



**HAL**  
open science

## The spatial evolution of young massive clusters

Anne Buckner, Zeinab Khorrami, Pouria Khalaj, Stuart Lumsden, Isabelle Joncour, Estelle Moraux, Paul Clark, René Oudmaijer, José Manuel Blanco, Ignacio de La Calle, et al.

► **To cite this version:**

Anne Buckner, Zeinab Khorrami, Pouria Khalaj, Stuart Lumsden, Isabelle Joncour, et al.. The spatial evolution of young massive clusters. *Astronomy and Astrophysics - A&A*, 2019, 622, pp.A184. 10.1051/0004-6361/201832936 . hal-02416588

**HAL Id: hal-02416588**

**<https://hal.science/hal-02416588>**

Submitted on 1 Sep 2021

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

# The spatial evolution of young massive clusters

## I. A new tool to quantitatively trace stellar clustering

Anne S. M. Buckner<sup>1</sup>, Zeinab Khorrami<sup>2</sup>, Pouria Khalaj<sup>3</sup>, Stuart L. Lumsden<sup>1</sup>, Isabelle Joncour<sup>3,5</sup>, Estelle Moraux<sup>3</sup>, Paul Clark<sup>2</sup>, René D. Oudmaijer<sup>1</sup>, José Manuel Blanco<sup>4</sup>, Ignacio de la Calle<sup>4</sup>, José M. Herrera-Fernandez<sup>4</sup>, Frédérique Motte<sup>3</sup>, Jesús J. Salgado<sup>4</sup>, and Luis Valero-Martín<sup>4</sup>

<sup>1</sup> School of Physics and Astronomy, University of Leeds, Leeds LS2 9JT, UK  
e-mail: a.s.m.buckner@leeds.ac.uk

<sup>2</sup> School of Physics and Astronomy, Cardiff University, The Parade CF24 3AA, UK

<sup>3</sup> Université Grenoble Alpes, CNRS, IPAG, 38000 Grenoble, France

<sup>4</sup> Quasar Science Resources, S.L., Edificio Ceudas, Ctra. de La Coruña, Km 22.300, 28232, Las Rozas de Madrid, Madrid, Spain

<sup>5</sup> Department of Astronomy, University of Maryland, College Park, MD 20742, USA

Received 2 March 2018 / Accepted 8 January 2019

### ABSTRACT

**Context.** There are a number of methods that identify stellar sub-structure in star forming regions, but these do not quantify the degree of association of individual stars – something which is required if we are to better understand the mechanisms and physical processes that dictate structure.

**Aims.** We present the new novel statistical clustering tool “INDICATE” which assesses and quantifies the degree of spatial clustering of each object in a dataset, discuss its applications as a tracer of morphological stellar features in star forming regions, and to look for these features in the Carina Nebula (NGC 3372).

**Methods.** We employ a nearest neighbour approach to quantitatively compare the spatial distribution in the local neighbourhood of an object with that expected in an evenly spaced uniform (i.e. definitively non-clustered) field. Each object is assigned a clustering index (“*I*”) value, which is a quantitative measure of its clustering tendency. We have calibrated our tool against random distributions to aid interpretation and identification of significant *I* values.

**Results.** Using INDICATE we successfully recover known stellar structure of the Carina Nebula, including the young Trumpler 14-16, Treasure Chest and Bochum 11 clusters. Four sub-clusters contain no, or very few, stars with a degree of association above random which suggests these sub-clusters may be fluctuations in the field rather than real clusters. In addition we find: (1) Stars in the NW and SE regions have significantly different clustering tendencies, which is reflective of differences in the apparent star formation activity in these regions. Further study is required to ascertain the physical origin of the difference; (2) The different clustering properties between the NW and SE regions are also seen for OB stars and are even more pronounced; (3) There are no signatures of classical mass segregation present in the SE region – massive stars here are not spatially concentrated together above random; (4) Stellar concentrations are more frequent around massive stars than typical for the general population, particularly in the Tr14 cluster; (5) There is a relation between the concentration of OB stars and the concentration of (lower mass) stars around OB stars in the centrally concentrated Tr14 and Tr15, but no such relation exists in Tr16. We conclude this is due to the highly sub-structured nature of Tr16.

**Conclusions.** INDICATE is a powerful new tool employing a novel approach to quantify the clustering tendencies of individual objects in a dataset within a user-defined parameter space. As such it can be used in a wide array of data analysis applications. In this paper we have discussed and demonstrated its application to trace morphological features of young massive clusters.

**Key words.** methods: statistical – stars: statistics – open clusters and associations: general – stars: general – stars: massive – ISM: individual objects: NGC 3372

## 1. Introduction

Massive stars are fundamental to the evolution of galaxies, profoundly impacting the interstellar medium through chemical enrichment (outflows, supernovae), mixing and turbulence (winds, outflows, supernovae), and heating/cooling (ionising radiation).

Unfortunately while isolated low mass star formation appears to be well described observationally (e.g. Shu et al. 1987, Andre et al. 2000, Luhman 2012), there is still little consensus about the formation of massive stars. This is largely due to observational reasons. High mass stars are rare, evolve rapidly, and have shorter lifetimes than low mass stars. They also emerge onto the main sequence still heavily embedded, having

formed almost exclusively in associations, groups and clusters (de Wit et al. 2005). The linked formation and evolution of both massive stars and clusters, and how they interact is clearly part of the “picture” for massive star formation, but much is still unknown (e.g. Zinnecker & Yorke 2007).

To discriminate between different models for cluster and/or massive star formation/evolution requires a multi-pronged analysis of the structure and dynamics of the stars and gas in young massive clusters (YMCs). To this end we created the StarFormMapper<sup>1</sup> (SFM) project.

One of the fundamental analytical techniques required is to study how the stars and gas “cluster” together. Here, we are

<sup>1</sup> <https://starformmapper.org/>

particularly interested in the study of the intensity, correlation and spatial distribution of point processes, which collectively help to define the distribution and clustering of those points (see, e.g. Møller & Waagepetersen 2007). We are not concerned in this paper with searching for stellar “sub-structure” (discrete star groupings), but rather for suitable statistical measures of the distribution of these point patterns. This is complicated in star formation regions as the distributions of stars and gas are inherently heterogeneous. Many of the best understood statistics from other fields are therefore not easily applied (or are simply invalid). In addition, we wish to use techniques which are valid in any number of dimensions and are applicable easily at different distances, whilst still being computationally simple.

Several global methods have been used in the past. The 2-point correlation function is well studied in cosmology, but has also been used in star forming regions (e.g. Gomez et al. 1993, Scalo & Chappell 1999). The Q parameter by Cartwright & Whitworth (2004) uses a very different technique, which compares the average length from the Minimum Spanning Tree (MST) with the average length from the complete graph of all points, and can distinguish between a smooth overall radial density gradient and multi-scale fractal sub-clustering in a region. It has successfully identified signatures of sub-structure in the Cygnus OB2 (Wright et al. 2014), Serpens, Ophiucus and Perseus star forming regions (Schmeja et al. 2008) and has been applied to assess the dynamical status of star clusters in numerical simulations (e.g. Parker et al. 2014). Similar methods have also been applied to the study of mass segregation (see e.g. Parker & Goodwin 2015). However these methods still suffer if heterogeneous structures are present (e.g. Cartwright & Whitworth 2009). In particular if we wish to compare observations and simulations, great care must be taken in such circumstances.

An alternative possible approach is suggested by the field of geostatistics where interest has also focused on the use of local indicators (e.g. Anselin 1995). In this case, rather than calculating a single parameter for a group of stars/gas as a whole, every unique point has its own derived value. These can then be used to characterise the distribution. That is the approach we will follow here.

In this paper we present our new statistical clustering tool INDICATE (INdex to Define Inherent Clustering And Tendencies) that we are currently implementing in the SFM project. Hopkins & Skellam (1954) established the Hopkins statistic to assess the global clustering tendency of a dataset by testing its spatial randomness through quantitative measurements of its uniformity. A single global value is calculated through a comparison of the mean k-nearest neighbour distances between objects within the dataset, and between points in the dataset and a similarly constructed uniform random sample. We propose instead to derive a similar index but for *every* point in the dataset individually.

This paper is structured as follows. Sect. 2 describes how our tool works and is calibrated. In Sect. 3 we use INDICATE to trace stellar morphological features of star formation, demonstrating its ability to cope with the complex, often poorly defined, spatial clusterings expected in young massive star forming regions/clusters. Our conclusions are presented in Sect. 4.

## 2. INDICATE: INdex to Define Inherent Clustering And Tendencies

### 2.1. General tool description

INDICATE is a tool to quantify the degree of association of each point in a 2+D discrete dataset. It requires no a priori knowledge

of – nor makes assumptions about – the sub-structure present in a dataset, since it is a *local* statistic. The separation of the spatial position of the  $j$ th point in the actual dataset with the  $N$ th nearest neighbour in an evenly spaced uniform (i.e. definitively non-clustered) control distribution is determined. The mean value of this separation,  $\bar{r}$ , is then derived. Finally, INDICATE assigns an index,  $I$ , to every point in the actual dataset, which is simply the ratio of the actual number of real neighbours within this mean separation,  $\bar{r}$ , and  $N$ , the nearest neighbour number. Since our tool fundamentally relies on properties linked to distance ratios in the real and uniform sample, it is itself region distance independent in principle at least (but see Appendix C). Because our tool uses an evenly spaced grid as the comparison it is computationally less intensive than a direct implementation of the Hopkins, or similar, statistics. Below we describe step-by-step how this index is derived for a simple 2D distribution. A future paper will deal with the implementation to a full 3D dataset using Gaia parallaxes.

#### 2.1.1. Step I: Define the bounds of the dataset

INDICATE is designed to be applicable to any desired N-dimensional parameter space. However, as outlined here, it assumes that all dimensions have the same scaling (e.g. J2000 sky coordinates should be converted to a local coordinate frame prior to beginning).

The bounds of the dataset parameter space are defined from the density distribution, and the area occupied by the data,  $A$ , measured. The shape of the delimited area has a negligible effect on the tool (Appendix B), but for clarity we will use a rectangular parameter space for all our explanatory datasets in the description of the tool. In practice, the shape is related to the problem to be studied, and is defined by the user.

