

CUBE: A MEMORY-LITE COSMOLOGICAL N -BODY ALGORITHM

HAO-RAN YU^{1,2}, UE-LI PEN^{2,3,4,5}, XIN WANG²

¹Tsung-Dao Lee Institute, Shanghai Jiao Tong University, Shanghai, 200240, China

²Canadian Institute for Theoretical Astrophysics, University of Toronto, Toronto, ON M5H 3H8 Canada

³Dunlap Institute for Astronomy and Astrophysics, University of Toronto, Toronto, ON M5S 3H4, Canada

⁴Canadian Institute for Advanced Research, CIFAR Program in Gravitation and Cosmology, Toronto, ON M5G 1Z8, Canada

⁵Perimeter Institute for Theoretical Physics, Waterloo, ON N2L 2Y5, Canada

Draft version September 22, 2017

Abstract

Cosmological large scale structure N -body simulations are computation-light, memory-heavy problems in supercomputing. Traditional N -body simulation algorithms use at least 24-byte memory per particle, of which six 4-byte single precision floating point numbers keep track phase space coordinates of each particle. Here we present the algorithm and accuracy of a new parallel, memory-lite, Particle-Mesh based N -body code, where each particle can occupy as low as 6-byte memory. This is accomplished by storing relative position and relative velocity of each particle, in the format of 1-byte-integer, respect to their averaged value of a mesh-grid. The remaining information is given by complimentary density and velocity fields, which are negligible in memory space, and proper ordering of particles, which gives no extra memory. Our numerical experiments show that this integer based N -body algorithm provides acceptable accuracy compared to traditional algorithm in cosmological simulations. This significant lowering of memory-to-computation ratio breaks the bottleneck of scaling up and speeding up large cosmological N -body simulations on multi-core and heterogenous computing systems.

1. INTRODUCTION

N -body simulation is a powerful tool to solve the highly nonlinear dynamic problems (Hockney & Eastwood 1988). It is widely used in cosmology to model the formation and evolution of the large scale structure (LSS). With the fast development of parallel supercomputers, we are able to simulate a system of more than a trillion (10^{12}) N -body particles. To date the largest N -body simulation in application is the “TianNu” simulation (Yu et al. 2017; Emberson et al. 2017) run on the TianHe-2 supercomputer by cosmological simulation code CUBEP3M (Harnois-Déraps et al. 2013). It uses nearly 3×10^{12} particles to simulate the cold dark matter (CDM) and cosmic neutrino evolution through the cosmic age.

The N -body simulations use considerable amount of memory, because the phase space coordinates (x, y, z, v_x, v_y, v_z) of each N -body particle must be stored as, at least, six single precision floating point numbers (24 bytes). Contrarily, their computing workload can be alleviated by many algorithms, like Particle-Mesh [cite] and Tree [cite], scale as $o(N \log N)$ or even $o(N)$. On the other hand, modern supercomputer systems use multi cores, many integrated cores (MIC) and even densely parallelized GPUs, bringing orders of magnitude higher computing power, whereas these architectures usually have limited memory allocation. Thus, the computation-light but memory-heavy applications, compared to matrix multiplication and decomposition calculations, are less suitable for fully usage of the computing power of modern supercomputers. For example, although native and offload modes of CUBEP3M are able to run on the Intel Xeon-PHI MIC architectures, with the require ment

of enough memory, TianNu simulation were done on TianHe-2 with only its CPUs – 73% of the total memory but only 13% of the total computing power.

We present a new N -body simulation code CUBE, using as low as 6 byte per particle (here after we use “bpp” referring to “byte per particle”). We show that it gives accurate results in cosmological LSS simulations. The method is presented in §2, and a comparison between this method and traditional method is shown in §3. Discussions and conclusions are in §4.

