

Publication Year	2021					
Acceptance in OA	2022-07-13T14:43:43Z					
Title	The SEDIGISM survey: molecular clouds in the inner Galaxy					
Authors	<ul> <li>Duarte-Cabral, A., Colombo, D., Urquhart, J. S., Ginsburg, A., Russeil, D., Schuller, F., Anderson, L. D., Barnes, P. J., BELTRAN SOROLLA, MARIA TERESA, Beuther, H., Bontemps, S., Bronfman, L., Csengeri, T., Dobbs, C. L., Eden, D., GIANNETTI, ANDREA, Kauffmann, J., Mattern, M., Medina, SN. X., Menten, K. M., Lee, MY., Pettitt, A. R., Riener, M., Rigby, A. J., TRAFICANTE, ALESSIO, Veena, V. S., Wienen, M., Wyrowski, F., Agurto, C., Azagra, F., CESARONI, Riccardo, Finger, R., Gonzalez, E., Henning, T., Hernandez, A. K., Kainulainen, J., Leurini, Silvia, Lopez, S., Mac-Auliffe, F., Mazumdar, P., MOLINARI, Sergio, Motte, F., Muller, E., Nguyen-Luong, Q., Parra, R., Perez-Beaupuits, JP., Montenegro-Montes, F. M., Moore, T. J. T., Ragan, S. E., Sánchez-Monge, A., Sanna, A., Schilke, P., SCHISANO, EUGENIO, Schneider, N., Suri, S., TESTI, Leonardo, Torstensson, K., Venegas, P., Wang, K., Zavagno, A.</li> </ul>					
Publisher's version (DOI)	10.1093/mnras/staa2480					
Handle	http://hdl.handle.net/20.500.12386/32471					
Journal	MONTHLY NOTICES OF THE ROYAL ASTRONOMICAL SOCIETY					
Volume	500					

# The SEDIGISM survey: molecular clouds in the inner Galaxy

A. Duarte-Cabral <sup>(b)</sup>, <sup>1\*</sup> D. Colombo, <sup>2\*</sup> J. S. Urquhart <sup>(b)</sup>, <sup>3\*</sup> A. Ginsburg, <sup>4</sup> D. Russeil, <sup>5</sup> F. Schuller <sup>(b)</sup>, <sup>2</sup> L. D. Anderson, <sup>6</sup> P. J. Barnes, <sup>7,8</sup> M. T. Beltrán, <sup>9</sup> H. Beuther, <sup>10</sup> S. Bontemps, <sup>11</sup> L. Bronfman, <sup>12</sup> T. Csengeri, <sup>11</sup> C. L. Dobbs, <sup>13</sup> D. Eden <sup>(b)</sup>, <sup>14</sup> A. Giannetti, <sup>2</sup> J. Kauffmann, <sup>15</sup> M. Mattern, <sup>2</sup> S.-N. X. Medina, <sup>2</sup> K. M. Menten, <sup>2</sup> M.-Y. Lee, <sup>2,16</sup> A. R. Pettitt <sup>(b)</sup>, <sup>17</sup> M. Riener, <sup>10</sup> A. J. Rigby <sup>(b)</sup>, <sup>1</sup> A. Traficante <sup>(b)</sup>, <sup>18</sup> V. S. Veena, <sup>19</sup> M. Wienen, <sup>2</sup> F. Wyrowski, <sup>2</sup> C. Agurto, <sup>20</sup> F. Azagra, <sup>20</sup> R. Cesaroni, <sup>9</sup> R. Finger, <sup>12</sup> E. Gonzalez, <sup>20</sup> T. Henning, <sup>10</sup> A. K. Hernandez, <sup>21</sup> J. Kainulainen, <sup>2,22</sup> S. Leurini, <sup>2,23</sup> S. Lopez, <sup>4</sup> F. Mac-Auliffe, <sup>20</sup> P. Mazumdar, <sup>2</sup> S. Molinari <sup>(b)</sup>, <sup>18</sup> F. Motte, <sup>24</sup> E. Muller, <sup>25</sup> Q. Nguyen-Luong, <sup>15</sup> R. Parra, <sup>20</sup> J.-P. Perez-Beaupuits <sup>(c)</sup>, <sup>20</sup> F. M. Montenegro-Montes, <sup>20</sup> T. J. T. Moore, <sup>14</sup> S. E. Ragan <sup>(b)</sup>, <sup>1</sup> A. Sánchez-Monge, <sup>19</sup> A. Sanna, <sup>2</sup> P. Schilke, <sup>19</sup> E. Schisano <sup>(c)</sup>, <sup>18</sup> N. Schneider <sup>(c)</sup>, <sup>19</sup> S. Suri, <sup>19</sup> L. Testi, <sup>19</sup> K. Torstensson, <sup>20</sup> P. Venegas, <sup>20</sup> K. Wang<sup>26</sup> and A. Zavagno<sup>5</sup>