The number density,  $n_{\text{obs}}$ , of the dataset is determined using,

$$n_{\text{obs}} = \frac{N_{\text{tot}}}{A}, \quad (1)$$

where  $N_{\text{tot}}$  is the total number of points in the dataset.

#### 2.1.2. Step II: Generate the control distribution

An adaptive evenly spaced uniform point distribution which we designate the *control distribution* is generated. The control distribution is rectangular (regardless of dataset shape defined in Step I), populates the bounded parameter space and has the same number density as the dataset i.e.

$$n_{\text{con}} = n_{\text{obs}} \quad (2)$$

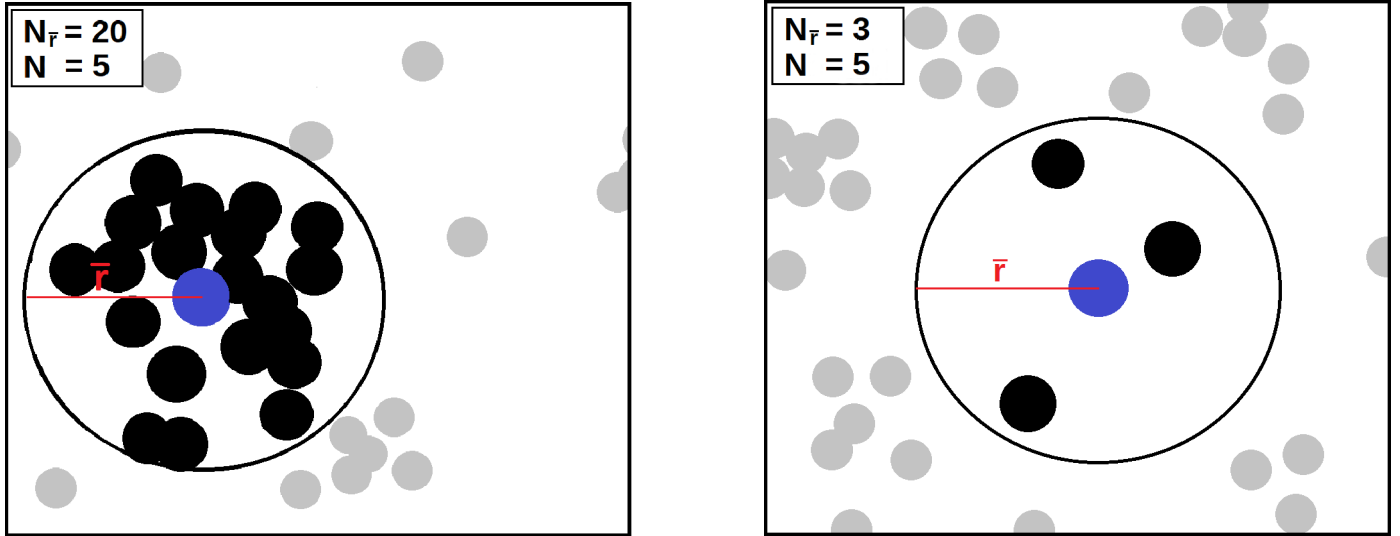
#### 2.1.3. Step III: Measure the mean nearest neighbour distance

The mean Euclidean distance,  $\bar{r}$ , of each point,  $j$ , in the dataset to its  $N$ th nearest neighbour in the control distribution is measured using:

$$\bar{r} = \frac{\sum_{j=1}^{N_{\text{tot}}} r_j}{N_{\text{tot}}}, \quad (3)$$

where,

$$r_j = \sqrt{(x_j - x_N^{\text{con}})^2 + (y_j - y_N^{\text{con}})^2}, \quad (4)$$



**Fig. 1.** Demonstration of how INDICATE defines the index  $I_{j,N}$  for a point. All points within a radius of  $\bar{r}$  of the selected point (marked in blue) are counted ( $N_{\bar{r}}$ ) and compared to the number of points expected within the same radius in an evenly spaced uniform point distribution with the same number density as the points parent sample ( $N$ ). The index of the blue point is calculated using Eq. (5) as  $I_5 = 4.0$  (left panel) and  $I_5 = 0.6$  (right panel).

and  $(x_j, y_j)$  are the respective  $x$  and  $y$  axis coordinates of point  $j$  and  $(x_N^{\text{con}}, y_N^{\text{con}})$  are the respective  $x$  and  $y$  axis coordinates of the  $N$ th nearest neighbour in the control distribution.

#### 2.1.4. Step IV: Calculate the index, $I$

The number of points,  $N_{\bar{r}}$ , closer than  $\bar{r}$  to each point,  $j$ , in the dataset is counted (see Fig. 1). The index of point  $j$  is then the ratio of the number of neighbours closer than  $\bar{r}$  in the dataset with that expected by a non-clustered distribution i.e.

$$I_{j,N} = \frac{N_{\bar{r}}}{N}, \quad (5)$$

where  $N$  is the  $N$ th nearest neighbour number (e.g. if  $\bar{r}$  is measured for the 5th nearest neighbour,  $N = 5$ ; 6th nearest neighbour,  $N = 6$ ...etc.). The ratio  $I_{j,N}$  is unitless and has a range of  $0 \leq I_{j,N} \leq \frac{N_{\text{tot}}-1}{N}$ , such that the higher its value the more spatially clustered point  $j$  is.

It is important to note the index is not a measure of local surface density – Eq. (5) describes the *local spatial distribution* of point  $j$ . Therefore although the index is proportional to the local point surface density of a dataset, it is possible for two datasets with significantly different densities to have identical index values if their points have the same spatial distribution. For example, Fig. 2A shows index values derived for Gaussian cluster with 100 members, using a nearest neighbour number of  $N = 5$ . On visual inspection the highest values have been assigned to stars with the highest degree of association. Figures 2B and 2C show the same cluster, but with an observed angular dispersion and surface density the cluster would have if its distance was a factor of 4 and 16 times larger than A respectively. Applying INDICATE under the same conditions as for Fig. 2A, the index values for each star remain unchanged in Fig. 2B and Fig. 2C i.e. for all members  $\Delta I_5 \equiv 0$  despite the surface density increasing by a factor of 256 between Fig. 2A and Fig. 2C, as the local spatial distribution of members in all three clusters is identical.

Hence the index can be used to directly analyse variation in the spatial distribution of points, in any desired parameter space,

**Table 1.** Statistics of index values derived for clusters' D, E, F (Fig. 3).

Cluster	Mo[ $I_5$ ]	max( $I_5$ )	% stars $I_5 > I_{\text{max}}$
D	8.2	10.6	81.2
E	1.6	4.8	17.5
F	3.0	6.4	34.1

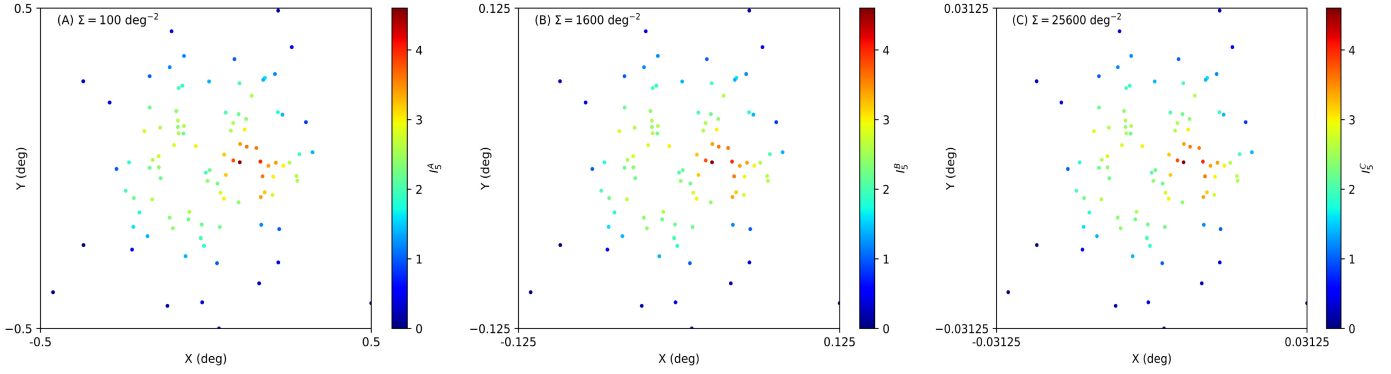
**Notes.** The percentage of members clustered above random ( $I_5 > I_{\text{max}}$ ) is given for each cluster.

(a) within a dataset and/or (b) comparatively between two or more datasets.

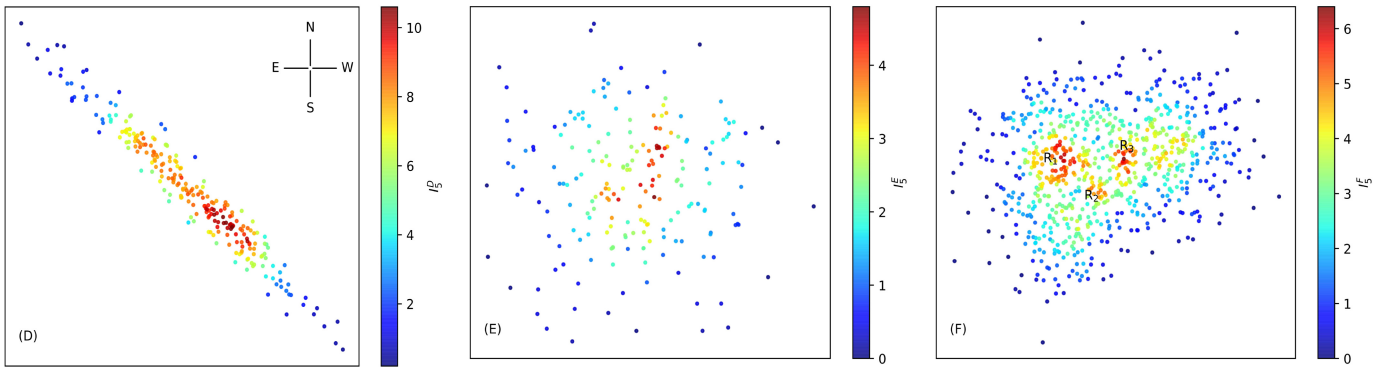
We conduct a series of statistical tests on randomly generated samples to calibrate, and investigate edge and field effects on, our index. The methodology and results of these tests are described in detail in Appendices A–C respectively. In brief:

1. There is a logarithmic relationship between the maximum index value for a random distribution and sample size.
2. The index is independent of a samples number density.
3. There is a relationship between the typical index value for a random distribution and the chosen  $N$ th nearest neighbour number.
4. The size of the control distribution is essentially arbitrary (but care should be taken when a point in the sample has an index value which is on the boundary of a chosen significance threshold).
5. Uniformly distributed interloping field stars (e.g. in observational datasets) typically do not significantly affect the index values of true cluster members.
6. If interloping field stars are distributed in a gradient, the index derived for true cluster members is independent of gradient shape for small nearest neighbour numbers ( $N = 3$ ).

