



# Carpooling to solve the cosmological simulation bottleneck

**CARPool**: fast, accurate computation of large-scale structure statistics by pairing costly and cheap cosmological simulations (arXiv:2009.08970)

N.Chartier, B.Wandelt, Y.Akrami and F.Villaescusa-Navarro

## **About observations and observables**

**Data sets of next generation galaxy surveys:** unmatched statistical power to constrain initial perturbations, cosmic structure growth and expansion history



Euclid Space Telescope, DESI, Rubin Observatory LSST, Subaru HSC & PFS, SKA, WFIRST, SPHEREx...

# Motivation

# We need theoretical predictions of large-scale structure (LSS) statistics.

#### **Possible routes**

- **Costly N-body codes**, unmatched for the non-linear regime of structure growth (*GADGET, GreeM*, *HACC*, PKDGRAV3...)
- Analytical computations with *LPT*, *SFT*, *EFT*...
- Approximate solvers (*surrogates*): Particle-Mesh codes (**PM**), emulators, Neural Networks...

Accuracy is traded for computational speed (especially in the non-linear regime), statistical unbiasedness not guaranteed...

# Statistics of observables from N-body simulations

Fractional overdensity field; z = 0.5



• Observables (bins) are collected into a vector  $\boldsymbol{y}$  of size p

• Estimator  $\hat{\boldsymbol{\mu}}$  of  $\mathbb{E}[\boldsymbol{y}]$ ?

#### **Random events and estimation**

 $y_1, \ldots, y_N$  are N independent random realizations sampled on seeds  $r_1, \ldots, r_N$  1  $\sum_{n=1}^{N}$ 

Estimation of the mean  $\mathbb{E}[\boldsymbol{y}] = \boldsymbol{\mu} \in \mathbb{R}^p$  with  $\bar{\boldsymbol{y}} = \frac{1}{N} \sum_{n=1}^N \boldsymbol{y_n}$ 

Standard deviation of each element  $\bar{y}_i$  decreases as  $\mathcal{O}(N^{-\frac{1}{2}})$ 





# Can we get the best of both worlds? An unbiased and faster estimator



# **"CARPool" Method:** Variance reduction with "N-body + surrogate" pairs

# **Numerical analysis** What we are going to use

 $\Lambda$ CDM cosmology, Redshift z=0.5

*Y* ----- N-body simulations – GADGET – are from the *Quijote Simulations* (Villaescusa-Navarro et al., 2019)

C — The cheap *surrogate* is L-PICOLA (Howlett, 2015 b), an *MPI* implementation of COLA (Tassev, 2013)

# An example

## Matter power spectrum 95 linearly spaced bins with :

 $k_{max} = 1.184 \ h \text{Mpc}^{-1}$   $\Delta k = 3.147 \text{e-}2 \ h \text{Mpc}^{-1}$ 

### **CARPool estimate Vs. N-body only**



#### **CARPool estimate Vs. N-body only**



# What is the trick?

### **Control Variates principle**

Observables from simulations:

N-body code"cheap" surrogate
$$\boldsymbol{y} = \begin{pmatrix} y_1 & y_2 & \dots & y_p \end{pmatrix}^T$$
 $\boldsymbol{c} = \begin{pmatrix} c_1 & c_2 & \dots & c_q \end{pmatrix}^T$  $\mathbb{E} [\boldsymbol{y}] = \boldsymbol{\mu}$  $\mathbb{E} [\boldsymbol{c}] = \boldsymbol{\mu}_{\boldsymbol{c}}$ Unknown truth(Un)known "wrong" truth

• Intuition with two random scalars:

$$\sigma_{y+c}^2 = \sigma_y^2 + \sigma_c^2 + 2\text{cov}(y,c)$$

#### **Control Variates for simulations**

• Scalar case ("bin per bin"):

WE DON'T CARE ABOUT THE BIAS OF THE CHEAP ESTIMATOR

#### **Control Variates for simulations**

•



Proof in Rubinstein & Marcus (1985)

Multivariate case:  $m{x}(m{eta}) = m{y}$  -  $m{eta} \left(m{c} - m{\mu_c}
ight), m{eta} \in \mathbb{R}^{p imes q}$ 

Error box

 $\frac{\det\left(\boldsymbol{\Sigma}_{\boldsymbol{x}(\boldsymbol{\beta})\boldsymbol{x}(\boldsymbol{\beta})}\right)}{\det\left(\boldsymbol{\Sigma}_{\boldsymbol{y}\boldsymbol{y}}\right)} = \prod_{i=1}^{s=rank(\boldsymbol{\Sigma}_{\boldsymbol{y}\boldsymbol{\alpha}})}$ 



Squared canonical crosscorrelations

- 1) The estimate is unbiased by construction
- 2) The control matrix/coefficient gives optimal variance reduction
- 3) The more correlated the full simulation and the surrogate statistics, the better

# CARPool

- In practice:
  - $oldsymbol{eta}^{\star}$  must be estimated with data
  - $\mu_c$  is unknown

Convergence Acceleration by Regression and Pooling

**1** Estimate  $\bar{\mu}_c$  from M fast surrogates

2 With N "simulation + surrogate" pairs, compute  $\bar{x}(\hat{\beta}) = \bar{y} - \hat{\beta} (\bar{c} - \overline{\mu}_{c})$ 



• N-body sims only

 $\bar{\mathbf{y}}$  $\mathbf{y_n}$ n=1



# Back to the first example

## Matter power spectrum 95 linearly spaced bins with :

 $k_{max} = 1.184 \ h \text{Mpc}^{-1}$   $\Delta k = 3.147 \text{e-}2 \ h \text{Mpc}^{-1}$ 

#### **CARPool estimate Vs. N-body only**



## **Confidence in CARPool estimate (Pk)**



## **Generalized variance reduction (Pk)**



### **Standard deviation reduction (Pk)**



# Matter Bispectrum 73 squeezed triangle configurations



### **Standard deviation reduction (Bk)**



# Matter reduced Bispectrum 40 equilateral triangle configurations

$$k_1 = k_2 = k_3$$

### **Confidence in CARPool estimate (Qk)**



## Matter PDF

# 70 bins $\rho/\bar{\rho} \in [0.08, 50]$

#### **Univariate CARPool for PDF**





## **Conclusion and discussion**

- CARPool reduces variance by factors 10 to 100, even in the nonlinear regime.
- With only 5 GADGET-III simulations, CARPool is able to compute Fourier-space two-point and three-point functions of the matter distribution at a precision comparable to 500 GADGET-III simulations.
- We have variance reduction even for the matter PDF. The remapping technique proposed by *Leclercq et al. (2013)*, that increases the correlation between LPT-evolved density fields and simulations, can improve the chosen surrogate.
- CARPool can be implemented with various "N-body + surrogate" pairs. All you need is:

 An inexpensive surrogate and statistics computation.
 Strong correlation with the costly simulations.

# Thank you for your attention!

# (backup slides)

### **Generalized variance reduction PDF**



### The smoothing trick