Affiliations are listed at the end of the paper

Accepted 2020 May 22. Received 2020 May 22; in original form 2019 October 10

# ABSTRACT

We use the <sup>13</sup>CO (2–1) emission from the SEDIGISM (Structure, Excitation, and Dynamics of the Inner Galactic InterStellar Medium) high-resolution spectral-line survey of the inner Galaxy, to extract the molecular cloud population with a large dynamic range in spatial scales, using the Spectral Clustering for Interstellar Molecular Emission Segmentation (SCIMES) algorithm. This work compiles a cloud catalogue with a total of 10663 molecular clouds, 10300 of which we were able to assign distances and compute physical properties. We study some of the global properties of clouds using a science sample, consisting of 6664 well-resolved sources and for which the distance estimates are reliable. In particular, we compare the scaling relations retrieved from SEDIGISM to those of other surveys, and we explore the properties of clouds with and without high-mass star formation. Our results suggest that there is no single global property of a cloud that determines its ability to form massive stars, although we find combined trends of increasing mass, size, surface density, and velocity dispersion for the sub-sample of clouds with ongoing high-mass star formation. We then isolate the most extreme clouds in the SEDIGISM sample (i.e. clouds in the tails of the distributions) to look at their overall Galactic distribution, in search for hints of environmental effects. We find that, for most properties, the Galactic distribution of the most extreme clouds is only marginally different to that of the global cloud population. The Galactic distribution of the largest clouds, the turbulent clouds and the high-mass star-forming clouds are those that deviate most significantly from the global cloud population. We also find that the least dynamically active clouds (with low velocity dispersion or low virial parameter) are situated further afield, mostly in the least populated areas. However, we suspect that part of these trends may be affected by some observational biases (such as completeness and survey limitations), and thus require further follow up work in order to be confirmed.

Key words: ISM: clouds - Galaxy: structure - stars: formation - galaxies: ISM, star formation.

### **1 INTRODUCTION**

The evolution of the gas that makes up the interstellar medium (ISM), and the ultimate means by which that gas gives way to star formation, involve the tight interplay of a wealth of physical processes. Our understanding of those processes has relied upon the statistical characterization of the molecular gas that is taking part in the star formation process. In particular, the star formation field has relied on a discretization of the molecular component of

the ISM into molecular clouds, across the Galactic disc, either as observed in 2D with dust continuum emission (e.g. the ATLASGAL survey, Schuller et al. 2009; the Hi-GAL survey, Molinari et al. 2010; or the Bolocam Galactic Plane Survey, Rosolowsky et al. 2010, Ginsburg et al. 2013), or with the 3D view of the Galactic plane from spectral-line observations, most commonly using the second-most abundant molecular species in the ISM, the CO molecule (and its isotopologues). Large survey observations of the Galactic plane in CO emission have allowed for a number of statistical studies of molecular clouds across the Galaxy (e.g. Scoville & Solomon 1975; Larson 1981; Solomon et al. 1987; Heyer et al. 2009; Roman-Duval et al. 2010; Rice et al. 2016; Miville-Deschênes, Murray & Lee 2017), and have provided a large-scale view of the distribution of gas

doi:10.1093/mnras/staa2480

<sup>\*</sup> E-mail: adc@astro.cf.ac.uk (ADC); dcolombo@mpifr-bonn.mpg.de (DC); J.S.Urquhart@kent.ac.uk (JSU)

in the Milky Way, crucial for our understanding of its spiral structure (e.g. Dame, Hartmann & Thaddeus 2001; Vallée 2014; Pettitt et al. 2014, 2015).