2. METHOD

The most memory consuming part of a N -body simulation is usually the phase space coordinates of N -body particles – 24 bpp (4 single precision floating numbers) must be used to store each particles’ 3D position and velocity vectors. CUBEP3M, an example of a memory-efficient parallel N -body code, can use as low as 40 bpp in sacrificing computing speed (Harnois-Déraps et al. 2013). This includes the phase coordinates (24 bpp) for particles in physical domain and buffered region, a linked list (4 bpp), and a global coarse mesh and local fine mesh. 4-byte real numbers are not necessarily adequate in representing the *global* coordinates in simulations. If the box size is many orders of magnitude larger than interactive distance between particles, especially in the field of resimulation of dense subregions, double precision (8-byte) coordinates are needed to avoid truncation errors. Another solution is to record *relative* coordinates for both position and velocity. CUBE replaces the coordinates and linked list 24+4=28 bpp memory usage with an integer based storage, reduces the basic memory usage from 28 bpp down to 6 bpp, described as following 2.1 and 2.2. The algorithm is described in 2.3.

2.1. Particle position storage

We construct a uniform mesh throughout the space and each particle belongs to its parent cell of the mesh. Instead of storing global coordinates of each particle, we store its offset relative to its parent cell which contains the particle. We divide the cell, in each dimension d , evenly into $2^8 = 256$ bins, and use a 1-byte (8 bits) integer $\chi_d \in \{-128, -127, \dots, 127\}$ to indicate which bin it locates in this dimension. The global locations of particles are given by cell-ordered format in memory space, and a complimentary number count of particle number in this mesh (density field) will give complete information of particle distribution in the mesh. Then the global coordinate in d th dimension x_d is given by $x_d = (n_c - 1) + (\chi_d + 128 + 1/2)/256$, where $n_c = 1, 2, \dots, N_c$ is the index of the coarse grid. The mesh is chosen to be coarse enough such that the density field takes negligible memory. This coarse density field can be further compressed into 1-byte integer format, such that a 1-byte integer show the particle number in this coarse cell in range 0 to 255. In the densest cells (rarely happened) where there are ≥ 255 particles, we can just write 255, and write the actual number as a 4-byte integer in another file.

In a simulation with volume L^3 and N_c^3 coarse cells, particle positions are stored with a precision of $L/(256N_c)$. The force calculation (e.g. softening length) should be configured much finer than this resolution, discussed in later sections. On the other hand, particle position can also be stored as 2-byte (16 bits) integers to increase the resolution. In this case, each coarse cell is divided into $2^{16} = 65536$ bins and the position precision is $L/(65536N_c)$, precise enough compared to using 4-byte global coordinates, see later results. We denote this case “x2” and denote using 1-byte integers for positions “x1”.

We collectively write the general position conversion formulae

$$\chi_d = [2^{8n_x}(x_d - [x_d])] - 2^{8n_x-1}, \quad (1)$$

$$x_d = (n_c - 1) + 2^{-8n_x}(\chi_d + 2^{8n_x-1} + 1/2), \quad (2)$$

where $[]$ is the operator to take the integer part. $n_x \in \{1, 2\}$ is the number of bytes used for each integer, x_d and χ_d are floating and integer version of the coordinate. The velocity counterpart of them are $n_v = 1, 2$, v_d and ν_d . $n_c = 1, 2, \dots$ is the coarse grid index, given by the ordering of particles and a integer based particle density field (see 2.3.1). The position resolution for n_x -byte integer, “ xn_x ”, is $2^{-8n_x}L/N_c$.

