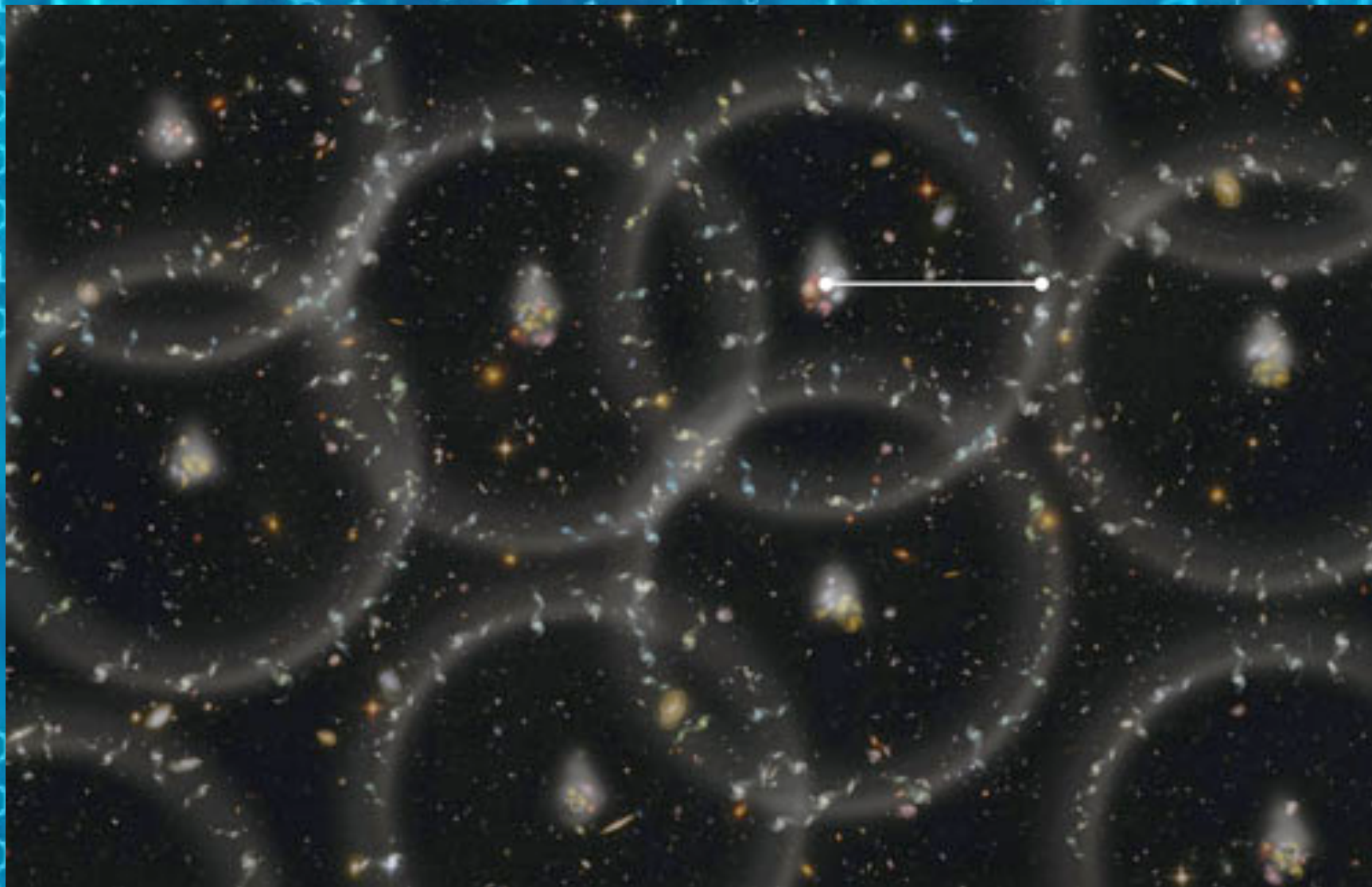# Graph Databases Solution for Higher Order Statistics in the Big Data Era of Astronomy

## Cristiano Sabiu（Yonsei University）

Juhan Kim（KIAS）, Xiao-Dong Li（Sun Yat-Sen）, Ben Hoyle（LMU）

CosKASI conference III "The Correlated Universe", Jeju, April 23-27

# Contents

- **Graph Databases**
    - **Introduction**
    - **Higher Order Spatial Clustering**
    - **Constructing Graph**
    - **Querying Graph**

