^* \cdot b = 0, \quad (\text{C7})$$

because the optimal values of both problems are given by (C5) and thus coincide. Furthermore, for any admissible q and x

$$\Phi(q) + \Phi^*(A^T x) - x \cdot b \geq 0, \quad (\text{C8})$$

because the right-hand sides of (C4) and (C6) cannot pass beyond their optimal values.

Moreover, let equality (C7) be satisfied for some admissible q^* and x^* ; then such q^* and x^* must solve the problems (C4) and (C6). Indeed, taking e.g. x^* for x in (C8) and using (C7), we see that for any other admissible q

$$\Phi(q^*) \leq \Phi(q), \quad (\text{C9})$$

i.e. that q^* solves the original optimization problem (C4).

Convex optimization problems with linear constraints considered in this section are called *convex programs*. Their close relatives are *linear programs*, namely optimization problems of the form

$$\inf_{Aq=b, q \geq 0} c \cdot q = \inf_{q \geq 0} \sup_x c \cdot q - x \cdot (Aq - b), \quad (\text{C10})$$

where the notation $q \geq 0$ means that all components of the vector q are non-negative. Proceeding essentially as above with $c \cdot q$ instead of $\Phi(q)$, we observe that in order not to obtain infinity when minimizing in q in (C5), we now have to require that $A^T x \leq c$ (i.e. $c - A^T x \geq 0$). The dual problem thus takes the form

$$\sup_{A^T x \leq c} x \cdot b \quad (\text{C11})$$

with an admissibility constraint on x . Instead of (C7) and (C8) we obtain

$$x^* \cdot b = c \cdot q^* \quad \text{or} \quad (A^T x^* - c) \cdot q^* = 0 \quad (\text{C12})$$

and

$$x \cdot b \leq c \cdot q \quad \text{or} \quad (A^T x - c) \cdot q \leq 0, \quad (\text{C13})$$

the latter inequality being automatically satisfied for any admissible x, q . Note that for linear programs, the fact that (C12) holds for some admissible q^*, x^* also implies that q^* and x^* solve their respective optimization problems.

For further background on optimization and duality, see, for example, Papadimitriou & Steiglitz (1982).

C3 Why the analogue computer of Section 4.2 solves the assignment problem

We suppose that the analogue computer described in Section 4.2 has settled into equilibrium, which minimizes its potential energy

$$U = \sum_{i=1}^N \alpha_i - \sum_{j=1}^N \beta_j \quad (\text{C14})$$

under the set of constraints

$$\alpha_i - \beta_j \geq C - c_{ij} \quad (\text{C15})$$

for all i, j . Our goal is here is to show that the set of equilibrium forces f_{ij} , acting on studs between row and column rods, solves the original linear programming problem of minimizing

$$\bar{I} = \sum_{i,j=1}^N c_{ij} f_{ij} \quad (\text{C16})$$

under constraints

$$f_{ij} \geq 0, \quad \sum_{k=1}^N f_{kj} = \sum_{k=1}^N f_{ik} = 1, \quad (\text{C17})$$

for all i, j and that in fact forces f_{ij} take only zero and unit values, thus providing the solution to the assignment problem.

Note first that if a row rod A_i and a column rod B_j are not in contact at equilibrium, then the corresponding force vanishes ($f_{ij} = 0$); if they are, then $f_{ij} \geq 0$. Now take a particular pair of rods A_i and B_j that are in contact. At equilibrium, the force f_{ij} must equal forces exerted on the corresponding stud by A_i and B_j . We claim that both of these forces must be integer. To see this, let us compute the force exerted by A_i . This rod contributes its weight, $+1$, possibly decreased by the force that it feels from other column rods that are in contact with A_i . Each of these takes -1 (its ‘buoyancy’) out of the total force, but we may have to add the force it feels in turn from other row rods with which it might be in contact. Proceeding in this way from one rod to another, we see that all contributions, whether positive or negative, are unity, so their sum f_{ij} must be integer. The same argument applies to rod B_j .

Does this process indeed finish or, at some stage, do we come back at an already visited stud and thus end up in an infinite cycle? In fact, for a general set of stud lengths $C - c_{ij}$, the latter cannot happen, because otherwise an alternating sum of some subset of stud lengths would give exactly zero – a zero probability event for a set of arbitrary real numbers.