We note that in samples which contain interloping field stars, that the field stars are also assigned index values by INDICATE, so care must be taken when interpreting the values and drawing conclusions on the physical origins of the clustering tendencies of stars in these samples.



**Fig. 2.** Panels A–C: index values derived by INDICATE for a synthetic Gaussian cluster with 100 members, using a nearest neighbour number of  $N = 5$  in all three instances. In B and C the angular dispersion of the cluster was reduced to simulate how it would be observed if its distance was a factor of 4 and 16 times larger than A, respectively. Despite the significant increase in point surface densities, the index values derived for each star is unchanged because there is no change in the relative spatial distribution of members.



**Fig. 3.** Index values derived by INDICATE for members of synthetic clusters D, E, F using a nearest neighbour number of  $N = 5$  and a standard control distribution (CDA – see Appendix B).

## 2.2. Implementation example

In Sect. 3 we apply our tool on a real stellar catalogue of NGC 3372. Here we demonstrate using synthetic datasets INDICATE’s ability to quantify the degree of association for each point in a 2D discrete dataset and suitability as a statistical measure for comparative analysis of the spatial distributions of points in multiple datasets.

Figure 3 shows three clusters (D, E and F) with different degrees of elongation, angular dispersion, sub-structure, surface density and number of members. We apply INDICATE to each using a  $N$ th nearest neighbour number of  $N = 5$  and a standard control distribution (CDA – see Appendix B). Table 1 shows their mode and maximum index values, and the percentage of members with  $I_5 > I_{\max}$  (Eq. (A.2)).

Cluster D is the most elongated, and its members have been identified by INDICATE as having the strongest clustering tendencies (highest index values) of the three clusters. We find that the greatest degree of association is within its central region (up to 53 neighbours within  $\bar{r}$ ) and spatial clustering of members is asymmetrical – stars to the NE of the highest index members have significantly higher index values than those to the SW.

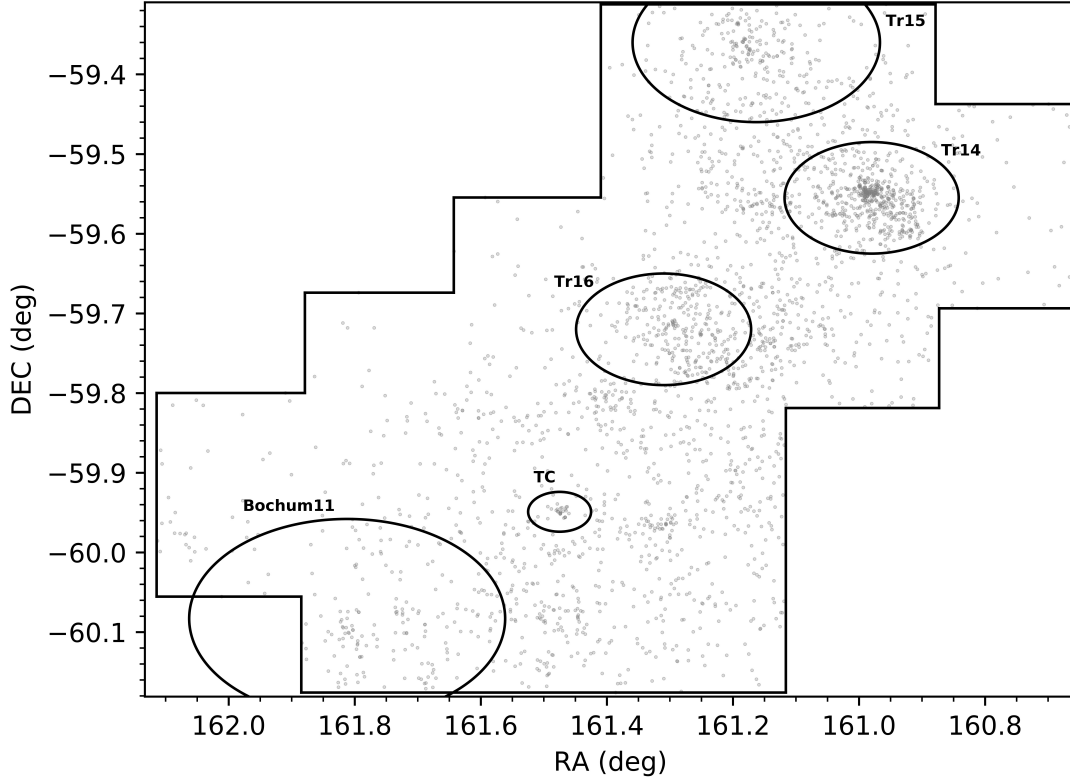
Members of cluster E have a markedly different spatial distribution and clustering tendencies to those of cluster D. The spatial distribution of members lacks a strong radial correlation and two concentrations of (relatively) high index stars are identified. Stars with the greatest degree of association have 24 neighbours within  $\bar{r}$  (a factor of 2 less than D), typically the degree of association of members of E is a factor of 6.6 less than members of cluster D

(i.e. members of E are significantly less tightly clustered than those of D). Of course, as the clusters are being analysed in a 2D parameter space it is conceivable that stars in cluster D and E may have very similar spatial distributions but are being viewed from different rotations around the 3D axis. We therefore advise caution when drawing conclusions from comparisons of the INDICATE values of two or more observational 2D datasets alone.

Cluster F has three concentrations of high index members (marked on Fig. 3 as  $R_1$ ,  $R_2$ ,  $R_3$ ). Stars that form part of  $R_1$  and  $R_3$  have similar index values, that is they have similar clustering tendencies, but those of  $R_2$  are less spatially clustered. If this were a real dataset, where clustering behaviours are dictated by underlying physics, these index values could be used as a starting point to explore the physical causes of the identified discrepancies between the clustering behaviours of the three concentrations (e.g. differences in evolutionary stage, initial conditions, stellar mass of members) and also the identified disparity of members clustering behaviours between clusters D, E and F. This form of quantitative analysis and comparison of spatial behaviours, as achieved here by INDICATE, is not possible with the discussed global methods and/or established clustering algorithms.

## 3. Tracing morphological features

The Carina Nebula (NGC 3372) is a massive star forming HII complex in the southern sky at a distance of 2.3 kpc, containing  $>10^5 M_{\odot}$  of gas+dust (Preibisch et al. 2011a). It is one



**Fig. 4.** Stellar catalogue by Kuhn et al. (2014), with stars represented as grey dots. The positions of the Trumpler 14 (Tr14), Trumpler 15 (Tr15), Trumpler 16 (Tr16), Treasure Chest (TC) and Bochum 11 clusters are overlaid as black ellipses (see text for details).

of the nearest and richest concentrations of OB stars ( $>130$ ; Kuhn et al. 2014) in the Galaxy and includes some of the most massive and luminous known single and binary stars (e.g. Eta Carinae, HD 93129, W25). The region is well studied and has considerable sub-structure including the Trumpler 14-16, Collinder 228, Collinder 232, Collinder 234, Bochum 10 and Bochum 11 clusters. Triggered star formation is ongoing in the complex, driven by massive star feedback (Smith et al. 2008, 2010; Preibisch et al. 2011b; Gaczkowski et al. 2013). Coupled with its low line of sight extinction, the Carina Nebula is therefore an ideal laboratory in which to study massive star formation. For a review of the region see Smith et al. (2008).

### 3.1. Identification of stellar structure

We apply INDICATE to the stellar catalogue of 2790 stars for the region as described by Kuhn et al. (2014), and plotted in Fig. 4, with a nearest neighbour number of  $N = 5$  and an extended control distribution (CDB – see Appendix B). This catalogue was selected because it covers an inner region of the Carina Complex ( $\sim 0.38^\circ$ ) which is rich in sub-structure, containing at least 20 sub-clusters (detected by the original authors) and includes the young Trumpler 14 (Tr14), Trumpler 15 (Tr15), Trumpler 16 (Tr16), Treasure Chest (TC) and Bochum 11 clusters. The position and radius of the TC cluster in Fig. 4 is as given by Dutra & Bica (2001), and the position and radius of Tr14-16 and Bochum 11 clusters by the MWSC catalogue (Kharchenko et al. 2013). On visual inspection, stars in the North West (NW) region of the catalogue have significantly higher degrees of association than those in the South East (SE) region.