- **Applications**
    - **Baryon Acoustic Oscillations**
    - **Cross-correlation with galaxy properties**
    - **Early Universe and Isocurvature Perturbations**
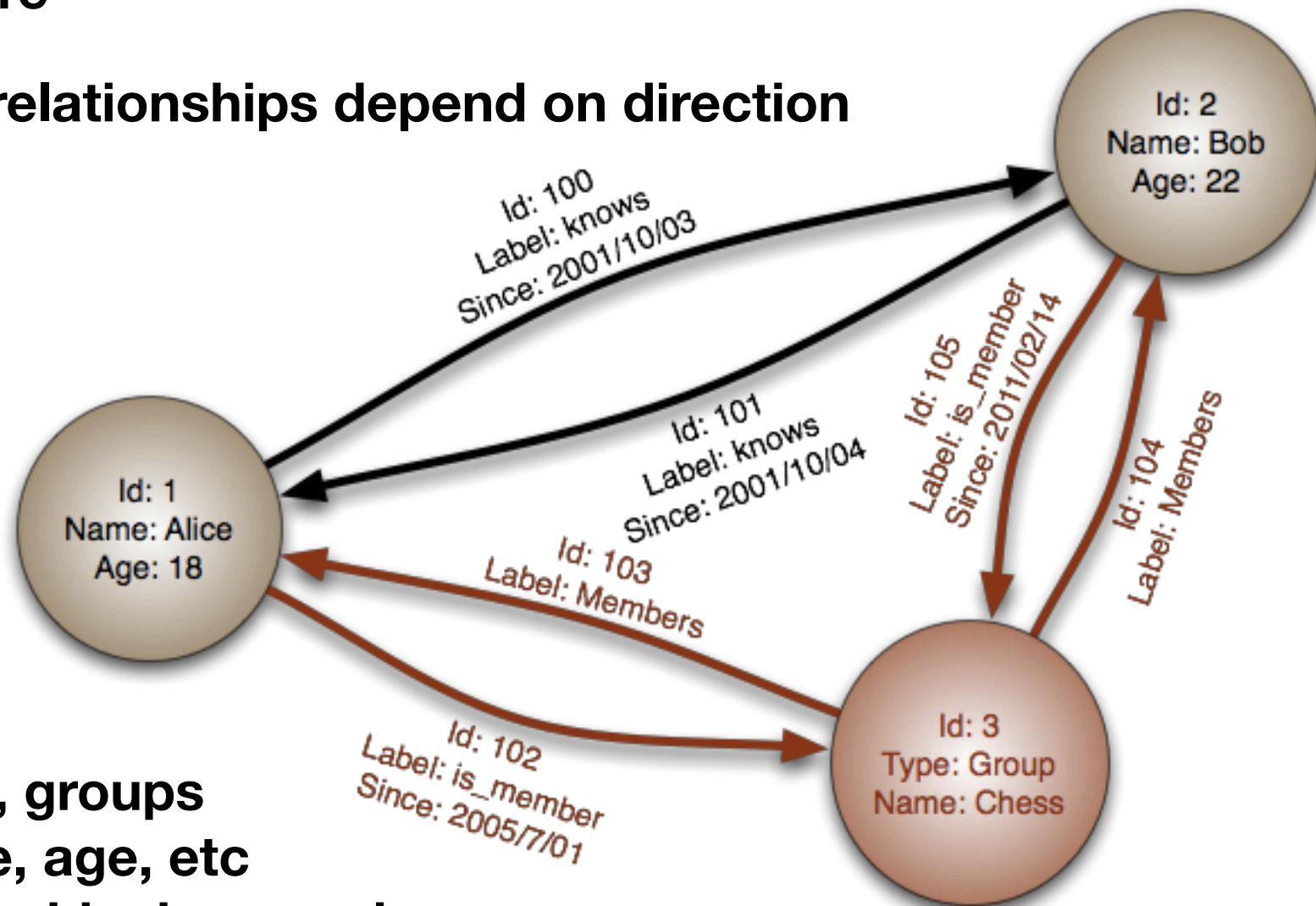
- **Conclusions**

# Graph Databases

A graph within graph databases is based on graph theory.

- <u>Nodes</u> represent entities or instances such as people, businesses, accounts, or any other item to be tracked. They are roughly the equivalent of the record, relation or row in a relational database, or MySQL item

- <u>Relationships</u>, are the lines that connect nodes to other nodes; representing the relationship between them. Meaningful patterns emerge when examining the connections and interconnections of nodes, properties and relationships. Relationships are the key concept in graph databases, representing an abstraction that is not directly implemented in a relational model like SQL

- <u>Properties</u> are germane information to nodes. For example, if *Wikipedia* were one of the nodes, it might be tied to properties such as website, reference material, or words that starts with the letter *w*, depending on which aspects of Wikipedia are germane to a given database.

# Graph Databases

**In 2009 <u>Facebook</u> gave up using their MySQL storage and moved to a graph structure**

**A <u>directed graph</u>: relationships depend on direction**
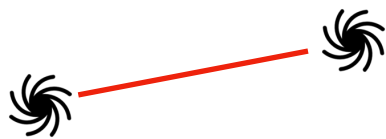


**<u>Nodes</u>: people, places, groups**
**Node <u>properties</u>: name, age, etc**
**<u>Relationships</u>: membership, known since**
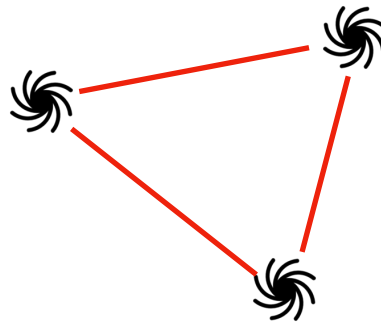
**<u>Graph for our purposes will be much simpler!</u>**

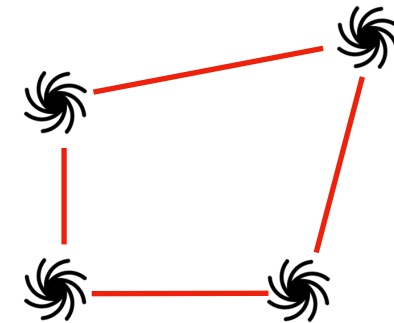# Graph Database solution for Galaxy Clustering statistics
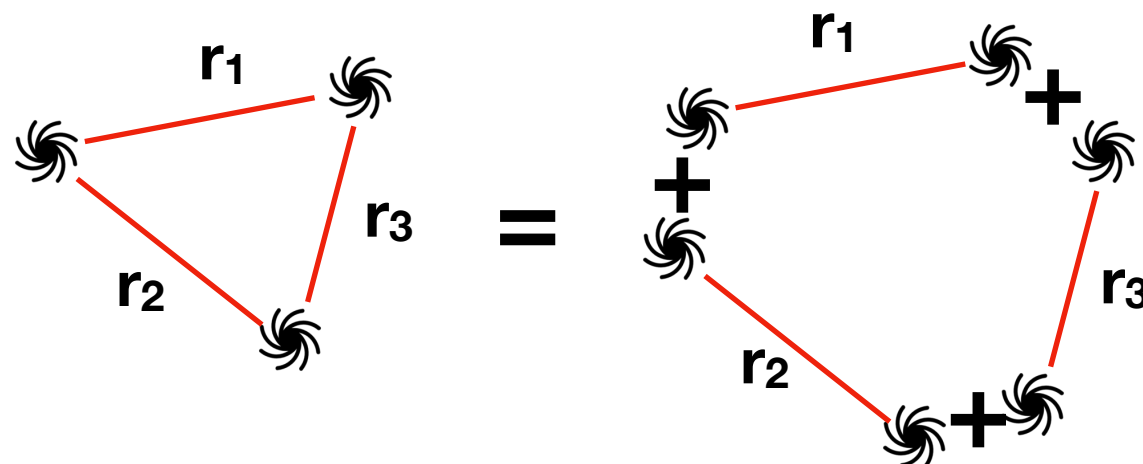
**2-point**        **3-point**        **4-point**