Consider now a row rod A_i . It is in contact with one or more column rods, for which the combined upward push must equilibrate the unit weight of A_i . Since any of the latter rods exerts a non-negative integer force, it follows that exactly one of these forces is unity, and all the other ones are zero. A similar argument holds for any column rod B_j .

We have thus shown that all f_{ij} in the equilibrium equal 1 or 0. One can of course ignore the vanishing forces. Then each row rod A_i is supported by exactly one column rod B_j , and each B_j supports exactly one A_i . This defines a one-to-one pairing, and we are only left with a check that this pairing minimizes (C16).

Observe that pushing a column rod down by some distance Δ and simultaneously increasing by Δ the length of all studs attached to this rod will have no effect on positions and constraints of all other rods, hence on the equilibrium network of contacts. Moreover, due to constraints (C17), the corresponding change in coefficients c_{ij} will not change the cost function (C16) in any essential way, except for just subtracting Δ .

We can use this observation to put all column rods at the same level, say at $z = 0$, adjusting c_{ij} to some new values c'_{ij} . Thus, for every i , the row rod A_i rests on the stud with the largest height $C - c'_{ij}$, so the equilibrium pairing maximizes the sum

$$\sum_{i,j=1}^N (C - c'_{ij}) f_{ij} \quad (\text{C18})$$

and thus minimizes (C16).²⁵

APPENDIX D: DETAILS OF THE VARIATIONAL TECHNIQUE FOR THE EULER–POISSON SYSTEM

In this appendix, we explain details of the variational procedure outlined in Section 6, which proves that prescription of the density fields at terminal epochs $\tau = 0$ and $\tau = \tau_0$ uniquely determines a regular and thus curl-free solution to the Euler–Poisson system (A17)–(A19).

The variational problem is posed for the functional

$$I = \frac{1}{2} \int_0^{\tau_0} dt \int d^3\mathbf{x} \tau^{3/2} \left(\rho |v|^2 + \frac{3}{2} |\nabla_x \phi_g|^2 \right) \quad (\text{D1})$$

with four constraints: the Poisson equation (A19), which we repeat here for convenience,

$$\nabla_x^2 \phi_g = \frac{\rho - 1}{\tau}, \quad (\text{D2})$$

the mass conservation (A18), also repeated here,

$$\partial_\tau \rho + \nabla_x \cdot (\rho v) = 0, \quad (\text{D3})$$

and the two boundary conditions

$$\rho(\mathbf{x}, 0) = 1 \quad \text{and} \quad \rho(\mathbf{x}, \tau_0) = \rho_0(\mathbf{x}). \quad (\text{D4})$$

In the following, we shall always denote by \iint the double integration over $0 \leq \tau \leq \tau_0$ and over the whole space domain in \mathbf{x} provided that the integrand vanishes at infinity sufficiently fast, or over the periodicity box in the case of periodic boundary conditions. A single integral sign \int will always denote the integration over the relevant space domain in \mathbf{x} .

First, we make this problem convex by rewriting the functional and constraints in a new set of variables with the mass flux $\mathbf{J}(\mathbf{x}, t) = \rho(\mathbf{x}, t) \mathbf{v}(\mathbf{x}, t)$ instead of the velocity \mathbf{v} . The mass conservation constraint, which was the only non-linear one in the old variables, now becomes linear:

$$\partial_\tau \rho + \nabla_x \cdot \mathbf{J} = 0, \quad (\text{D5})$$

and one can check that the density of kinetic energy takes the form

$$\frac{1}{2} \rho |v|^2 = \frac{1}{2\rho} |\mathbf{J}|^2 = \max_{c, \mathbf{m}: c + |\mathbf{m}|^2/2 \leq 0} (\rho c + \mathbf{J} \cdot \mathbf{m})$$

or

$$\frac{|\mathbf{J}|^2}{2\rho} = \max_{c, \mathbf{m}} [\rho c + \mathbf{J} \cdot \mathbf{m} - F(c, \mathbf{m})], \quad (\text{D6})$$