These Galactic plane surveys, alongside some resolved studies of molecular clouds in nearby spiral galaxies, have also suggested a number of scaling relations (namely between the sizes of clouds, their linewidths, and their mass surface densities, e.g. Larson 1981; Solomon et al. 1987; Heyer et al. 2009; Sun et al. 2018), as well as some differences in the mass spectra of clouds towards different environments (e.g. Colombo et al. 2014; Rice et al. 2016; Miville-Deschênes et al. 2017). All of these findings have implications in our interpretation of the global properties of molecular clouds, and how they might evolve. Most of these surveys, however, were finding and describing molecular clouds that had typical sizes close to their resolution element - which can bias the interpretation of the results - and given their lower resolution they could also potentially suffer from severe blending of the emission along the same line of sight, especially with our edge-on perspective of the Milky Way (e.g. Duarte-Cabral et al. 2015; Duarte-Cabral & Dobbs 2016).

With the advent of new high-resolution and large-scale spectroscopic surveys of the Galactic plane (such as the Structure, Excitation, and Dynamics of the Inner Galactic InterStellar Medium survey - SEDIGISM, Schuller et al. 2017; the CO High Resolution Survey – COHRS, Dempsey, Thomas & Currie 2013; the <sup>13</sup>CO/C<sup>18</sup>O (J = 3-2) Heterodyne Inner Milky Way Plane Survey – CHIMPS, Rigby et al. 2016; the Three-mm Ultimate Mopra Milky Way Survey - ThrUMMS, Barnes et al. 2015; or the Galactic Census of High and Medium-mass Protostars - CHaMP, Barnes et al. 2011), not only are these shortcomings now greatly minimized, but we can start to explore the details of the substructure within molecular clouds where star formation is actively taking place, and the clouds' link to the large-scale Galactic environment. This opens a new and exciting era in the study of star formation in a Galactic context. Given that molecular clouds are highly hierarchical systems, it is essential to be able to define molecular clouds with a large dynamic range in spatial scales (e.g. as in Colombo et al. 2019), and this is at the heart of this work. In this paper, we explore the global properties of molecular clouds from the high-resolution <sup>13</sup>CO (2-1) emission from the SEDIGISM survey, covering the inner Galactic plane (from  $+300^{\circ} < \ell < +18^{\circ}$ , Schuller et al. 2017), which is described in Section 2. Section 3 contains the details of the method used for the extraction of molecular clouds from this data set, along with a description of all the derived properties and data products released with the molecular cloud catalogue. In Section 4, we describe the methods used to determine the distances and distinguish between derived near/far kinematic distances to all the clouds in the catalogue, essential to derive the physical properties. In Section 5, we explore the distributions of the global properties of the SEDIGISM clouds, and also compare these with other samples in the literature. In Section 6, we explore possible indications of environmental dependence of cloud properties, by isolating the most extreme clouds (i.e. clouds in the tails of the distributions), and comparing their Galactic distribution with that of the entire cloud population. Finally, our findings are summarized in Section 7.

### 2 DATA

In this paper, we use data from the SEDIGISM survey conducted with the Atacama Pathfinder Experiment 12 m submillimetre telescope (APEX; Güsten et al. 2006). In particular, we use the <sup>13</sup>CO (2–1) to extract and characterize the molecular clouds towards the inner Galaxy. The complete details on the observations, data reduction

and data-quality checks can be found in the survey overview papers (Schuller et al. 2017,2020).

In summary, the SEDIGISM survey observed a total of 84 deg<sup>2</sup>, covering from  $-60^{\circ} \le \ell \le +18^{\circ}$ , and  $|b| \le 0.5^{\circ}$ , plus a few extensions in *b* towards some regions, as well as an additional field towards the W43 region ( $+29^{\circ} \le \ell \le +31^{\circ}$ ). The <sup>13</sup>CO (2–1) data that we use here is the DR1 data set (fully described in Schuller et al., in preparation), which has a typical  $1\sigma$  sensitivity of 0.8–1.0 K (in  $T_{\rm mb}$ ) per 0.25 km s<sup>-1</sup> channel, and an FWHM beam size,  $\theta_{\rm MB}$ , of 28 arcsec.