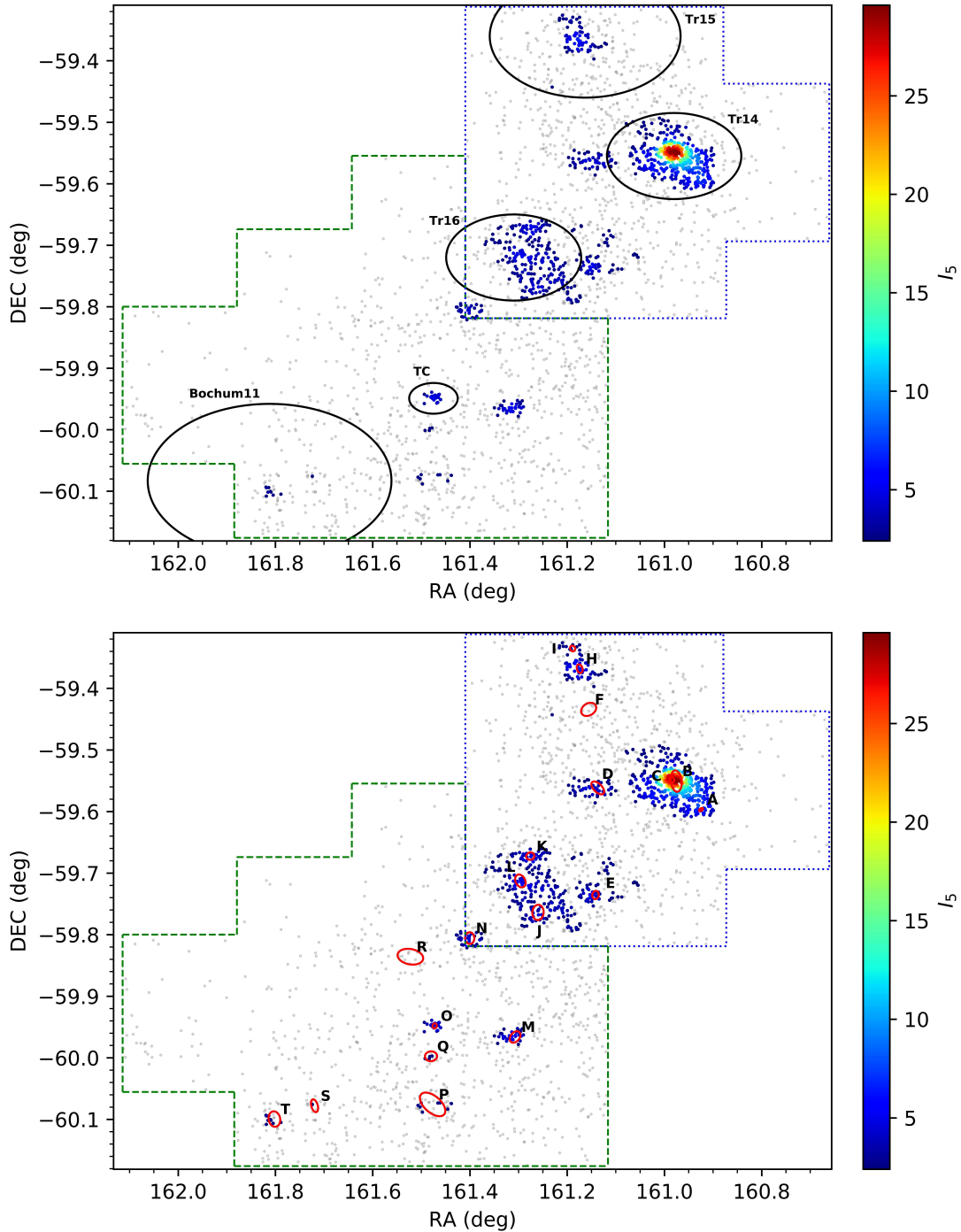
The top plot in Fig. 5 and left plots of Fig. 7, show the distribution of stars with their index values and the boundaries

**Table 2.** Statistics for stars within the radial boundaries of the Tr14, Tr15, Tr16, TC and Bochum 11 clusters with an index value of  $I_5 > I_{\text{sig}}$ .

Cluster	Total	$\bar{I}_5$	max $I_5$
Tr14	470 (85.2%)	12.4	29.6
Tr15	75 (29.1%)	3.9	6.6
Tr16	212 (73.4%)	3.6	5.6
TC	25 (83.3%)	4.3	4.8
Bochum 11	10 (5.8%)	2.7	3.0

of the Tr 14-16, TC, Bochum 11 clusters overlaid. We define a significance threshold – that is the value of  $I_5$  above which a star is significantly clustered above random – of three standard deviations above the mean value expected from a random distribution of the same size evaluated with  $N = 5$  and CDB, such that  $I_{\text{sig}} = \bar{I}_5^{\text{random}} + 3\sigma = 2.3$ . All five clusters are clearly identified by stars within their radial boundaries having an index above the defined significance threshold. This is an expected result as by definition the spatial distribution of cluster stars – particularly at their centres – should display a higher degree of clustering than a random (and background) field.

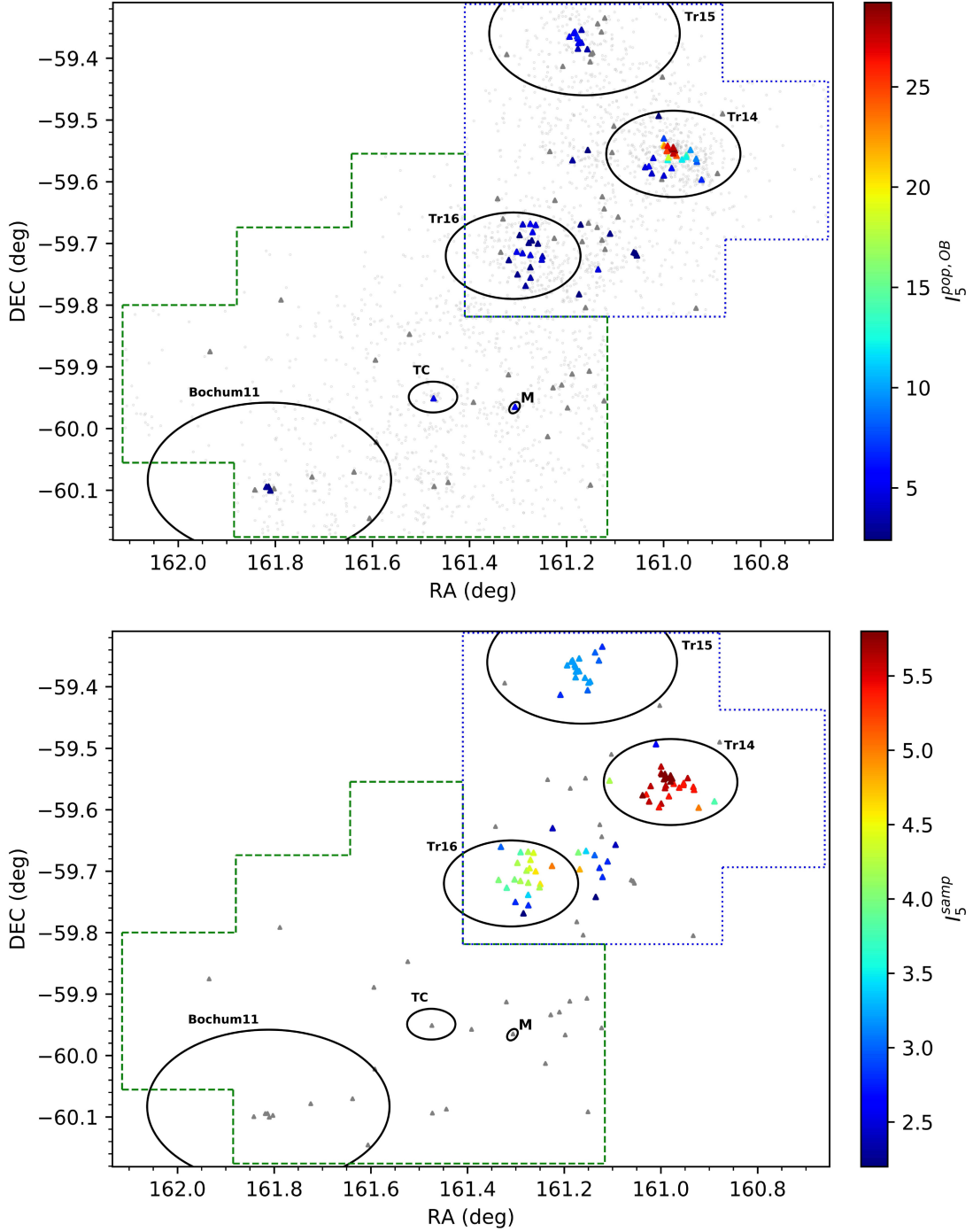
Table 2 gives statistics on the index values derived for each cluster. More than 80% of the stars within the bounds of Tr14 and TC clusters are clustered above random, which is markedly larger than Tr15 (29.1%) and Bochum 11 (5.8%). Trumpler 16 has a comparable proportion of stars with index values above the threshold to that of the Tr14 and TC clusters (73.4%) but unlike the other four clusters these stars are not centrally concentrated and instead are in less compact concentrations across the cluster region, which is consistent with results of previous studies of



**Fig. 5.** Index values,  $I_5$ , calculated by INDICATE for the Carina region with positions of the Tr14, Tr15, Tr16, TC and Bochum 11 clusters overlaid as black ellipses (*top panel*) and 19 sub-clusters identified by Kuhn et al. (2014) overlaid as red ellipses (*bottom panel*; see text for details). The borders of our designated NW and SE regions are marked with blue dotted and green dashed lines respectively. Stars with an index value above the significance threshold ( $I_5 > 2.3$ ) coloured as described by the colour bar. Grey dots are stars with  $I_5 < 2.3$ .

the cluster's structure (e.g. Wolk et al. 2011). Interestingly, stars clustered above random in Tr15, Tr16 and TC have similar mean index values – that is they have similar degrees of association and clustering tendencies. Stars within the bounds of Tr14 display the highest degree of clustering behaviour with a mean index value a factor of 3 larger than those of the Tr15, Tr16 and TC clusters. In addition, Tr14 also contains the most spatially clustered stars in the Carina region, centrally concentrated at its core, with stars here having up to an additional 137 stars in their local neighbourhoods above that expected in a spatially random distri-

bution. By contrast, stars above the threshold within the bounds of Bochum 11 display the lowest degree of clustering behaviour of the five clusters – having a maximum of just 3 – 4 stars in their local neighbourhoods above random. The high/low proportion of stars within the radial boundary of Tr14/Bochum 11 with an index value above the significance threshold, suggests these clusters are the most and least tightly clustered respectively in the region. In the absence of kinematic data however, we refrain from drawing any conclusions as to the physical origin of this trend.



**Fig. 6.** Index values calculated for the OB population by INDICATE when applied to the entire stellar catalogue,  $I_5^{\text{pop,OB}}$  (top panel); and OB sub-sample,  $I_5^{\text{samp}}$  (bottom panel; see text for details). Positions of the Tr14, Tr15, Tr16, TC, Bochum 11 and M (sub-) clusters are overlaid as black ellipses. The borders of our designated NW and SE regions are marked with blue dotted and green dashed lines respectively. Coloured triangles represent OB stars with an index value above the respective significance thresholds. Grey triangles are OB stars below the respective significance thresholds.