²⁵ Those readers familiar with linear programming will recognize that the proof just presented is based on two ideas: (i) the total unimodularity of the matrix of constraints in terms of which the equalities in (C17) can be written and (ii) the complementary slackness (see, e.g., Papadimitriou & Steiglitz 1982, sections 3.2 and 13.2).

where

$$F(c, \mathbf{m}) = \begin{cases} 0 & \text{if } c + |\mathbf{m}|^2/2 \leq 0 \\ +\infty & \text{otherwise.} \end{cases} \quad (\text{D7})$$

Note that in (D6) the variables c, \mathbf{m} , as well as ρ, \mathbf{J} , are functions of (\mathbf{x}, τ) . The action functional may now be written as

$$I = \frac{1}{2} \iint \left(\frac{1}{\rho} |\mathbf{J}|^2 + \frac{3}{2} |\nabla_{\mathbf{x}} \phi|^2 \right) \tau^{3/2} d^3 \mathbf{x} d\tau, \quad (\text{D8})$$

and turns out to be convex.

To see this, first note that the operation of integration is linear and thus preserves convexity of the integrand. The integrand is a positive quadratic function of $\nabla_{\mathbf{x}} \phi$ and therefore is convex in ϕ ; furthermore, (D6) implies that it is also convex in (ρ, \mathbf{J}) , since the kinetic energy density $|\mathbf{J}|^2/2\rho$ is the Legendre transform of the function $F(c, \mathbf{m})$, which itself is convex.

Note also that by representing the kinetic energy density in the form (D6), we may safely allow ρ to take negative values: the right-hand side being in that case $+\infty$, it will not contribute to minimizing (D1).

We now derive the dual optimization problem. We introduce the scalar Lagrange multipliers $\psi(\mathbf{x}, t)$, $\vartheta_{\text{in}}(\mathbf{x})$, $\vartheta_0(\mathbf{x})$ and $\theta(\mathbf{x}, t)$ for the Poisson equation (D2), the boundary conditions (D4) and the constraints of mass conservation (D5), respectively, and observe that the variational problem may now be written in the form

$$\begin{aligned} \inf_{\rho, \mathbf{J}, \phi} \sup_{\substack{c, \mathbf{m}, \theta, \psi, \vartheta_0, \vartheta_{\text{in}}, \\ c + |\mathbf{m}|^2/2 \leq 0}} \iint d^3 \mathbf{x} d\tau \left[\frac{3}{2} \psi \left(\nabla_{\mathbf{x}}^2 \phi - \frac{\rho - 1}{\tau} \right) \right. \\ \left. + \theta (\partial_{\tau} \rho + \nabla_{\mathbf{x}} \cdot \mathbf{J}) + \tau^{3/2} \left(\rho c + \mathbf{J} \cdot \mathbf{m} + \frac{3}{4} |\nabla_{\mathbf{x}} \phi|^2 \right) \right] \\ + \int \vartheta_{\text{in}}(\mathbf{x}) (\rho(\mathbf{x}, 0) - 1) d^3 \mathbf{x} \\ - \int \vartheta_0(\mathbf{x}) [\rho(\mathbf{x}, \tau_0) - \rho_0(\mathbf{x})] d^3 \mathbf{x}. \end{aligned} \quad (\text{D9})$$

To see that (D9) is indeed equivalent to minimizing (D1) under the constraints (D3) or (D5), (D2) and (D4), observe that for those ρ, \mathbf{J}, ϕ that do not satisfy the constraints, the sup operation over $\theta, \psi, \vartheta_{\text{in}}, \vartheta_0$ will give positive infinity; the sup will be finite (and thus contribute to the subsequent minimization) only if all constraints are satisfied. (This argument is the functional version of what is explained in Appendix C2 for the finite-dimensional case.)