In this paper, we will use the complete contiguous data set (i.e. the entire survey data except for the W43 field). This consists of 77 datacubes of roughly  $2^{\circ} \times 1^{\circ}$  (note that the latitude range is sometimes larger than  $1^{\circ}$ ), centred at all integer longitudes between  $\ell = 301^{\circ}$  and  $\ell = 17^{\circ}$  (i.e. spaced by  $1^{\circ}$  in longitude). This provides a  $1^{\circ}$  overlap in longitude between consecutive tiles, which ensures all sightlines (except for the first and last fields) are contained in two tiles. The velocity ranges from  $-200 \text{ to } +200 \text{ km s}^{-1}$  in all datacubes, and the pixel size is of 9.5 arcsec. Fig. 1 shows the full  $\ell v$  map of the contiguous data set from the SEDIGISM survey, that we use here.

### **3 MOLECULAR CLOUD EXTRACTION**

### 3.1 The method: SCIMES

In order to decompose the <sup>13</sup>CO emission from the SEDIGISM survey into discrete clouds, we use the Spectral Clustering for Interstellar Molecular Emission Segmentation (SCIMES) algorithm (v.0.3.2).<sup>1</sup> The original algorithm is fully described in Colombo et al. (2015), and the improvements included in the version we use here are detailed in Colombo et al. (2019). In brief, SCIMES brings a significant advancement with respect to other more commonly used cloud-extraction algorithms (e.g. Clumpfind by Williams, de Geus & Blitz 1994, Gaussclumps by Stutzki & Guesten 1990, or Fellwalker by Berry 2015), as it is a fully automated method that uses spectral clustering and graph theory to analyse the dendrogram of the emission, and decompose the hierarchical structure of the ISM into 'clusters' of molecular gas emission (i.e. molecular clouds, considering the resolution of SEDIGISM). Unlike other cloudextraction algorithms, SCIMES relies on the natural transitions in the emission to define discrete structures, and it is robust against changes in the input parameters (as demonstrated in Colombo et al. 2015).

The cloud extraction with SCIMES was performed on each of the 77 tiles of  $2^{\circ} \times 1^{\circ}$ . We ran SCIMES on these relatively small cubes because it would be extremely computationally expensive (and memory intensive) to generate a single dendrogram from the full SEDIGISM data set, and perform SCIMES's affinity matrix analysis, where each cluster is equivalent to an additional dimension in the clustering space.

### 3.1.1 Input parameters and files

In order to optimize the performance of the SCIMES clustering algorithm, we have performed a few preparation steps on the original DR1 data. First, we enhanced the signal-to-noise (S/N) ratio of the data set prior to running SCIMES by smoothing the data in velocity. This was done by binning the data into  $0.5 \text{ km s}^{-1}$  channels. We then re-sampled these binned datacubes back into  $0.25 \text{ km s}^{-1}$  channels (using linear interpolation), simply so that the SCIMES assignment



**Figure 1.** Longitude–velocity ( $\ell v$ ) map of the <sup>13</sup>CO peak intensity (in grey-scale) for the SEDIGISM coverage analysed in this paper. The peak intensity map was built after masking out voxels of the <sup>13</sup>CO datacube with intensities  $< 2.5\sigma_{rms}$  (estimated locally for each line of sight). The clouds extracted with SCIMES are overlaid as colours, where each cloud has a different (random) colour.



Figure 2. Histogram of the rms noise level of the entire survey, from the velocity-smoothed datacubes that we use for the SCIMES extraction, showing that it peaks at  $\sim 0.7$  K, with a median value of 0.78 K.

masks (Section 3.2) kept the same format as the original emission datacubes from DR1 (essential to have straight forward voxel-by-voxel match between the DR1 emission maps and the clouds' assignment masks). We have performed some tests on the science demonstration field (Schuller et al. 2017), with binned and nonbinned data, and this step allows us to remove high-frequency noise spikes, speeding up the dendrogram construction and the SCIMES clustering, with minimal loss in the information retrieved.