As a $n_x = 1$, 1D ($d = 1$), 4-coarse-cell ($N_c = 4$) example, if

$$\chi_1 = (-128, 127, 0, 60),$$

particle number density

$$\rho_c^{1D} = (1, 0, 2, 1),$$

then in unit of coarse cells,

$$x_1 = (0.001953125, 2.998046875, 2.501953125, 3.736328125).$$

2.2. Particle velocity storage

Similarly, actual velocity in d th dimension v_d is decomposed into an averaged velocity field on the same coarse grid v_c and a residual Δv relative to this field:

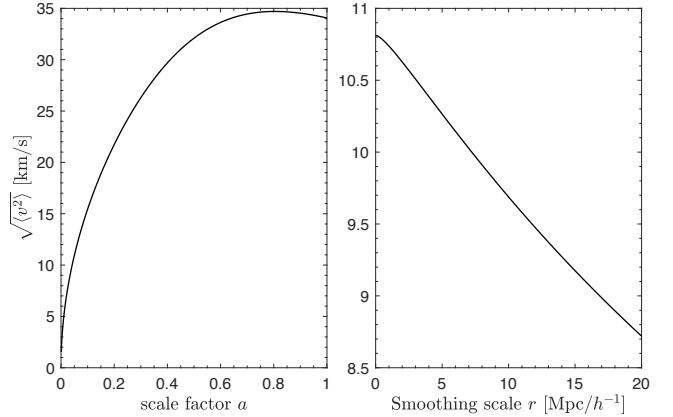


FIG. 1.— please check units

$v_d = v_c + \Delta v$. v_c is always recorded and kept updated, and occupies negligible memory. We then divide velocity space Δv into uneven bins, and use a n_v -byte integer to indicate which Δv bin the particle is located.

The reason why we use uneven bins is that, slower particles are more abundant compared to faster ones, and one should better resolve slower particles tracing at least linear evolution. On the other hand, there could be extreme scattering particles (in case of particle-particle force), and we can safely ignore or less resolve those non-physical particles. One of the solution is that, if we know the probability distribution function (PDF) $f(\Delta v)$ we divide its cumulative distribution function (CDF) $F(\Delta v) \in (0, 1)$ into 2^{8n_v} bins to determine the boundary of Δv bins, and particles should evenly distribute in the corresponding uneven Δv bins. Practically we find that either $f(v_d)$ or $f(\Delta v)$ is close to Gaussian, so we can use Gaussian CDF, or any convenient analytic functions which are close to Gaussian, to convert velocity between real numbers and integers.

The essential parameter of the velocity distribution is its variance. On non-linear scale, the velocity distribution function is non-Gaussian. However, to the first order approximation, we simply assume it as Gaussian and characterized by the variance

$$\sigma_v(a, k) = \frac{1}{3}(aHfD)^2 \int_0^k d^3 q \frac{P_L(q)}{q^2}, \quad (3)$$

where $a(z)$ is the scale factor, $H(z)$ the Hubble parameter, D is the linear growth factor, $f = d \ln D / d \ln a$, and $P_L(k)$ is the linear power spectrum of density contrast at redshift zero. $\sigma_v(a, k)$ is a function of cosmic evolution a and a smoothing scale k , or r (see Figure 1). Δv is the velocity dispersion relative to the coarse grid, so we approximate its variance as

$$\sigma_\Delta^2(a) = \sigma_v^2(a, r_c) - \sigma_v^2(a, r_p), \quad (4)$$

where r_c is the scale of coarse grid, and r_p is the scale of average particle separation. In each dimension of 3D velocity field, we use $\sigma_\Delta^2(a)/3$ according to the equipartition theorem. On different scales, we measure the statistics of v_d , v_c and Δv and find good agreement with the above model.

The simulation results are very insensitive if we manually tune the variance of the model σ_Δ within an order of magnitude. However, the method of using



FIG. 2.— Spacial decomposition in CUBE in a 2D analogy. In this example, there are 2 images per dimension ($\text{nn} = 2$), and 2 tiles per image per dimension ($\text{nnt} = 2$). The orange boxes show the overlapped physical+buffer regions, inside of which one physical regions is indicated with green.