Performing an integration by parts in the τ variable in (D9) and using the boundary conditions on the mass density (D4), we find that $\vartheta_{\text{in}}(\mathbf{x}) = \theta(\mathbf{x}, 0)$ and $\vartheta_0(\mathbf{x}) = \theta(\mathbf{x}, \tau_0)$. Integrating further by parts in the \mathbf{x} variable, assuming that boundary terms at infinity vanish (or that we have periodic boundary conditions in space) and rearranging terms, we obtain

$$\begin{aligned} \inf_{\rho, \mathbf{J}, \phi} \sup_{\substack{c, \mathbf{m}, \theta, \psi: \\ c + |\mathbf{m}|^2/2 \leq 0}} \iint d^3 \mathbf{x} d\tau \left[\rho \left(c \tau^{3/2} - \partial_{\tau} \theta - \frac{3}{2\tau} \psi \right) \right. \\ \left. + \mathbf{J} \cdot (\mathbf{m} \tau^{3/2} - \nabla_{\mathbf{x}} \theta) + \frac{3}{4\tau^{3/2}} |\nabla_{\mathbf{x}} \psi - \tau^{3/2} \nabla_{\mathbf{x}} \varphi_{\text{g}}|^2 \right. \\ \left. - \frac{3}{4\tau^{3/2}} |\nabla_{\mathbf{x}} \psi|^2 + \frac{3}{2\tau} \psi \right] \\ - \int \theta(\mathbf{x}, 0) d^3 \mathbf{x} + \int \theta(\mathbf{x}, \tau_0) \rho_0(\mathbf{x}) d^3 \mathbf{x}. \end{aligned} \quad (\text{D10})$$

Performing minimization with respect to ρ, \mathbf{J}, ϕ first, as in (C5) of Appendix C2, we see that the following two equalities must hold (remember that ρ need not be positive at this stage):

$$c = \frac{1}{\tau^{3/2}} \left(\partial_{\tau} \theta + \frac{3\psi}{2\tau} \right), \quad \mathbf{m} = \frac{1}{\tau^{3/2}} \nabla_{\mathbf{x}} \theta, \quad (\text{D11})$$

so that terms linear in ρ and \mathbf{J} vanish in (D10). It follows that c and \mathbf{m} are determined by θ and ψ and that the constraint $c + |\mathbf{m}|^2/2 \leq 0$ can be written as

$$\partial_{\tau} \theta + \frac{1}{2\tau^{3/2}} |\nabla_{\mathbf{x}} \theta|^2 + \frac{3}{2\tau} \psi \leq 0. \quad (\text{D12})$$

Also, the inf with respect to ϕ is straightforward and gives

$$\tau^{3/2} \nabla_{\mathbf{x}} \varphi_{\text{g}} = \nabla_{\mathbf{x}} \psi. \quad (\text{D13})$$

Using (D11) and (D13) in (D10), we arrive at the optimization problem of maximizing

$$\begin{aligned} J = \iint \left(\frac{3}{2\tau} \psi - \frac{3}{4\tau^{3/2}} |\nabla_{\mathbf{x}} \psi|^2 \right) d^3 \mathbf{x} d\tau \\ + \int \theta(\mathbf{x}, \tau_0) \rho_0(\mathbf{x}) d^3 \mathbf{x} - \int \theta(\mathbf{x}, 0) d^3 \mathbf{x} \end{aligned} \quad (\text{D14})$$

under constraint (D12). Equations (D14) and (D12) constitute a variational problem *dual* to the original one.

As both the original and the dual variational problems have the same saddle-point formulation (D9) or (D10), the optimal values of the two functionals (D1) and (D14) are equal. Let $(\rho, \mathbf{J}, \varphi_{\text{g}})$ be a solution to the original variational problem and θ, ψ be a solution to the dual one. Subtracting the (equal) optimal values from each other, we may now write, similarly to (C7),

$$\begin{aligned} \iint \left(\frac{\tau^{3/2}}{2\rho} |\mathbf{J}|^2 + \frac{3\tau^{3/2}}{4} |\nabla_{\mathbf{x}} \varphi_{\text{g}}|^2 \right. \\ \left. + \frac{3}{4\tau^{3/2}} |\nabla_{\mathbf{x}} \psi|^2 - \frac{3}{2\tau} \psi \right) d^3 \mathbf{x} d\tau \\ + \int \theta(\mathbf{x}, 0) d^3 \mathbf{x} - \int \theta(\mathbf{x}, \tau_0) \rho_0(\mathbf{x}) d^3 \mathbf{x} = 0. \end{aligned} \quad (\text{D15})$$