Secondly, given that the noise in the survey is not perfectly uniform (due to different observing weather conditions), it is also essential

to mask the datacubes using the local noise level, in order to prevent high-noise regions from being used in the dendrogram tree, and incorrectly identified as clouds. For this purpose, we estimated the local noise level at each pixel (i.e. each line of sight) in the velocitysmoothed datacubes, by taking the first 50 channels (which are line-free, and on the high-frequency end, i.e. at negative systemic velocities), and computing the  $1\sigma$  standard deviation. Fig. 2 shows the distribution of this local  $1\sigma$ -rms noise level for all the pixels in our velocity-smoothed data set, showing that it peaks at ~0.7 K. We then create a mask of each datacube, by setting any 3D pixels (voxels) whose emission is lower than  $2\sigma$  of the local noise to zero. Note that since we already go down to  $2\sigma$  of the local noise, we do not perform any dilation of the masks after this step (which is a technique sometimes used to remove potential breaks in clouds in low S/N areas).

Using these masked datacubes, we computed the dendrogram tree of the 3D structures in the data (using the ASTRODENDRO<sup>2</sup> implementation, which is based on the original IDL procedures from Rosolowsky et al. 2008). The dendrogram is composed of three types of structures: *leaves*, which are at the top of the hierarchy and contain no substructure, i.e. they are associated with local peaks of emission; *branches*, which split into multiple substructures; and the *trunk*, which is at the bottom of the hierarchy (i.e. it has no parent structure), and comprises all *branches* and *leaves*. We built our dendrograms using the same input parameters as in the science

<sup>2</sup>http://www.dendrograms.org

demonstration field (Schuller et al. 2017): we considered a noise level ( $\sigma_{\rm rms}$ ) of 0.7 K for all tiles (corresponding to the peak of the noise distribution in Fig. 2), a  $4\sigma_{\rm rms}$  value as the minimum difference between two peaks for them to be considered as separate structures, and a lower threshold for detection of  $2\sigma_{\rm rms}$ , to maximize the connections between different structures at contiguous lower intensity levels.<sup>3</sup> Note that we specifically chose to use a single fixed value of  $\sigma_{\rm rms}$  to build the dendrograms across the entire survey (rather than using a local S/N ratio approach) so that we could define our structures using a uniform criterion throughout. This not only makes it easier to replicate our results using other data sets, but it also ensures that the type of structures we extract are equivalent throughout the entire survey, and not dependent on the local noise conditions.<sup>4</sup> The choice of a  $4\sigma_{\rm rms}$  for the significance of individual peaks, coupled with the fact that the dendrogram was built from datacubes that had been masked based on the pixel-based noise level, was so as to maximize the retained detailed information encoded in the survey, whilst minimizing the inclusion of noise spikes. In addition, we set a minimum number of voxels for a structure to be considered as real to be six times the number of pixels per beam  $(N_{\text{ppbeam}} = 9)$ , so that structures are both resolved spatially (i.e. at least three beams), and in velocity (spanning at least two channels, which corresponds to our effective velocity resolution in the smoothed datacubes). Note, however, that the ASTRODENDRO implementation that we use to build the dendrogram does not separate the spatial axes from the spectral axis. In practice, this means that this criterion will still allow some clouds to be retained whilst being unresolved in one of the axes. Those sources are dealt with in a post-processing step (see Section 3.1.3).

Once the dendrograms were constructed for each tile, we ran SCIMES using both the 'volume' and the 'flux' (which in our case, refers simply to the surface brightness) as the clustering properties (cf. Colombo et al. 2015). This extraction recovered a total of 20 387 gas clusters from the 77 tiles, but most of these are duplicated due to the overlap between consecutive fields. In order to build the final catalogue (and respective assignment masks), we performed a cleaning up procedure to handle clouds in overlapping areas. This is described in the following section.