- The computational expensive part of measuring higher order statistics is in making sure the distances between data points satisfy our specific criteria i.e. $r_1$, $r_2$, $r_3$, etc

- However we notice that all higher order statistics are just complex combinations 2-point statistics

$r_1$  $r_3$  $r_2$  $=$  $+$  $r_1$  $+$  $r_3$  $r_2$  $+$

So rather than treat data positions as the important feature, rather treat each data pair (relationship) as the main information!

# Graph Database solution for Galaxy Clustering statistics
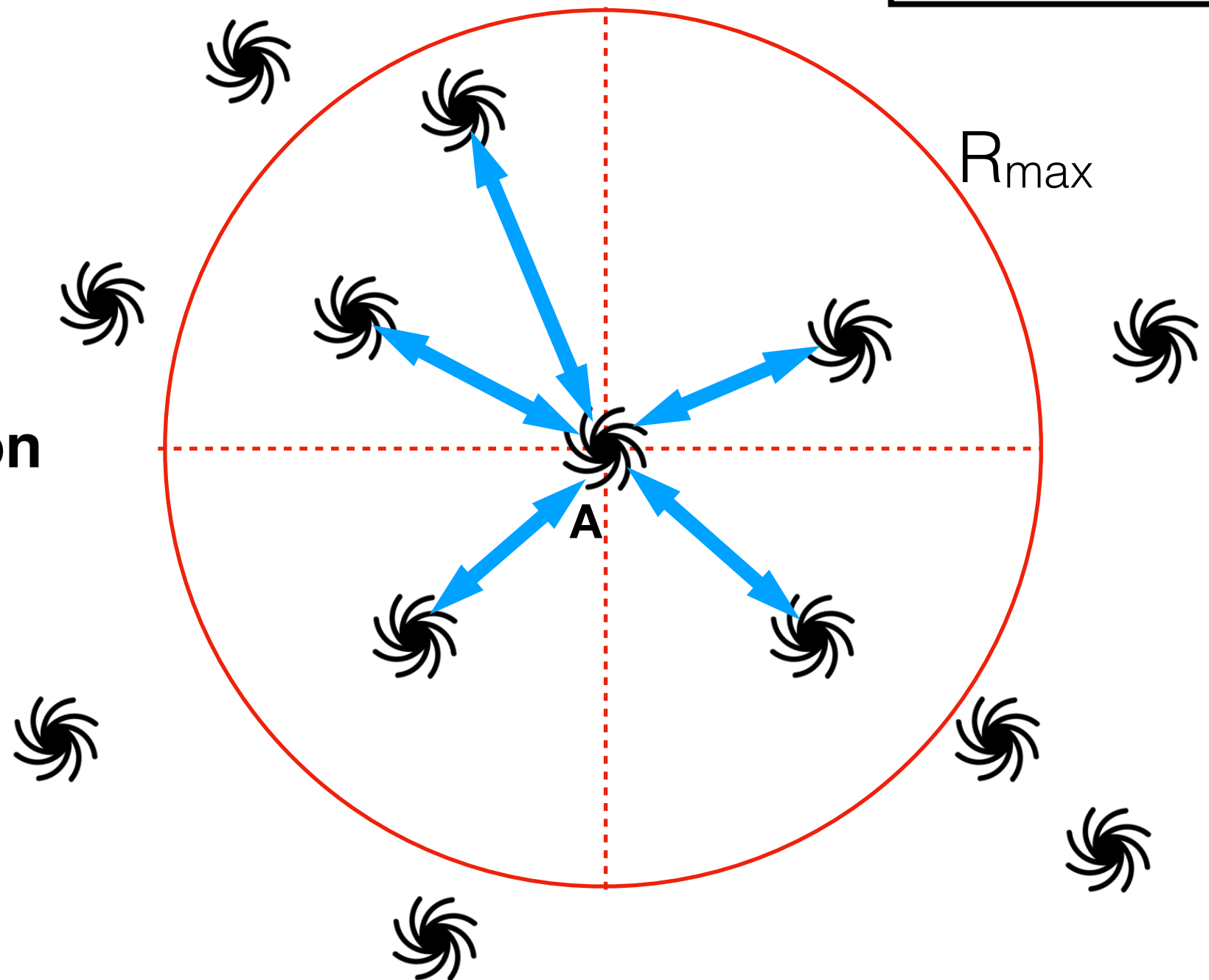
**This is a Node**

**Relationship**

Each galaxy (or random) point is a **node** which may have <u>relationships</u> to there nodes

For our purposes the relationship information is the distance to neighbours within **$R_{max}$**