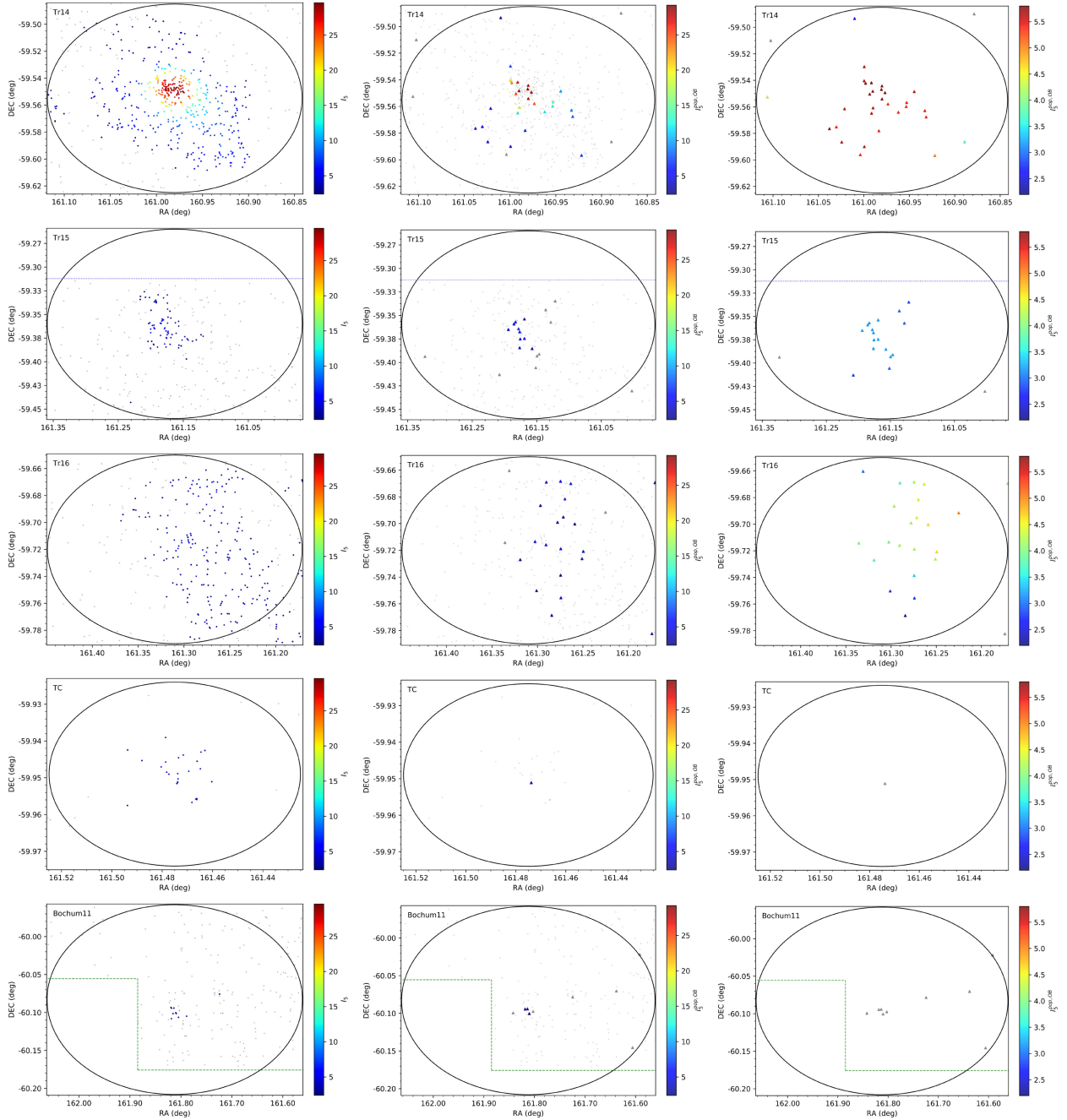
The bottom plot in Fig. 5 shows the distribution of stars with their index values and the positions of the 19 sub-clusters<sup>2</sup> found by Kuhn et al. (2014), as stellar overdensities using finite mixture models, overlaid. Fifteen sub-clusters are clearly identified with a significant number of members that have an index value above the defined threshold. Four sub-clusters (*F*, *P*, *R*,

<sup>2</sup> The 20th sub-cluster, their “G” cluster, is ignored in our analysis due to its large angular extension across the centre of the region.

*S*) do not contain any (or very few) stars spatially clustered above random, suggesting these sub-clusters may not be real clusterings but instead fluctuations in the dispersed population field.

We now look at the clustering tendencies of individual stars across the Carina Nebula. A total of 35.2% of stars in the catalogue are clustered above random. Stars in the NW region typically have higher index values than the SE region, with 49.9% and 9.1% clustered above random respectively. To gauge the





**Fig. 7.** Zoomed-in plots of clusters Tr14, Tr15, Tr16, TC and Bochum11 as shown in (Left panel) top of Fig. 5, (Middle panel) top of Fig. 6 and (Right panel) bottom of Fig. 6.

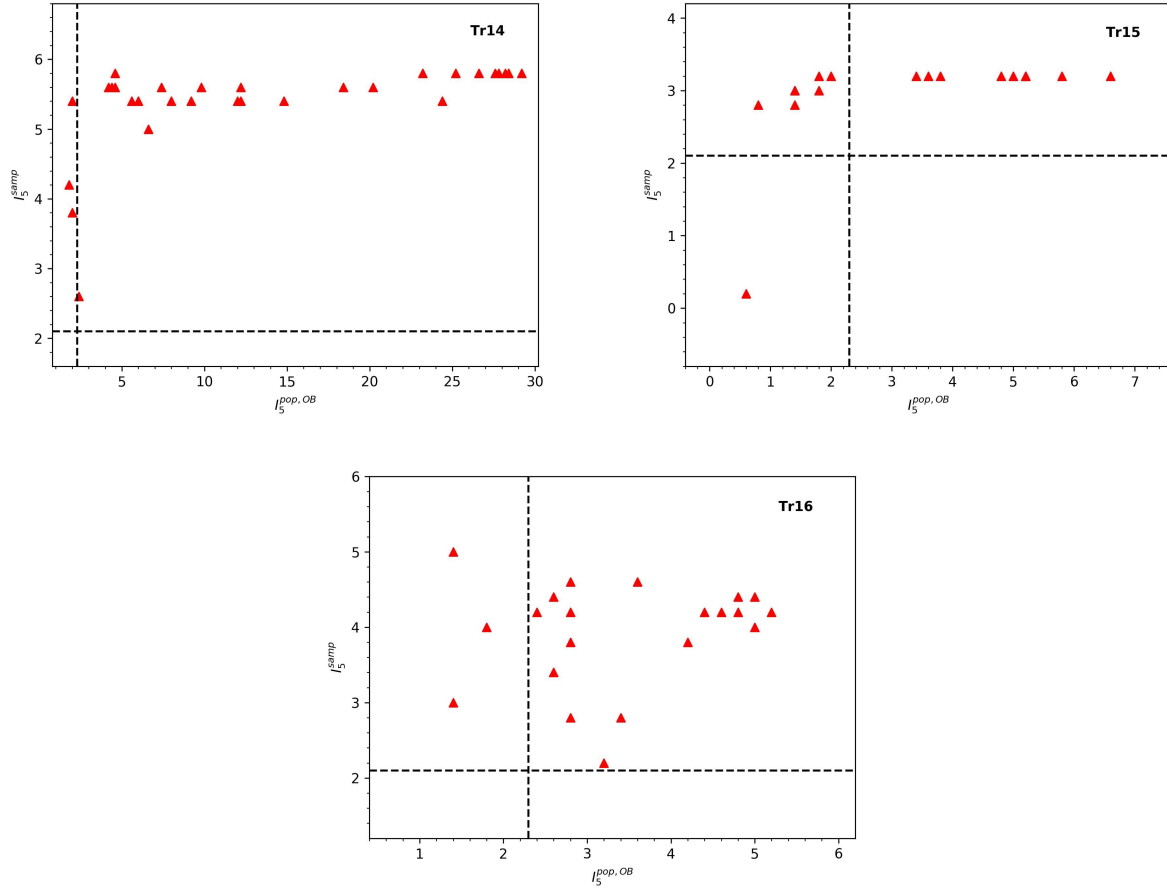
significance of this difference between the NW and SE regions we run a 2 sample K–S test of the index values for all stars in the NW region against those for all stars in the SE region, with a strict significance boundary of  $p < 0.01$ , finding a value of  $p \ll 0.001$  i.e. stars in the NW and SE regions have significantly different clustering tendencies. This result is not entirely unexpected as the NW region is heavily sub-structured, containing 3/5 of the young clusters, and 12/19 of the sub-clusters detected by Kuhn et al. (2014) – whereas the SE region is comparatively sparsely populated and being shaped by radiative winds of the Tr14 and Tr16 clusters (Smith et al. 2008). Therefore the disparity of the SE and NW regions clustering tendencies is reflective of differences in the apparent star formation activity in these regions.

### 3.2. Clustering tendencies of the OB population

We create a sub-sample of the OB stars in the catalogue. To identify OB stars a cross-match search of the catalogue with the SIMBAD<sup>3</sup> database was performed, finding 134 stars listed as either O or B spectral type. Thirteen of these have an ambiguous spectral type or are flagged as being higher order systems so are excluded from the sub-sample, leaving a final selection of 121 stars.

The term “mass segregation” is used interchangeably in the literature to describe two quite different realisations. The first definition (hereafter *Type 1*) refers to a system in which the massive stars are concentrated together at its centre; whereas the second definition (hereafter *Type 2*) refers to a system in which the

<sup>3</sup> <http://simbad.u-strasbg.fr/simbad/>



**Fig. 8.** Index values calculated for the OB population by INDICATE when applied to the entire stellar catalogue,  $I_5^{\text{pop,OB}}$ , and the sub-sample,  $I_5^{\text{samp}}$ . Dashed black lines represent the respective significance thresholds (see text for details). Red Triangles represent OB stars within the radius of Trumpler 14 (*top left*) Trumpler 15 (*top right*) Trumpler 16 (*bottom panel*).

massive stars are in stellar concentrated regions, but are not necessarily concentrated together.