uneven bins gets much better results than simply using equal bins between minimum and maximum values [$\min(\Delta v)$, $\max(\Delta v)$]. So, one can safely use a standard Λ CDM (Cold Dark Matter with a cosmological constant) for slightly different cosmological models, in equation (3). In CUBE, the velocity conversion takes the formula

$$\nu_d = \left\lfloor (2^{8n_\nu} - 1)\pi^{-1} \tan^{-1} \left(v_d \sqrt{\pi/2\sigma_\Delta^2} \right) \right\rfloor, \quad (5)$$

$$v_d = \tan \left(\frac{\pi\nu_d}{2^{8n_\nu} - 1} \right) \sqrt{2\sigma_\Delta^2/\pi}, \quad (6)$$

where $\lfloor \cdot \rfloor$ is the operator to take the nearest integer. Tangent functions are convenient and compute very fast. Compared to error functions used in Gaussian case, they take the same variance at $v_d = 0$ but resolve high velocities relatively better. Note again that proper choice of conversion formulae and σ_Δ only optimizes the velocity space sampling, but does not affect the physics.

Initially, particles are generated by initial condition generator, at a higher redshift. The coarse grid density field v_c is also generated at this step by averaging all particles in the coarse cell. A global σ_Δ is calculated by equation (4), where linear approximation is hold. Then velocities are stored by equation (5). During the simulation, v_c is updated every time step, and a nonlinear σ_Δ is measured directly from the simulation, and can be simply used in the next time step, after scaled by the ratio of growth factors between two adjacent time steps. More details see section 2.3.2.

2.3. Code overview

CUBE uses a 2-level PM force calculation, same as CUBEP3M. However, in order to apply the integer based format to the N -body simulation, substantial structural changes need to be done. CUBE is written in Coarray Fortran, where Coarray features implement MPI com-

munication between computation nodes/images¹. The algorithm is described in this language.

2.3.1. Spacial decomposition

CUBE uses cubic decomposition structures. The global simulation volume is decomposed into nn^3 cubic sub-volumes, and each of these are assigned to a coarray *image*. Inside of an image, the sub-volume is further decomposed into nnt^3 cubic *tiles*. Each tile is surrounded by a *buffer* region which is ncb coarse cells thick. The buffer is designed for two reasons: (1) computing the fine mesh force, whose cut-off $\text{nforce_cutoff} \leq \text{ncb}$, and (2) collecting all possible particles travelling from a tile's buffer region to its center, *physical* region. Thus, a integer-based coarse mesh particle number density (ρ_c) is defined as

```
integer(1) rho_c(nex,nex,nex,nnt,nnt)[nn,nn,*]
```

where $\text{nex} = \text{nt} + 2 \times \text{ncb}$ covers the buffer region on both sides, nnt is the tile dimensions, and nn is the image *codimensions*². In the buffer region, phase space coordinates of particles xp and vp are copied from adjacent tiles and images, when necessary. Figure 2 shows the spacial decomposition in a 2-dimensional analogy, with $\text{nn} = 2$ and $\text{nnt} = 2$.

The particle position and velocity arrays xp and vp (equivalent to χ_d and ν_d in 2.1 and 2.2 respectively) are required to be sorted according to the same memory layout of `rho_c`, such that n_c and thus global positions of particles x_d can be obtained by equation (2).

2.3.2. Algorithm

Figure 3 shows the overall structure of the code and these subroutines are described in following paragraphs.

initialize and **read_particles** – The subroutine `initialize` creates necessary FFT plans and read in configuration files telling the program at which redshifts we need to do checkpoints, halofinds, or stop the simulation. Force kernels are also created or read in here. In `read_particles`, for each image, we read in all particles in *physical* regions of every tile (indicated in Figure 2), i.e., their initial positions xp , velocities vp in the compressed format, and their corresponding, physical part of the number density field `rho_c`. These are obtained by the initial condition generator (see Appendix A). Because physical regions of tile are *complete* and *disjoint* in space, particles at this stage are also complete and disjoint. We call it “disjoint state”. At this stage, `rho_c`’s buffer regions of each tile are 0, and the elements of xp and vp beyond physical number can be arbitrary and are not used.

buffer_density, **buffer_x** and **buffer_v** – In order to use the integer based format, xp and vp must always be ordered, and their number density field `rho_c` must always be present. In updating arrays of xp and vp (not simultaneously) of a local tile, a vicinity (buffer) region of the physical region is also needed. First, buffer regions of `rho_c` is synchronized between tiles and images by subroutine `buffer_density`. Then, by subroutines

¹ Images are the concept of computing nodes or MPI tasks in Coarray Fortran. We use this terminology in this paper.