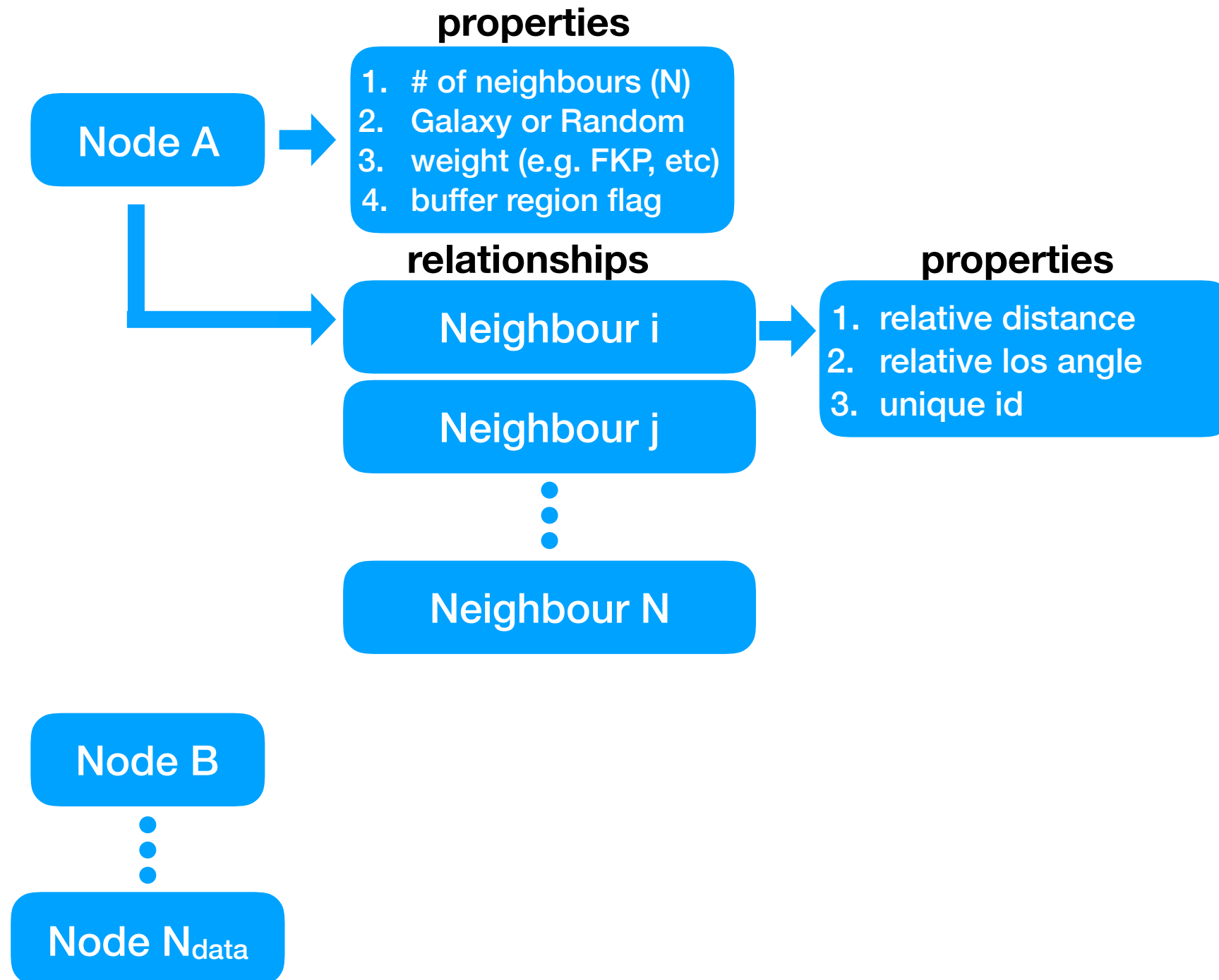Neighbour list can be obtained quickly using a **kd-tree** algorithm