### 3.1.2 Handling clouds in overlapping regions

In order to handle the clouds that appear in overlapping areas, we have followed a procedure similar to that used by Colombo et al. (2019). This procedure is schematically described in Fig. 3. In essence, we have split our data set into a *main* run (which is composed of all tiles centred at odd longitudes), and a *secondary* run (which is composed of all tiles centred at even longitudes). We then exclude all objects that touch a tile edge on the longitude axis, since their contours are not closed, and they should be fully recovered in the complementary run. We only made an exception for objects that touch the first and

<sup>3</sup>These values are solely defined by the data quality, but tests using slight variations for the different parameters for the dendrogram construction were performed as part of our work on the science demonstration field (Schuller et al. 2017). Those tests have shown that the SCIMES clustering algorithm is robust against small differences on the parameters used to construct the dendrogram.

<sup>4</sup>Although the choice of a unique value for  $\sigma_{\rm rms}$  does require an extra postprocessing step to ensure that detected structures also have a high local S/N ratio (see Section 3.1.3), doing the cloud extraction on an S/N map with the SCIMES algorithm (which does not segment clouds at a fixed brightness threshold) means that we could end up with structures identified across the survey using different criteria, which is non-ideal.



**Figure 3.** Schematic sketch of the procedure used to decide which clouds to retain from the two overlapping runs, namely the removal of clouds that touch the edge of the tiles (which are then recovered in full in the complementary run), and the selection of the larger clouds for overlapping cases between the two runs.

last longitude edges of the contiguous coverage (i.e. at  $\ell = 18^{\circ}$  and  $\ell = 300^{\circ}$ ), which are retained in the final catalogue with a tag that indicates that they are edge clouds. Similarly, we also retain clouds that touch the survey's upper and lower latitude edges, and tag them as being edge clouds. Finally, we proceed to checking the matches between the *main* and *secondary* runs. We start by including all objects that do not overlap between the two runs, and whenever two (or more) clouds overlap, we simply retain the larger object between the two runs. After this procedure, we have compiled a total of 11 638 unique molecular clouds.

### 3.1.3 Removal of spurious sources

As mentioned in Section 3.1.1, despite our best efforts to avoid having any noisy spikes in the dendrogram (by imposing a noise level threshold) or unresolved sources (by imposing a minimum number of voxels), some spurious sources still persist to the dendrogram construction and into our final catalogue. One of the reasons for this is the fact that we have applied an average noise  $\sigma_{\rm rms}$  for the entire survey (so that the dendrogram for all fields was built upon a fixed physical value of emission intensity). This means that in areas where the local noise level is higher than this average  $\sigma_{\rm rms}$ , some noisy peaks would have been considered as robust emission peaks. Most of these sources are located near the noisier edges of the observed fields, and are relatively small (close to the beam size). We therefore applied the following selection criteria to remove spurious sources from the final catalogue: (1) any source touching an edge that has a projected (footprint) size of less than five beams<sup>5</sup> (where the angular size of the beam is taken to be  $\Omega_{\rm mb} = \theta_{\rm mb}^2 \pi / (4ln(2)) \approx 888 \, {\rm arcsec}^2$ , e.g. Kauffmann et al. 2008); (2) any source whose projected footprint size is less than two beam sizes; and (3) any source whose S/N ratio was less than 3.5 (estimated by taking the peak of emission and comparing it to the local noise level). Most spurious sources were successfully removed with this set of criteria, but some remained, in particular towards the noisier high-velocity end of the spectrum [at

<sup>5</sup>This size was determined by inspecting the datacubes. Unlike in the middle of the map where the noisy spikes are of the order of a beam size, the noisy spikes in the edges are typically much larger than a beam size due to gridding/convolution of the data whilst doing the data reduction. We also consider that even if some real sources were to be included in this criterion, those clouds would be both small and incomplete (since they touch an edge), and therefore their properties would be highly unusable.

positive velocities, which can be clearly seen on the  $\ell v$  plots of Figs 1 and A1 (available online)]. Therefore, we applied another criterion to our removal procedure: (4) any source outside the Galactic centre (GC) region (i.e. outside a +355 <  $\ell$  < 10° range), and with a centroid velocity  $v_{lsr} > 160 \, \mathrm{km \, s^{-1}}$ . The resulting catalogue contains 10 663 molecular clouds (whose masks are shown as colours in Fig. 1).