As the index quantifies the degree of association of stars, it by definition identifies (and quantitatively measures) Type 2 mass segregation as values are assigned to stars based on the degree of spatial clustering in their local neighbourhood. Figures 6 and 7 show the index distribution of the OB stars with the positions of the Tr14-16, TC, Bochum 11 and M (sub) clusters overlaid. We find the Tr14, Tr15, Tr16, TC, Bochum 11 and M (sub) clusters have signatures of Type 2 mass segregation. A total of 57.0% of OB stars are clustered above random, which is notably higher than the general populations 35.2% (Sect. 3.1) i.e. cluster concentrations are more frequent around massive stars than typical for stars in this region. Massive stars in the NW region have notably different clustering tendencies to those in the SE region, with 68.1% and 18.5% clustered above random respectively. To gauge the significance of this difference we run a 2 sample K-S test of the index values for all sub-sample stars in the NW region against those for all sub-sample stars in the SE region, with a strict significance boundary of  $p < 0.01$ , finding a value of  $p \ll 0.001$ , which confirms OB stars in the NW and SE regions have significantly different clustering tendencies. These results show signatures of Type 2 mass segregation are present across Carina but are primarily found in the NW region.

It is also possible to use our tool to find signals of the “classical” Type 1 mass segregation. We apply INDICATE to the sub-sample of OB stars with a nearest neighbour number of

$N = 5$  using an extended control distribution (CDB) and define a significance threshold – that is the value of  $I_5$  above which a star is significantly clustered above random – of three standard deviations above the mean value expected from a random distribution of the same size evaluated with  $N = 5$  and CDB, such that  $I_{\text{sig}} = \bar{I}_5^{\text{random}} + 3\sigma = 2.1$ . As the index is a quantitative measure of the degree of clustering of OB stars with other OB stars, it is a local measure of Type 1 mass segregation. We find the massive population is notably more self-clustered than is typical amongst the general stellar population with a total of 64.5% of stars in the sub-sample clustered above random ( $\sim$  a factor of two larger than for the general population). The Tr14, Tr15 and Tr16 clusters have signatures of Type 1 mass segregation with a significant number of OB members more clustered than expected for a random distribution, and mean index values of OB stars within their cluster radii of  $I_5^{\text{samp}} = 5.4, 3.0$  and  $3.9$  respectively. Neither the TC or Bochum 11 clusters have signatures of mass segregation, with mean index values of 0.2 and 0.8 respectively. Massive stars in the NW region have completely different clustering tendencies than the SE region, with 83.0% and 0.0% clustered above random respectively. To gauge the significance of this difference we run a 2 sample K-S test of the index values for all sub-sample stars in the NW region against those for all sub-sample stars in the SE region, with a strict significance boundary of  $p < 0.01$ , finding a value of  $p \ll 0.001$ , which confirms OB stars in the NW and SE regions have significantly different clustering tendencies. These results clearly show that signatures of Type 1 mass segregation are present in the NW

region but not in the SE region – massive stars here are not spatially concentrated together above random.

Finally, we look for correlations in the clustering behaviour of OB stars – is there a relation between the stellar concentrations around massive stars and the self-concentration of the OB population in the Carina region? Figure 8 shows a comparison of the index values derived for the OB population of Tr14, Tr15 and Tr16 from the application of INDICATE to (1) the entire stellar catalogue (Sect. 3.1) and (2) the OB sub-sample (Sect. 3.2). In both Tr14 and Tr15 there is a clear trend between the concentration of OB stars and the concentration of (lower mass) stars around OB stars: while there is a maximum degree of association an OB star can have w.r.t. other OB stars, stellar concentrations around an OB star may continue to increase. We find that Tr16 does not follow this trend, which is consistent with what is known about the structure of the Trumpler clusters. Unlike the Tr14 and Tr15 clusters, Tr16 does not have a strong central concentration but instead is irregularly shaped and heavily sub-structured with multiple sub-clusters (Ascenso et al. 2007, Wang et al. 2011, Wolk et al. 2011). Thus the index values of Tr16 reflect that the OB stars are not clustered together in a single concentration with a (near) constant degree of clustering, but are instead scattered across a region with local concentrations of stars and a variable degree of association.

#### 4. Conclusions

We have developed a powerful novel statistical clustering tool called INDICATE (INDEX to Define Inherent Clustering And Tendencies) to study the intensity, correlation and spatial distribution of point processes in discrete astronomical datasets. The tool assesses the clustering tendency of each object in a dataset and assigns it an index  $I_{j,N}$  (Eq. (5)), using a nearest neighbour approach by comparing the spatial distribution of objects in its local neighbourhood with that expected in an evenly spaced uniform (i.e. definitively non-clustered) distribution. INDICATE requires no a priori knowledge of, and makes no assumptions about, the shape of a distribution, presence/number of clusters and/or sub-structure of a dataset.

For any application of INDICATE there are three variable parameters: (1) size and (2) number density of the distribution it is being applied to; and (3) the  $N$ th nearest neighbour number used by the tool. We calibrated our tool against random distributions to define statistically significant values of  $I_{j,N}$  (see appendices) finding:

1. There is a logarithmic relationship between the maximum  $I_{j,N}$  value for a random distribution and sample size (Eq. (A.2)).
2.  $I_{j,N}$  is independent of a distributions number density.
3. There is a relationship between the typical modal  $I_{j,N}$  value for a random distribution and the chosen  $N$ th nearest neighbour number (Eq. (A.3)).
4. The size of the control distribution is essentially arbitrary, as the modal difference between  $I_{j,N}$  values calculated for a random distribution using the standard size (Sect. 2.1) and an expanded size (Appendix B) is inversely proportional the the  $N$ th nearest neighbour number (Eq. (B.3)). However, care should be taken when including/excluding points which are on the boundary of a chosen significance threshold value of  $I_{j,N}$  based solely on their index values during an analysis, particularly for indices derived for small sample sizes using large nearest neighbour numbers.
5. Uniformly distributed (interloping) field stars in observational datasets typically do not significantly affect the index

values of true cluster members. The error on the index value derived for true cluster members is given by Eqs. (C.1)–(C.3).

6. If interloping field stars are distributed in a gradient, the index derived for true cluster members is independent of gradient shape for small nearest neighbour numbers ( $N = 3$ ). However, as field stars are also assigned index values, care must be taken when drawing conclusions on the physical origins of the clustering tendencies of stars in the dataset.

One of the primary strengths of our tool is its versatility and flexibility to be applied to a user-defined analysis. In this paper we demonstrated one potential application of the tool – to look for signals of mass segregation and trace variations in degree of stellar association in star forming regions/clusters.

Arguably the three most popular established methods to identify mass segregation are: (i) Radial Mass Functions (e.g. Sagar et al. 1988), (ii) the  $\Lambda_{MSR}$  parameter (Allison et al. 2009), and (iii) the Local Density Ratio (Maschberger & Clarke 2011; Küpper et al. 2011; Parker 2014). Each have their respective strengths and weaknesses (see Parker & Goodwin 2015 for a discussion), but primarily the decision of which method one employs is based upon what type of mass segregation one is searching for. In the literature, “mass segregation” is used interchangeably to describe two quite different realisations: (1) the concentration of massive stars together at a system’s centre and (2) (lower mass) stellar concentrations around the massive stars in a system but which are not necessarily concentrated together. In our case, we are interested in better understanding the role of local and global environmental conditions in massive star formation. Our aim therefore was to measure the degree of association (or lack thereof) of each high mass star with the general stellar population and with each other, in young (<5 Myr) regions, i.e. look for signatures of both types of mass segregation. For our purpose, it is ideal to employ a single method to search for and quantify signatures of both types in a given region, so that they can be directly compared and a quantitative analysis of the impact of local environment formation conditions on spatial structures undertaken. This is possible using INDICATE, as demonstrated in Sect. 3.2.

Additional strengths of INDICATE are:

- Our tool does not require a priori knowledge of the centre, and works independently of the shape, of the distribution.
- The index has been calibrated against random distributions, so statistically significant values are easily identified.
- As  $I_{j,N}$  is a measure of spatial association (not density), the clustering behaviour (index values) of massive stars in two or more regions can be directly compared, regardless of differences in their distances, average angular separation of sources and/or field sizes.
- It can provide both a global *and* local measure of Type 1 mass segregation. By definition  $I_{j,N}$  is a local measure, and a global measure can be obtained for the subset by e.g. calculating the mean index value of the massive stars and comparing it to that expected by a random distribution (Sect. 3.2).
- Conclusions on Type 1 mass segregation in a system are not based on the larger spatial distribution of other stars in the system as a whole. Index values for high mass stars are derived through comparison to a control distribution, not internally with other sub-samples of the system (low mass stars), and significant values are determined through comparison to those expected in a random distribution. Therefore high mass stars index values are independent of the completeness of the resolved low mass population census.

- As a local measure, INDICATE is robust against outliers as they (a) will not influence the index values of the other members in a subset, (b) are easily identifiable by their comparative low index values and as such (c) in an global analysis of a system to find signatures of Type 1 mass segregation will have a statistically negligible effect on the overall conclusions drawn for the subset.

We applied our tool to the stellar catalogue of the Carina Nebula (NGC 3372) by [Kuhn et al. \(2014\)](#), a region chosen because of its known high mass stellar content (>130 OB stars) and extensive sub-structure:

1. We recover known stellar structure in the region, including the Tr14-16, Treasure Chest and Bochum 11 clusters.
2. We find members of the 4/19 sub-clusters identified by [Kuhn et al. \(2014\)](#) as stellar overdensities are more clustered than typical for the extended distribution of stars in the Carina region, but contain no, or very few, stars with a degree of association above random. This suggests these sub-clusters may be fluctuations in the dispersed population field rather than real clusters.
3. Stars in the NW and SE regions have significantly different clustering tendencies. The NW region is known to be heavily sub-structured, whereas the SE is more sparsely populated and being shaped by radiative winds of the Tr14 and Tr16 clusters ([Smith et al. 2008](#)). Therefore this result is reflective of differences in the apparent star formation activity in these regions. Further study is required to ascertain the physical origin of that difference.
4. The different clustering properties between the NW and SE regions are also seen for OB stars and are even more pronounced.
5. There are no signatures of classical (Type 1) mass segregation present in the SE region – massive stars here are not concentrated together above random.
6. Stellar concentrations are more frequent around massive stars than typical for the general population, particularly in the young Tr14 cluster.
7. For Tr14 and Tr15 we find a relation between the concentration of OB stars and the concentration of (lower mass) stars around OB stars. This relation is notably absent from Tr16. Unlike the Tr14 and Tr15 clusters, Tr16 does not have a strong central concentration but instead is irregularly shaped and heavily sub-structured with multiple sub-clusters ([Ascenso et al. 2007](#), [Wang et al. 2011](#), [Wolk et al. 2011](#)). Therefore this result reflects the known structure of the

clusters: Tr14 and Tr15 are centrally concentrated, whereas in Tr16 the OB stars are not clustered together in a single concentration with a (near) constant degree of clustering, but are instead scattered across a region with local concentrations of stars and a variable degree of association.