$R_{max}$

A

# Graph Database solution for Galaxy Clustering statistics

## Graph Database Structure



**properties**

**Node A**
1. # of neighbours (N)
2. Galaxy or Random
3. weight (e.g. FKP, etc)
4. buffer region flag

**relationships**

**Neighbour i**

**properties**
1. relative distance
2. relative los angle
3. unique id

**Neighbour j**

**Neighbour N**
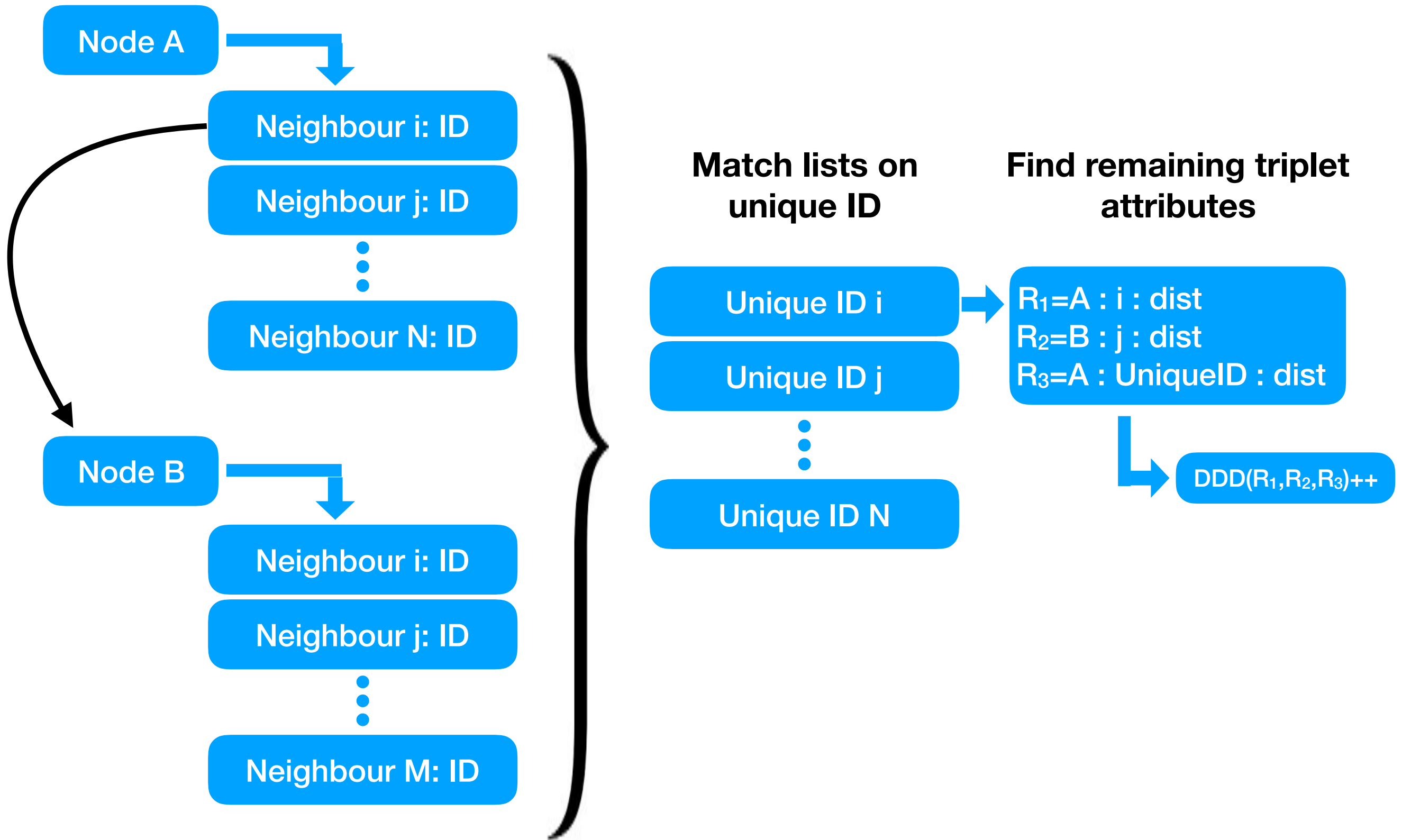
**Node B**

**Node $N_{data}$**

# Graph Database solution for Galaxy Clustering statistics

## Query Graph: 3PCF

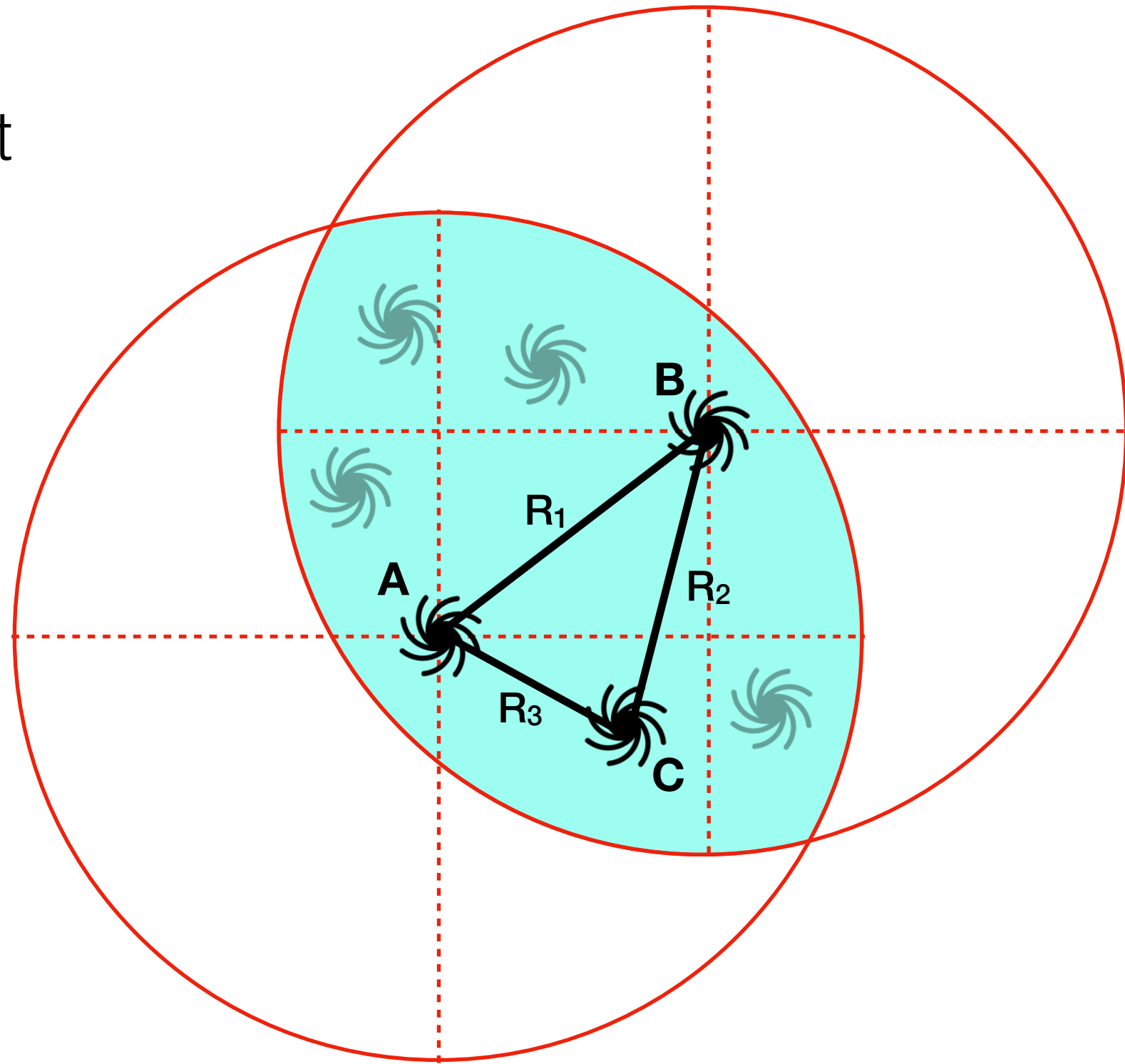# Graph Database solution for Galaxy Clustering statistics