By comparing the integrated intensities inside the cloud masks with the total integrated intensities along each sightline, we estimate that the extracted clouds contain ~70 per cent of the total integrated flux above  $3\sigma_w$  (similar to Barnes et al. 2016), and ~50 per cent of the flux above  $2\sigma_w$ , where  $\sigma_w$  is the standard deviation of the total integrated intensity map, defined as  $\sigma_w = \sqrt{N_c}\sigma_{rms}\Delta v$ , with  $N_c$  being the total number of channels used for the integration,  $\sigma_{rms}$  the average noise level per channel (i.e. 0.7 K), and  $\Delta v$  the channel width (i.e. 0.25 km s<sup>-1</sup>). This suggests there is a non-negligible amount of molecular gas in a relatively diffuse inter-cloud medium. In addition, from the datacubes with the cloud masks, we find that of all  $\ell b$  pixels with clouds, we have ~82 per cent of sightlines with a single cloud assignment, meaning that only ~18 per cent of the lines of sight have multiple clouds (~16 per cent with two clouds, ~2 per cent with three clouds, and <1 per cent with more than three clouds).

# 3.2 Data products: Cloud masks and catalogues

From our SCIMES extraction, we have produced two main data products: a catalogue with the properties of all the molecular clouds; and the respective assignment datacubes in the same format as the input 3D datacubes of emission. These data products are made publicly available alongside the data release of the survey.<sup>6</sup>

In the assignment datacubes, each voxel holds the unique ID number of the cloud it has been assigned to by SCIMES, and the voxels with no assigned cloud take the value -1. These assignment datacubes are particularly useful for performing further studies on specific clouds, as they can be used to assign voxels to clouds, and therefore pull out the entire 3D structure of clouds from the original emission datacubes. Fig. 4 shows an example of the results from the cloud extraction towards a small portion of the survey, with the <sup>13</sup>CO peak intensity map in grey-scale, and the cloud masks overlaid as colours. In online Appendix A, we show the same images for the entire survey coverage (from online Fig. A1).

All the properties held in the catalogue of molecular clouds produced whilst running SCIMES are listed in online Table A1 (in online Appendix A). In essence, the catalogue contains two sets of properties: the directly measured quantities, and the physical properties derived from these after a distance has been assigned (see Section 4). Note that all the quantities we present in the catalogue were estimated using the default 'bijection' paradigm, which is the most appropriate for characterizing substructures within the nested dendrogram tree (Rosolowsky et al. 2008). Amongst the directly measured properties are the ID number, the cloud name, the clouds' centroid longitude  $(\ell)$ , latitude (b), and velocity (v), the velocity dispersion  $(\sigma_v)$ , the projected footprint area (*Area*) and the respective equivalent radius (*R*), the average integrated intensity ( $\langle I_{13}_{CO} \rangle$ ), and the peak intensity  $(T_{13}^{\text{peak}})$ . We also include a tag (edge) to indicate whether a cloud touches an edge of the survey coverage, in which case it is an incomplete object.

Given that some clouds will be close to the resolution element of our survey, a beam deconvolution on the sizes is needed. This will only affect the smaller objects, and has only very marginal effects on the statistical properties that we derive. Nevertheless, in the catalogue we also provide the equivalent radius deconvolved from the beam  $(R^d)$ .