We are going to show that the left-hand side of (D15) may be given the form of a sum of three non-negative terms, each of which will therefore have to vanish. First, we rewrite the last two integrals, using the mass conservation constraint (D5) and integrations by parts, in the form

$$- \iint \partial_{\tau} (\theta \rho) d^3 \mathbf{x} d\tau = - \iint (\partial_{\tau} \theta \rho + \nabla_{\mathbf{x}} \theta \cdot \mathbf{J}) d^3 \mathbf{x} d\tau.$$

Secondly, we note that

$$\begin{aligned} \iint \left(\frac{3\tau^{3/2}}{4} |\nabla_{\mathbf{x}} \varphi_{\text{g}}|^2 + \frac{3}{4\tau^{3/2}} |\nabla_{\mathbf{x}} \psi|^2 \right) d^3 \mathbf{x} d\tau \\ = \iint \left[\frac{3}{4\tau^{3/2}} |\tau^{3/2} \nabla_{\mathbf{x}} \varphi_{\text{g}} - \nabla_{\mathbf{x}} \psi|^2 - \frac{3}{2\tau} \psi (\rho - 1) \right] d^3 \mathbf{x} d\tau, \end{aligned}$$

which follows from the Poisson constraint (D2). Taking all this into account in (D15), we obtain, after a rearrangement of terms,

$$\begin{aligned} \iint \frac{\rho}{2\tau^{3/2}} \left| \frac{\tau^{3/2}}{\rho} \mathbf{J} - \nabla_{\mathbf{x}} \theta \right|^2 d^3 \mathbf{x} d\tau \\ + \iint -\rho \left(\partial_{\tau} \theta + \frac{1}{2\tau^{3/2}} |\nabla_{\mathbf{x}} \theta|^2 + \frac{3}{2\tau} \psi \right) d^3 \mathbf{x} d\tau \\ + \iint \frac{3}{4\tau^{3/2}} |\tau^{3/2} \nabla_{\mathbf{x}} \varphi_{\text{g}} - \nabla_{\mathbf{x}} \psi|^2 d^3 \mathbf{x} d\tau = 0. \end{aligned} \quad (\text{D16})$$

The left-hand side is a sum of three non-negative terms (the second is so by D12), all of which must thus vanish. This gives

$$\mathbf{v} = \frac{1}{\rho} \mathbf{J} = \frac{1}{\tau^{3/2}} \nabla_x \theta, \quad \nabla_x \varphi_g = \frac{1}{\tau^{3/2}} \nabla_x \psi \quad (\text{D17})$$

and

$$\partial_\tau \theta + \frac{1}{2\tau^{3/2}} |\nabla_x \theta|^2 + \frac{3}{2\tau} \psi = 0, \quad (\text{D18})$$

wherever ρ is non-vanishing (otherwise the left-hand-side is non-positive by D12). The last equality turns into the Euler equation

$$\partial_\tau \mathbf{v} + (\mathbf{v} \cdot \nabla_x) \mathbf{v} = -\frac{3}{2\tau} (\mathbf{v} + \nabla_x \varphi_g) \quad (\text{D19})$$

by taking the gradient and using (D17).

By (D17) and (D18), any two hypothetically different minimizing solutions for either variational problem give rise to the same velocity potential and to the same gravitational potential (up to in-

significant constants) and thus define the same solution $(\rho, \mathbf{v}, \varphi_g)$ to the Euler–Poisson equations with the boundary conditions (D4) and the condition of curl-free velocity.

Moreover, for any such solution $(\rho, \mathbf{v}, \varphi_g)$, one can use (D17) to define θ and ψ that satisfy (D18) and thus (D12). By (D16), the values of functionals I and \bar{I} evaluated at these functions will coincide; together with convexity this implies, by an argument similar to that given in Appendix C2 concerning (C9), that such $(\rho, \mathbf{v}, \varphi_g)$ and (θ, ψ) in fact minimize both functionals under the corresponding constraints.

This means that a (curl-free) velocity field, a gravitational field and a density fields $(\mathbf{v}, \varphi_g, \rho)$ will satisfy the Euler–Poisson equations (A17)–(A19) (repeated as D19, D3, and D2 in this Appendix) and the boundary conditions (D4) if and only if they minimize (D1) under the corresponding constraints. This establishes uniqueness.

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