**Query Graph: 3PCF**

The intersection of two sets (**A**,**B**) contains a list of points, **C**, that will complete the triangle A,B,C such that,
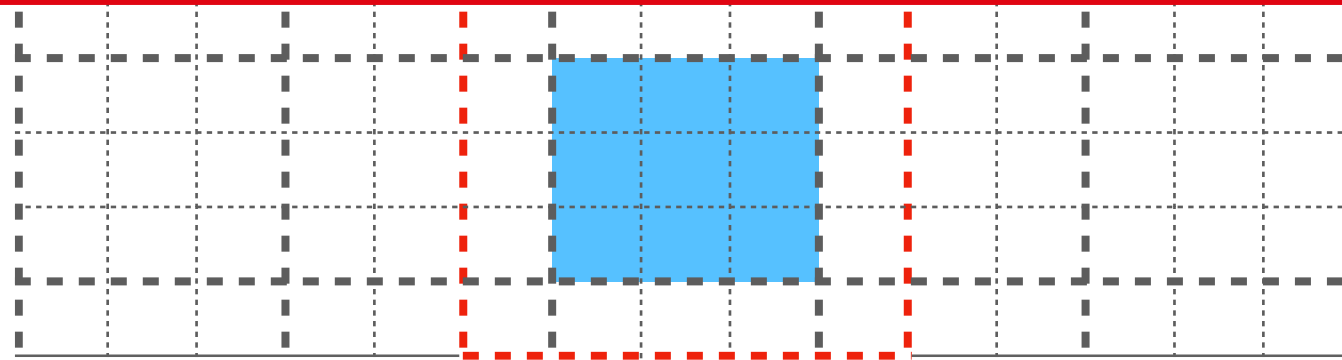**$0 < R_1, R_2, R_3 < R_{max}$**

The intersection must be computed ~ $N^2$ times!

For 2 **ordered lists**, the intersection can be computed very <u>quickly!</u>

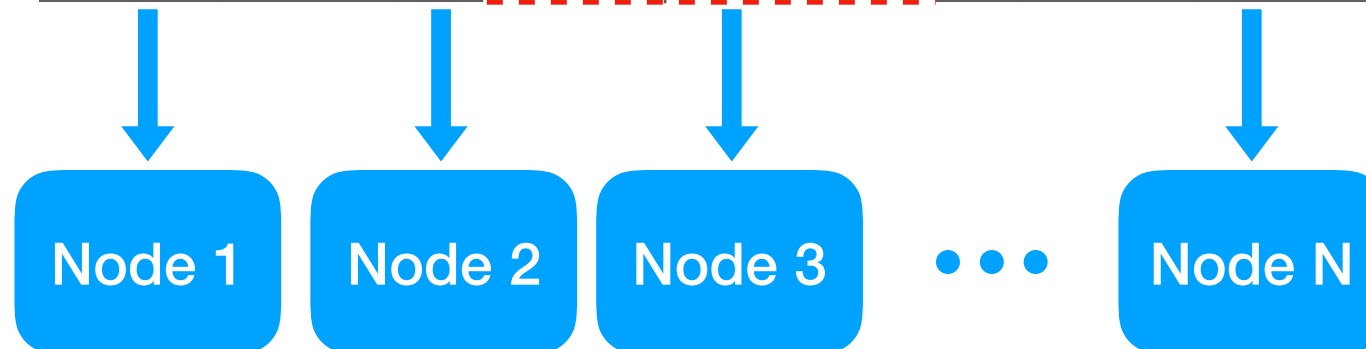# Graph Database solution for Galaxy Clustering statistics
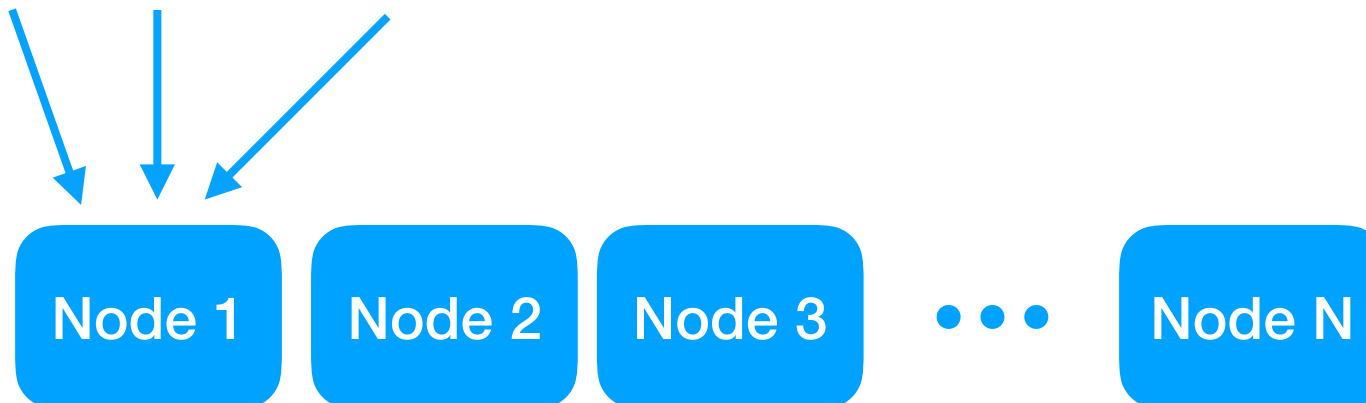
**Domain decomposition is precomputed**
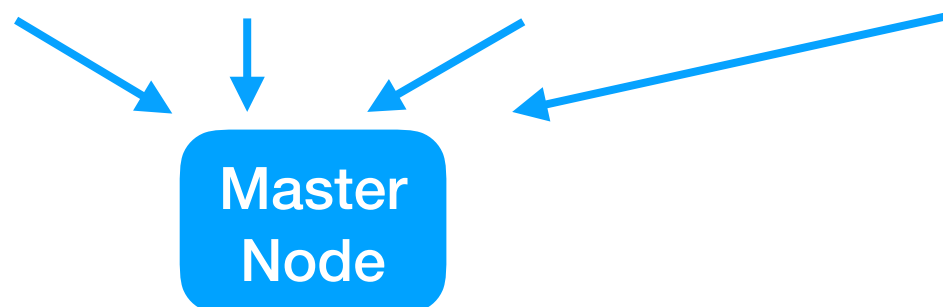
**Spawn *N* MPI processes**