*Acknowledgements.* The Star Form Mapper project has received funding from the European Union’s Horizon 2020 research and innovation programme under grant agreement No 687528. We would like to thank Simon Goodwin and Lee Mundy for suggestions for the improvement of our tool and paper respectively.

## References

- Allison, R. J., Goodwin, S. P., Parker, R. J., et al. 2009, *MNRAS*, **395**, 1449
- Andre, P., Ward-Thompson, D., & Barsony, M. 2000, *Protostars and Planets IV*, 59
- Anselin, L. 1995, *Geog. Anal.*, **27**, 93
- Ascenso, J., Alves, J., Vicente, S., & Lago, M. T. V. T. 2007, *A&A*, **476**, 199
- Cartwright, A., & Whitworth, A. P. 2004, *MNRAS*, **348**, 589
- Cartwright, A., & Whitworth, A. P. 2009, *A&A*, **503**, 909
- de Wit, W. J., Testi, L., Palla, F., & Zinnecker, H. 2005, *A&A*, **437**, 247
- Dutra, C. M., & Bica, E. 2001, *A&A*, **376**, 434
- Gaczkowski, B., Preibisch, T., Ratzka, T., et al. 2013, *A&A*, **549**, A67
- Gomez, M., Hartmann, L., Kenyon, S. J., & Hewett, R. 1993, *AJ*, **105**, 1927
- Hopkins, B., & Skellam, J. 1954, *Ann. Bot.*, **18**, 213
- Kharchenko, N. V., Piskunov, A. E., Schilbach, E., Röser, S., & Scholz, R.-D. 2013, *A&A*, **558**, A53
- Kuhn, M. A., Feigelson, E. D., Getman, K. V., et al. 2014, *ApJ*, **787**, 107
- Küpper, A. H. W., Maschberger, T., Kroupa, P., & Baumgardt, H. 2011, *MNRAS*, **417**, 2300
- Luhman, K. L. 2012, *ARA&A*, **50**, 65
- Maschberger, T., & Clarke, C. J. 2011, *MNRAS*, **416**, 541
- Møller, J., & Waagepetersen, R. P. 2007, *Scand. J. Stat.*, **34**, 643
- Parker, R. J. 2014, *MNRAS*, **445**, 4037
- Parker, R. J., & Goodwin, S. P. 2015, *MNRAS*, **449**, 3381
- Parker, R., Wright, N., Goodwin, S., & Meyer, M. 2014, *MNRAS*, **438**, 620
- Preibisch, T., Schuller, F., Ohlendorf, H., et al. 2011a, *A&A*, **525**, A92
- Preibisch, T., Ratzka, T., Kuderna, B., et al. 2011b, *A&A*, **530**, A34
- Sagar, R., Miakutin, V. I., Piskunov, A. E., & Dluzhnevskaja, O. B. 1988, *MNRAS*, **234**, 831
- Scalo, J., & Chappell, D. 1999, *ApJ*, **510**, 258
- Schmeja, S., Kumar, M. S. N., & Ferreira, B. 2008, *MNRAS*, **389**, 1209
- Shu, F. H., Adams, F. C., & Lizano, S. 1987, *ARA&A*, **25**, 23
- Smith, N., & Brooks, K. J. 2008, in *The Carina Nebula: A Laboratory for Feedback and Triggered Star Formation*, ed. B. Reipurth, 138
- Smith, N., Povich, M. S., Whitney, B. A., et al. 2010, *MNRAS*, **406**, 952
- Wang, J., Feigelson, E. D., Townsley, L. K., et al. 2011, *ApJS*, **194**, 11
- Wolk, S. J., Broos, P. S., Getman, K. V., et al. 2011, *ApJS*, **194**, 12
- Wright, N. J., Parker, R. J., Goodwin, S. P., & Drake, J. J. 2014, *MNRAS*, **438**, 639
- Zinnecker, H., & Yorke, H. W. 2007, *ARA&A*, **45**, 481

## Appendix A: Calibration of the index

**Table A.1.** Constants of Eq. (A.2) for a  $N$ th nearest neighbour number of  $N = 3, 5, 7$  and  $9$  with their respective fit correlation coefficient (R) and standard error (SE).

$N$	$C_1$	$C_2$	R	SE
3	2.508	0.489	0.831	0.047
5	2.291	0.361	0.851	0.032
7	2.244	0.206	0.796	0.023
9	2.093	0.197	0.777	0.023

We conduct a series of baseline tests to aid interpretation and identification of significant index values. These tests (a) define the threshold at which an index value becomes significant i.e. the value above which it can reasonably be assumed point  $j$  was not drawn from a random distribution and (b) quantify the impact of dataset parameters and the choice of  $N$ th nearest neighbour number on the distribution and range of index values INDICATE generates.

### (i) Sample size

We generate random samples of size,  $S$ , in the range  $50 \leq S < 100\,000$ . For each sample size 100 realisations are created with a constant number density of  $n_{\text{obs}} = 1$  object per unit area and INDICATE is implemented with a  $N$ th nearest neighbour number of  $N = 5$ . We keep the number density and  $N$ th nearest neighbour number for every sample constant to ensure any identified trends or patterns in samples' index values can be attributed to sample size alone.

There is no dependence between the index and the size of a sample, with typical modal and mean values of  $\text{Mo}[I_5] = 0.8$  and  $\bar{I}_5 = 1.0$  for random distributions under the above stated conditions. However, we find there is a logarithmic relationship between the upper range limit<sup>4</sup> of  $I_{j,N}$  and sample size for random distributions, i.e.

$$0.0 \leq I_{j,N} \leq I_{\text{max}}, \quad (\text{A.1})$$

where

$$I_{\text{max}} = C_1 + C_2 \times \log_{10} S, \quad (\text{A.2})$$

and  $C_1, C_2$  are constants which are dependant on the  $N$ th nearest neighbour number (see Table A.1). Equation (A.2) defines as a function of sample size the threshold value above which we can definitively assume a point does not have a spatially random distribution.

### (ii) Field density

We generate random samples of number density,  $n_{\text{obs}}$ , in the range  $10^{-6} \leq n_{\text{obs}} \leq 10^6$ , in increments of an order of magnitude. For each value of number density 100 realisations are created, with a constant sample size of 10 000 and INDICATE is implemented with a  $N$ th nearest neighbour number of  $N = 5$ . We find there is no dependence of  $I_{j,N}$  on field density.

### (iii) $N$ th nearest neighbour number, $N$ .

We generate a 100 realisations of random samples of size  $S = 10\,000$  and number density  $n_{\text{obs}} = 1$ . For each sample INDICATE is implemented with a  $N$ th nearest neighbour number of  $N = 3, 5, 7$  and  $9$ . There is a relationship between the

upper range limit of  $I_{j,N}$ , sample size and  $N$ th nearest neighbour number (Eq. (A.2), Table A.1). The typical modal index value,  $\text{Mo}[I_{j,N}]$ , of randomly distributed samples vary as a function of  $N$ :

$$\text{Mo}[I_{j,N}] \equiv \frac{N-1}{N}. \quad (\text{A.3})$$

The typical mean index values of randomly distributed samples are  $0.9 \leq \bar{I}_N \leq 1.0$ .

## Appendix B: Investigation of edge effects

Section 2.1.2 described how the control distribution used by INDICATE is generated. Here we investigate whether the proximity of a point in a dataset to its delimited boundaries and/or the total length of each axis of the control distribution influences a sample's index values. We repeat the calibration tests (Appendix A) using two different types of control distribution:

1. *Control distribution A* (CDA) – occupies the same bounded parameter space and has the same number density (Eq. (2)) as the test sample;
2. *Control distribution B* (CDB) – occupies the same, and is extended beyond the, bounded parameter space of the test sample; such that area of the control distribution is a factor of four times larger than the test sample (see Fig. B.1). Increasing the area of the control distribution by a factor of four ensures that the  $r_j$  of edge points in the test sample (Eq. (4)) is not calculated using edge points of the control distribution (which in principle could subsequently increase  $\bar{r}$ , and decrease  $I_{j,N}$ ). It has the same number density as the test sample (Eq. (2)).<sup>5</sup>

We define an ‘‘edge point’’ as any point in the sample dataset whose  $(x,y)$  position is less than that of the second smallest  $x$  and/or  $y$  positions and/or greater than the second largest  $x$  and/or  $y$  positions of points in CDA (i.e. where the measured nearest neighbour distance of the sample point to the control distribution points would be affected due to lack of control points in any given direction in the control distribution for point  $j$ ).

For  $N > 5$  the modal index value of edge points,  $\text{Mo}[I_{j,N}^E]$ , deviates from that of the sample as a whole using CDA (Eq. (A.3)), such that

$$\text{Mo}[I_{j,7}^E] \equiv \frac{N-2}{N} \quad \text{for} \quad N = 7, \quad (\text{B.1})$$

$$\text{Mo}[I_{j,9}^E] \equiv \frac{N-3}{N} \quad \text{for} \quad N = 9. \quad (\text{B.2})$$

For a proportion of all (edge and non-edge) points in the test samples' there is a statistically small discrepancy between the index values calculated using CDA and CDB. The modal difference between the two sets of indices is inversely proportional to  $N$  i.e.