² Coarray Fortran concept. Codimensions can do communications between images.

```

program CUBE
  call initialize
  call read_particles
  sync all
  call buffer_density
  call buffer_xp
  call buffer_vp
  do
    call timestep
    call update_xp
  sys
    call buffer_xp
    call update_vp
    call buffer_vp
    if(checkpoint_step) then
      call update_xp
      call checkpoint
      if (final_step) exit
      call buffer_density
      call buffer_xp
      call buffer_vp
    endif
  enddo
  call finalize
end

```

FIG. 3.— Overall structure of CUBE.

buffer_x and **buffer_v** respectively, **xp** and **vp** are updated to contain common, buffered particles, in an order according to the new, buffered **rho_c**. We call this stage “buffered state”. After **particle_initialize** is done by all images (synchronized by “**sync all**”, which is equivalent to a **mpi_barrier**), these three subroutines are called and particles are converted from disjoint state to buffered state.

timestep – We operate a Runge-Kutta 2 method for time integration. i.e., we update position (D =drift) and velocity (K =kick) every half time step. For n time steps, the operation would be $(DKKD)^n$ which is 2nd order accurate. The actual simulation applies varied time steps. In each iteration of the main loop, we firstly call **timestep**, where a increment of time **dt** is controlled by particles’ maximum velocities, accelerations, and cosmic expansions.

update_xp – According to **dt**, subroutine **update_xp** updates the positions of particles in a “gather” algorithm (in contrast, CUBEP3M uses “scatter” algorithm) tile by tile. Because each tile is in the buffered state with buffer depth **ncb**, we are able to collect all possible particles whose $vp^*dt < ncb$ from physical+buffer region to its physical region. In order to keep particles ordered, we have to execute $xp = xp + vp * dt$ twice, first time to obtain an updated density field on the tile, **rho_c_new**, and second time to generate a new, local particle list **xp_new** and **vp_new** on the tile. The reason for this repeated execution is the dependence of **x_new** and **v_new**’s value and ordering on both the old **xp**, **vp** and the new **rho_c_new**. However, $xp = xp + vp * dt$ scales as $O(N)$ and is computational inexpensive. Although **xp_new** and **vp_new** arrays take extra memory, they do not appear simultaneously for multiple tiles, and compared to **xp** and **vp**, the extra memory overhead is by a factor of $1/nnt^3$. This requires the “gather” algorithm to move particles – as we synchronize the correct **{rho_c, xp, vp}** in the buffer region of each tile, and the particles in the extended (physical+buffer) region are a superset of particles locating in the physical region in the next time step,

given that $vp * dt$ is not greater than the buffer depth. In **update_xp**, **{rho_c_new, xp_new, vp_new}** can be set to collect only the physical region of the tile, with the rest of this subset discarded. At the end, **{rho_c, xp, vp}** is replaced by **{xp_new, vp_new, rho_c_new, vfield}**, and the memory of latter is freed up. After this step, the tile does not contain buffered regions, being a disjoint state.

update_vp – After drift D , we need to update velocities (kick K) of particles in physical region of each tile. Before this step, we need to call **buffer_density** and **buffer_xp** in order that particle positions are in buffered state. CUBE uses a 2-level particle mesh scheme (Harnois-Déraps et al. 2013). Local fine forces have a force cut-off **nforce_cutoff** $\leq ncb$. So, if we apply a fine-grid particle-mesh on an extended tile with buffered-depth **ncb** (see Figure 2 regard it periodic), the force on physical regions does not depend on the false periodic boundary assumption, and the resulting fine force and velocity update on physical region is correct.