Node 1   Node 2   Node 3  •••  Node N

**OpenMP threading**

CPU 1   CPU 2  •••  CPU N

**Reduce OMP arrays for each node/domain**

Node 1   Node 2   Node 3  •••  Node N

**Collect MPI arrays and output**

Master Node

# Graph Database solution for Galaxy Clustering statistics

## **Benchmarking**



★C. Sabiu, B. Hoyle, J. Kim, X-D Li
★https://arxiv.org/abs/1901.00296

## **BAO in the 3-point function**

Taking the SDSS CMASS DR12 galaxies (~1M)

We look at equilateral triangles

We see evidence of the BAO peak



Isotropic, equilateral 3PCF

★C. Sabiu, B. Hoyle, J. Kim, X-D Li
★https://arxiv.org/abs/1901.00296

# Graph Database solution for Galaxy Clustering statistics

## 4-point correlation function



★C. Sabiu, B. Hoyle, J. Kim, X-D Li
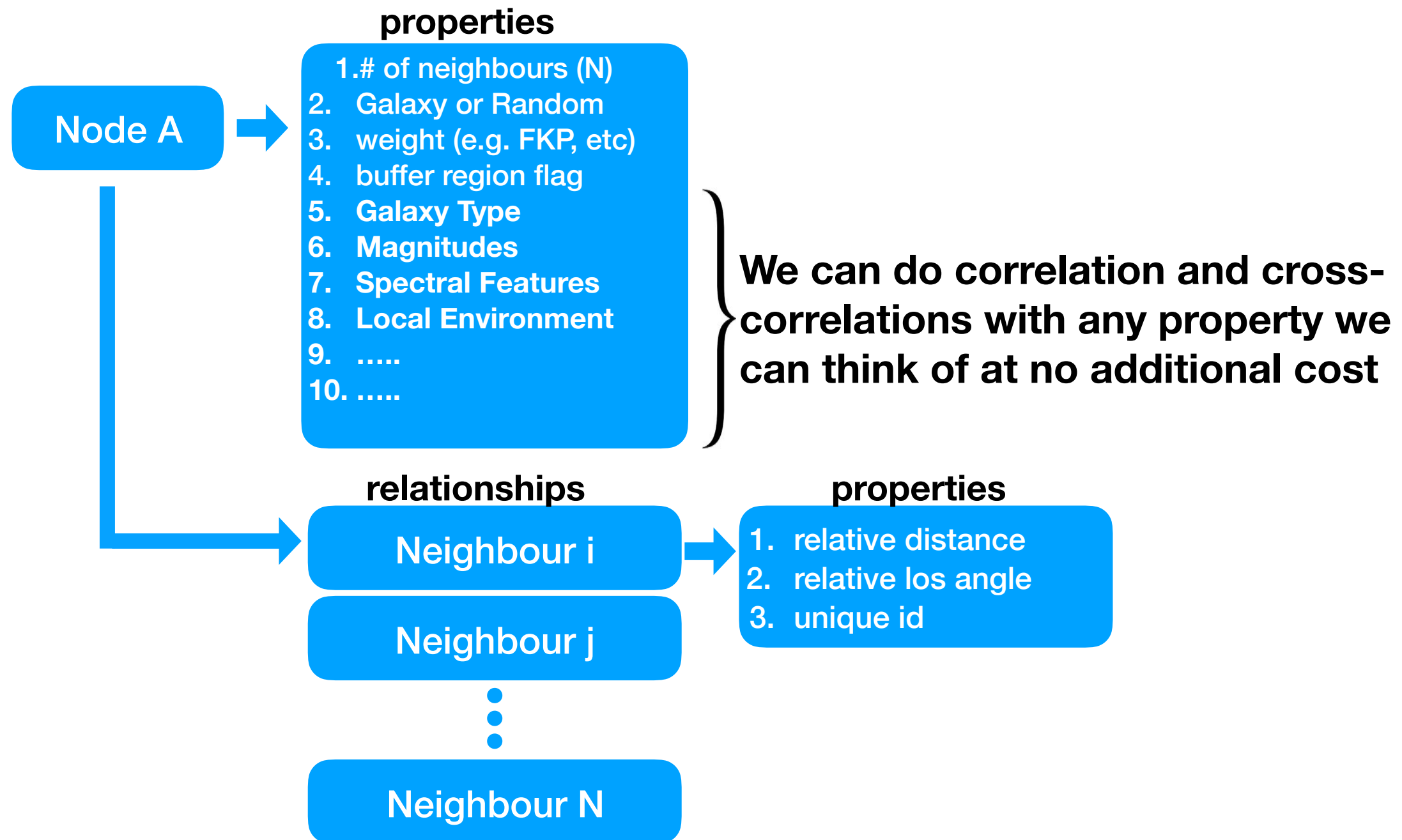
★https://arxiv.org/abs/1901.00296

# Graph Database for exploring new quantities in galaxy evolution, cosmology

## Graph Database Structure

# Graph Database for exploring new quantities in galaxy evolution, cosmology

## Graph Database Structure



**properties**

**Node A** →
1. # of neighbours (N)
2. Galaxy or Random
3. weight (e.g. FKP, etc)
4. buffer region flag
5. **Galaxy Type**
6. **Magnitudes**
7. **Spectral Features**
8. **Local Environment**
9. .....
10. .....