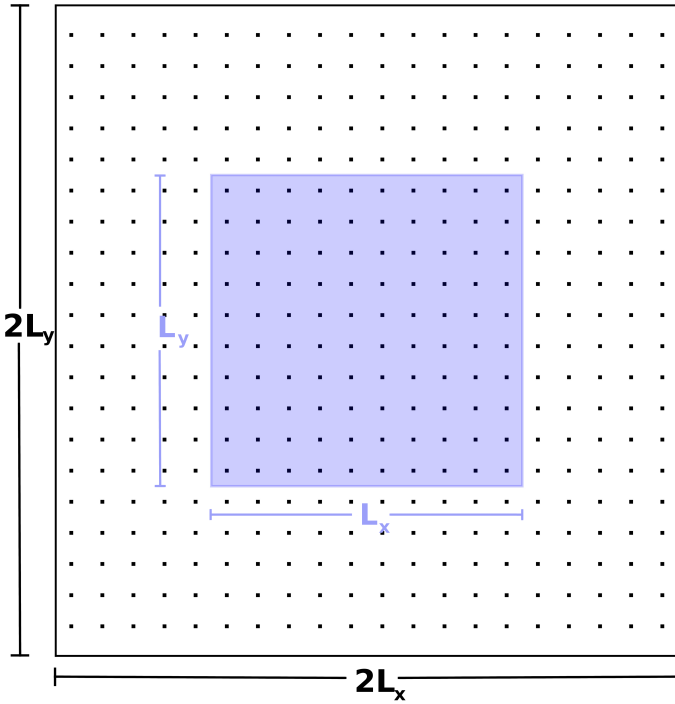
$$\text{Mo}[\Delta I_{j,N}] \equiv \text{Mo}[I_{j,N}^{\text{CDA}} - I_{j,N}^{\text{CDB}}] \equiv \frac{1}{N} \quad \text{for} \quad N \geq 3, \quad (\text{B.3})$$

where  $I_{j,N}^{\text{CDA}}$  and  $I_{j,N}^{\text{CDB}}$  are the index values calculated for each sample point  $j$  using CDA and CDB respectively. For any given point  $j$  if

$$\Delta I_{j,N} > 0 \quad \leftrightarrow \quad I_{j,N}^{\text{CDA}} > I_{j,N}^{\text{CDB}}. \quad (\text{B.4})$$

<sup>5</sup> For samples with non-rectangular delimited areas this distribution should always be used.

<sup>4</sup> Derived as the maximum value over all realisations.



**Fig. B.1.** Dimensions of control distributions CDA (blue shaded) and CDB (all visible) as described in Appendix B, where  $L_x$  and  $L_y$  are the length of a test sample's  $x$  and  $y$  axis respectively; and the black dots are the points of the control distributions.

The proportion of all points with  $\Delta I_{j,N} > 0$  increases with decreasing sample size and increasing nearest neighbour number, reaching  $\sim 90\%$  for sample size of  $S = 50$  using  $N = 9$ ; it is independent of field density. The number of edge points with  $\Delta I_{j,N} > 0$  is proportionally lower than non-edge points i.e. expanding the control distribution has less of an effect on edge points than non-edge points. This is because the  $r_j$  measured for edge points in CDB is (slightly) smaller than in CDA (as it is no longer artificially increased due a lack of control points in any given direction), which subsequently causes a small decrease in  $\bar{r}$  (Eq. (3)). In both control distributions a radius of  $\bar{r}$  from an edge point can partially encompass an area outside the bounds of the dataset (where there can be no neighbouring points), but for non-edge points a radius of  $\bar{r}$  always encompasses an area within the bounds of the dataset (neighbouring points can be present in any given direction within  $\bar{r}$ ). Thus a small decrease in  $\bar{r}$  is more likely to exclude a nearest neighbour (decrease  $N_{\bar{r}}$ , and subsequently  $I_{j,N}$  – Eq. (5)) for a non-edge point than an edge point.

To conclude, as the typical  $\Delta I_{j,N}$  for any given point between the two control distributions is very small, choice of control distribution type (CDA or CDB) is essentially arbitrary, but care should be taken when including/excluding points which are on the boundary of a chosen significance threshold value of  $I_N$  during an analysis – particularly indices derived for small sample sizes using large nearest neighbour numbers.

### Appendix C: Investigation of field effects

To ascertain the influence of interlopers on the index values of true cluster members we conduct additional calibration tests. A dataset consisting of a Gaussian cluster with 500 members is generated and the index value of each member determined using the steps outlined in Sect. 2.1.

**Table C.1.** Constants of Eqs. (C.2), (C.3) for a  $N$ th nearest neighbour number of  $N = 3, 5, 7$  and  $9$ .

$N$	$C_3$	$C_4$	$C_5$	$C_6$	$C_7$
3	2.549	0.510	4.725	0.830	-6.077
5	2.033	0.438	3.394	2.133	-3.844
7	1.617	0.320	3.715	1.467	-4.140
9	1.421	0.306	3.156	1.867	-3.559

In our first test, field stars are introduced to the dataset with incrementally increasing frequency, such that the number of interloping field stars at any given time is equal to a fraction,  $F$ , of total cluster members in the range  $0.01 \leq F \leq 1.0$ . The positions of the field stars are randomly drawn from a uniform distribution. For each fraction of field stars 100 realisations are made, and for each realisation the difference,  $\Delta I_{j,N}$ , between the index values derived for cluster members in the dataset that does not contain field stars and the current level of field star contamination is measured for a  $N$ th nearest neighbour number of  $N = 3, 5, 7$  and  $9$ . As we are simulating an observational dataset for which cluster membership is uncertain,  $N_{\text{tot}} = S = 500 + (F \times 500)$ .

We find the modal difference for all combinations of  $F$  and  $N$  is  $\text{Mo}[\Delta I_{j,N}] = 0$ , i.e. typically the index values of true cluster members are unaffected by the presence of interloping field stars. The proportion of cluster members with  $\Delta I_{j,N} \neq 0$  increases with increasing  $F$  and  $N$ , reaching a maximum of  $\sim 95\%$  for  $F = 1.0$  and  $N = 9$ . In observationally obtained datasets the error on the index value derived for true cluster members is therefore

$$I_{j,N} + F2 \leq I_{j,N} \leq I_{j,N} + F1, \quad (\text{C.1})$$

where

$$F1 = \max[\Delta I_{j,N}] = C_3 + C_4 \times \log(F), \quad (\text{C.2})$$

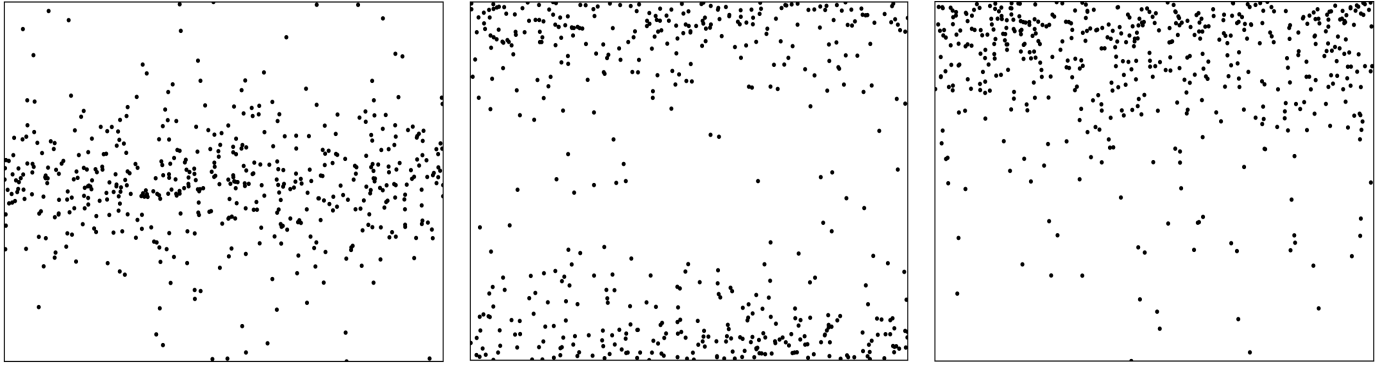
$$F2 = \min[\Delta I_{j,N}] = C_5 \times \exp(C_6 \times F) + C_7, \quad (\text{C.3})$$

and  $C_{3-7}$  are constants dependant on the  $N$ th nearest neighbour number (see Table C.1).

In our second test  $F = 1.0$  field stars are distributed in three large scale gradient patterns (Fig. C.1) which are randomly generated in the same parameter space as the Gaussian cluster. For each gradient 100 realisations are made, and for each realisation the difference,  $\Delta I_{j,N}$ , between the index values derived for cluster members in the dataset that does not contain field stars and the current level of field star contamination is measured for a  $N$ th nearest neighbour number of  $N = 3, 5, 7$  and  $9$ . As we are simulating an observational dataset for which cluster membership is uncertain,  $N_{\text{tot}} = S = 500 + (F \times 500) = 1000$ .

We find the modal difference for all gradients with  $N = 3$  is  $\text{Mo}[\Delta I_{j,N}] = 0$ , i.e. for small values of  $N$  the index derived for cluster members is independent of gradient shape. This is expected as the index is a local measure, and the value of  $N$  essentially defines its resolution (the smaller  $N$ , the higher the resolution). Thus index values are more susceptible to the effects of variation in the degree of field star association within the gradient when larger values of  $N$  are employed.

As noted previously, INDICATE is distance independent for a fully resolved dataset. However, in practice, clearly INDICATE cannot detect unresolved binaries and higher order systems in datasets nor a priori know any difference between a member of a grouping and a fore- or background field star. Even with



**Fig. C.1.** A realisation of the three gradient field population shapes tested in Appendix C.

best efforts, not all field stars will be removed from observationally obtained datasets before analysis, so consideration must be given before drawing conclusions about the clustering tendencies of region stars. In particular, when a pronounced large scale 2D spatial distribution gradient of the field population is present, and cluster membership is uncertain, caution must be taken when drawing conclusions about the physical origins of the clustering tendencies of stars – as field stars within the denser

regions of the gradient naturally will have a higher degree of association and thus index. Similar care must be taken when interpreting index values for 2D datasets in which a smaller angular resolution cluster is superimposed onto a larger angular resolution cluster, or that contains two clusters at significantly different distances. Simulations and bootstrapping techniques can be used to test the magnitude of such effects on individual datasets.