The compensating coarse grid force is globally computed by using a coarser (usually by factor of 4) mesh by dimensional splitting – a distributed-memory cubic decomposed 3D coarse density field is interpolated by particles, and we Fourier transform data in consecutive three dimensions with global transposition in between (known as the pencil decomposition). After the multiplication of force kernels, the inverse transform takes place to get the cubic distributed coarse force field, upon which velocities are updated again.

An optional particle-particle (PP) force can be called to increase the force resolution and the velocities are updated last time according to this. **v field**. During each of these force calculations, maximum accelerations are collected serving for the **timestep** in the next iteration, controlling next **dt**.

After all possible velocity updates, **vp** in physical regions are correctly updated. Particle locations in the buffer regions has updated before PM and remained unchanged. So we simply call **buffer_v** again such that the **update_x** in the next iteration will be done correctly. We call **update_vp** to convert **vp** into buffered state.

checkpoint – If a desired redshift is reached, we execute the last drift step in the $(DKKD)^n$ operation by **update_xp**, and call **checkpoint** to save the disjoint state of **{xp_new, vp_new, rho_c_new, vfield}** on disk. Related operations, like run-time halo finder, projections are also done at this point. If final desired redshift is reached, **final_step** let us exit the main loop. These corresponding logical variables are controlled in **timestep**.

finalize – Finally, in **finalize** subroutine we destroy all the FFT plans and finish up any timing or statistics taken in the simulation.

2.4. Memory layout

The memory overhead of fine force calculation is controlled by $(1/nnt)^3$. The memory usage of coarse arrays are negligible.

Here we summarize the memory usage in unit of bpp (byte per particle).

3. RESULTS

We use the same initial conditions, and use CUBEP3M and CUBE to simulate the large scale structure formation separately and compare their results.

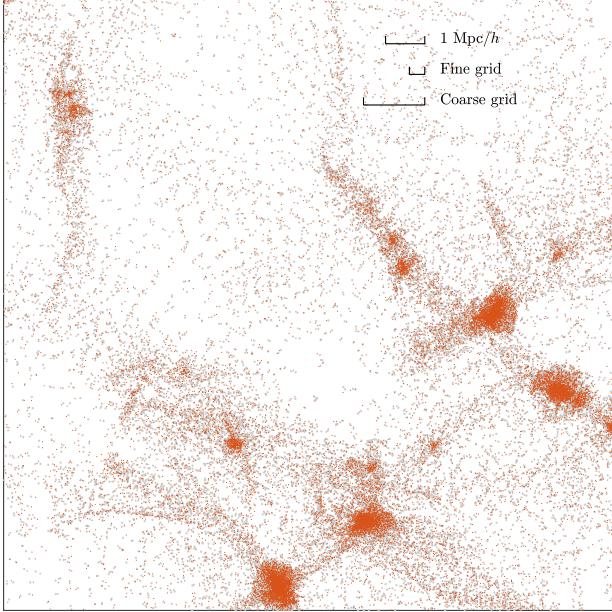


FIG. 4.— particles.

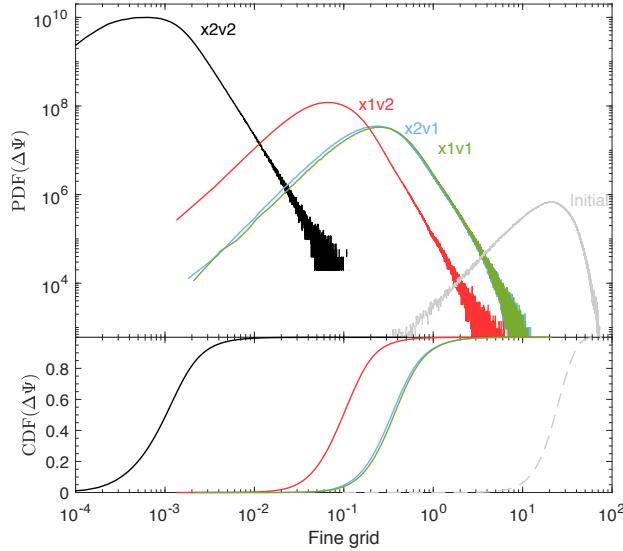


FIG. 5.— Displacement.