**We can do correlation and cross-correlations with any property we can think of at no additional cost**

**relationships**

**properties**

**Neighbour i** →
1. relative distance
2. relative los angle
3. unique id

**Neighbour j**

**Neighbour N**

theoretical $\xi_l(r)$ with b1=1.82,b2=0.22 and b3=26

measured $\xi_n$ and $\xi_l$

measured $\xi_n$ and $\xi_l$

measured $\xi_l(r) - \xi_n(r)$

Luminosity

Number

BAO

# Graph Database for exploring new quantities in galaxy evolution, cosmology

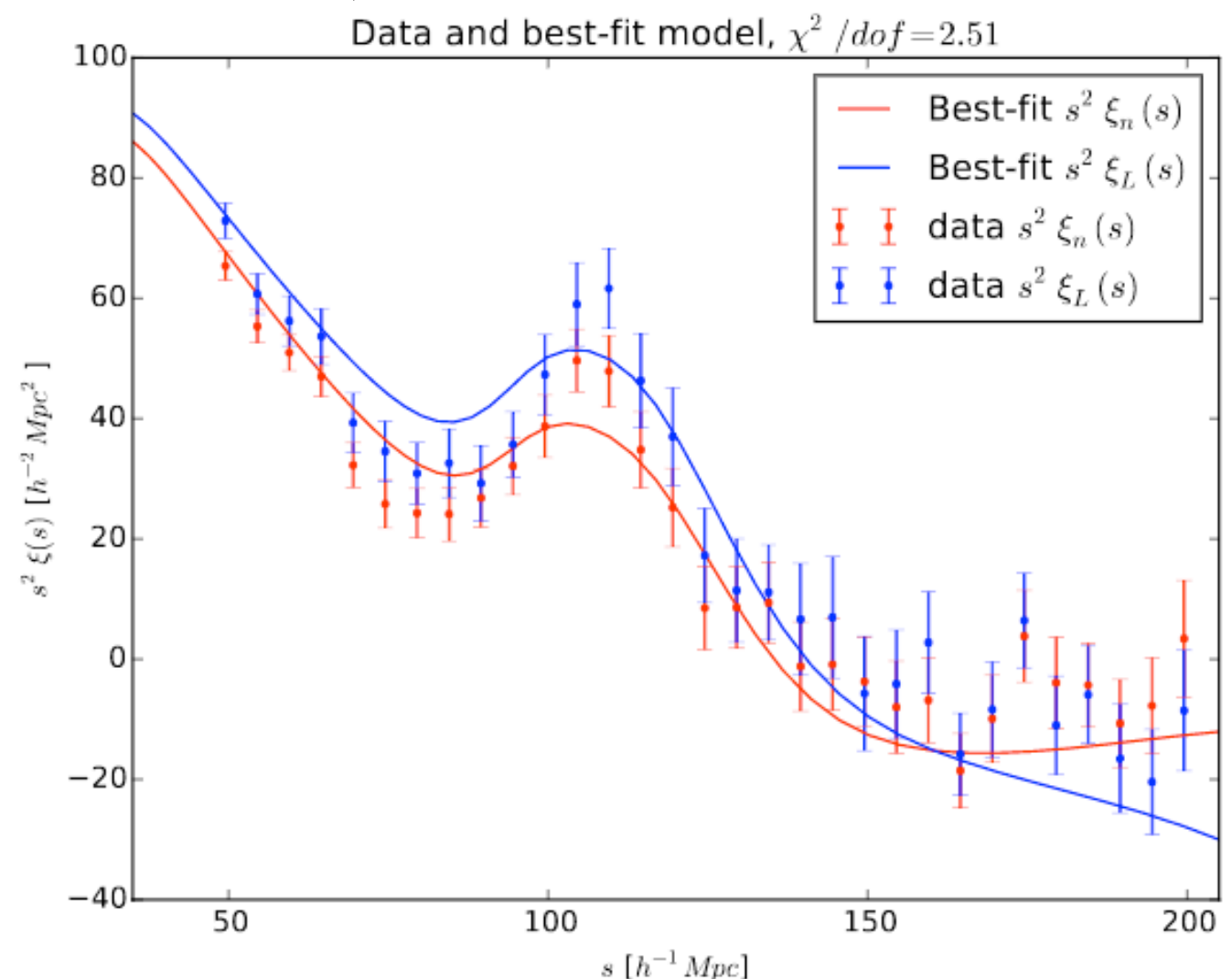**Following the theoretical work of Barkana & Loeb 2010**

**We develop a model for the luminosity weighted correlation function of galaxies:**

$$\xi_L = B_{L,t}^2 \xi_{\text{tot}} + 2B_{L,t}B_{L,\Delta}\xi_{\text{add}} + B_{L,\Delta}^2 B_{\text{CIP}}\hat{\xi}_{\text{CIP}},$$

**This equation has dependence on:**
- **A linear bias with dark matter**
- **large scale clustering of baryons, potentially a new quantity to consider in galaxy evolution**
- **Compensated Isocurvature Perturbations (CIP) between baryons and dark matter in the early universe**

**Looking at spatial cross-correlations with different quantities unlocks new physical interpretation of the data**



Data and best-fit model, $\chi^2/dof = 2.51$

Legend:
- Best-fit $s^2 \xi_n(s)$
- Best-fit $s^2 \xi_L(s)$
- data $s^2 \xi_n(s)$
- data $s^2 \xi_L(s)$

y-axis: $s^2 \xi(s) \ [h^{-2}\ Mpc^2]$
x-axis: $s \ [h^{-1}\ Mpc]$

M. T. Soumagnac, R. Barkana, C. G. Sabiu, A. Loeb, et.al **PRL 2016**
M. T. Soumagnac, C. G. Sabiu, R. Barkana, and J. Yoo    **MNRAS 2019**

# Conclusions



**GRAMSCI**

**v1 out soon!**

- **We introduce a new clustering algorithm, soon publicly available under a GNU public release licence**

- **GRAph Made Statistics for Cosmological Information: GRAMSCI available soon from: http://bitbucket.org/csabiu/gramsci**

- **GRAMSCI performs much better than purely tree based approaches**

- **We show the performance by measuring all possible 3pCF unto and beyond the BAO scale with current SDSS BOSS data**

- **We make the first measurements of the 4-point function of SDSS galaxies**

- **We show the flexibility of adopting a Big Data Analytic approach. As an example the luminosity-number density cross correlation has the potential to unlock new information in the galaxy data that depends on baryonic physics and compensated isocurvature originating in the very early Universe.**

## Thank You 감사합니다