Results of `izipx=1`
Results of `izipx=2`

4. DISCUSSION AND CONCLUSION

Cosmological simulation, PM method.

PID.

Summarize memory.

Phi, GPU.

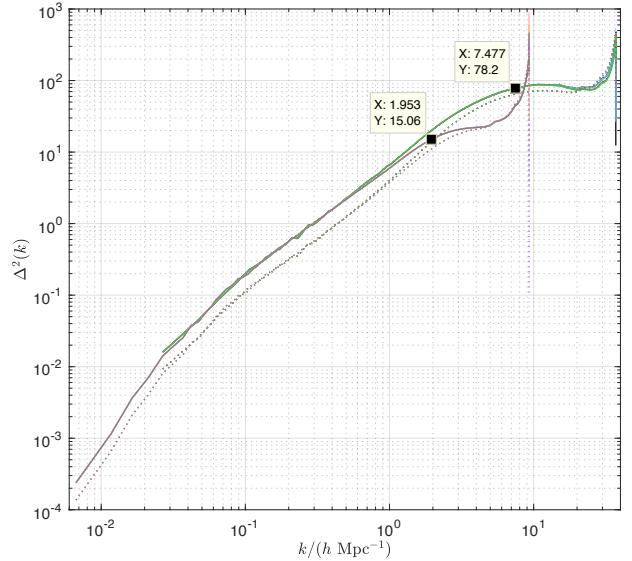


FIG. 6.— power spectrum.

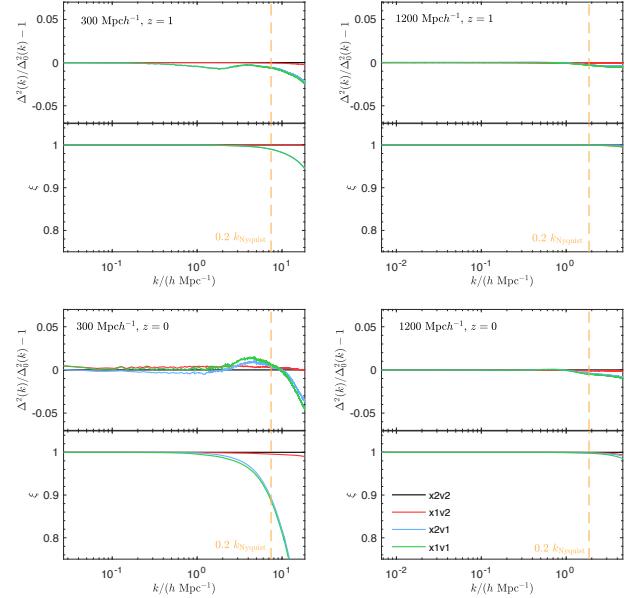


FIG. 7.— power spectrum and cross correlation.

APPENDIX

A. INITIAL CONDITION GENERATOR

REFERENCES

- Emberson, J. D. et al. 2017, Research in Astronomy and Astrophysics, 17, 085, 1611.01545
Harnois-Déraps, J., Pen, U.-L., Iliev, I. T., Merz, H., Emberson, J. D., & Desjacques, V. 2013, MNRAS, 436, 540, 1208.5098
Hockney, R. W., & Eastwood, J. W. 1988, Computer simulation using particles
Yu, H.-R. et al. 2017, Nature Astronomy, 1, 0143, 1609.08968