In addition to the properties already described, we also estimated some basic parameters to characterize the clouds' morphology. First, we estimated the projected semimajor and semiminor axes from the second moment of the emission in 2D, weighted by the intensity (major and minor), along with the respective position angle (PA), and the aspect ratio ( $AR_{mom} = major/minor$ ). However, this moment method is relatively limited in providing a good approximation of a cloud's morphology, and can easily underestimate the true aspect ratio. Therefore, we also determined the projected geometrical medial axis of the clouds, which is the longest running spine along the 2D-projected cloud's mask, which is farthest away from the external edges (any internal holes in the cloud's masks are filled before determining the medial axis). From that, we include in the catalogue also the medial axis length ( $length_{MA}$ ), as well as the medial axis width as being twice the average distance to the cloud edge (width<sub>MA</sub>), and the corresponding aspect ratio  $(AR_{MA} = length_{MA}/width_{MA})$ . Fig. 5 shows an example of this medial axis for a cloud in our sample. Note that this is a purely geometrical medial axis (i.e. it is built on the assignment masks, with no information on the actual structure of the emission), and thus it is only a first approximation of the possible filamentary nature of clouds. A more accurate description of filamentary structures detected with ATLASGAL using the SEDIGISM survey data has been performed by Mattern et al. (2018), and shall be expanded to the entire SEDIGISM survey in future work.

The determination of the physical properties of the clouds requires a distance to be assigned. In Section 4, we detail the procedures that we followed to determine distances to the SEDIGISM clouds. Once the distances have been assigned, we can compute the physical properties of clouds. In the catalogue, besides the measured sizes in angular scales, we also present the sizes in physical scales, i.e. already converted using the assigned distance.

We then estimated a few other physical properties, which required using an  $X_{13}_{CO(2-1)}$  conversion factor between the integrated intensities of <sup>13</sup>CO (2-1) and the H<sub>2</sub> column densities. We adopted  $X_{^{13}\text{CO}(2-1)} = 1^{+1}_{-0.5} \times 10^{21} \text{ cm}^{-2} (\text{K km s}^{-1})^{-1}$ , as estimated in the SEDIGISM science demonstration field (Schuller et al. 2017), by comparing the SEDIGISM <sup>13</sup>CO emission to the H<sub>2</sub> column densities as derived from the Hi-GAL survey data (Molinari et al. 2010).7 With this  $X_{13}CO(2-1)$ , and assuming a mean molecular weight  $\mu_{H_2}$ of 2.8 (Kauffmann et al. 2008), we derived the clouds' masses (*M*), average gas surface densities ( $\Sigma$ ), and virial parameter ( $\alpha_{vir}$ ), defined as  $\alpha_{\rm vir} = 5\sigma_v^2 R/GM$  (Bertoldi & McKee 1992), where G is the gravitational constant,  $\sigma_v$  the velocity dispersion, and R is the equivalent radius. This formulation of  $\alpha_{vir}$  assumes a spherical geometry and a uniform density, and it only takes into account the balance between kinetic and gravitational energies. Thus,  $\alpha_{vir}$  is a very simplistic tool, and it should not be taken as a strict measurement of the gravitationally bound state of a cloud (e.g. Bertoldi & McKee 1992; Kauffmann, Pillai & Goldsmith 2013; Traficante et al. 2018a, b). However, given its wide usage in the literature, we estimate it here to allow a direct comparison of our results with those of other surveys.

<sup>&</sup>lt;sup>7</sup>The Hi-GAL column density maps for this calibration were built by fitting a pixel-by-pixel grey body curve to the spectral energy distribution from 160 to 500 µm (Elia et al. 2013), assuming a dust-to-gas ratio of 1:100, and an opacity law with a fixed spectral index  $\beta = 2$ , and  $\kappa_0 = 0.1 \text{ cm}^2 \text{ g}^{-1}$  at  $\nu_0 = 1200 \text{ GHz}$  (Hildebrand 1983).



Figure 4. Example of the SCIMES cloud extraction results, on a small section of the SEDIGISM survey. The top panel shows the  $\ell b$  map with the <sup>13</sup>CO peak intensity in grey-scale, and the SEDIGISM cloud masks overlaid as colours, where each cloud has a different (random) colour. The bottom panel shows the  $\ell v$  map of the same field, with the same colour-scheme as the top panel.

Finally, in the catalogue we provide the surface density and the virial parameters using both the measured *R* (noted as  $\Sigma$  and  $\alpha_{vir}$ ), and the deconvolved  $R^d$  (noted as  $\Sigma^d$  and  $\alpha_{vir}^d$ ). For the analysis presented in this paper, we will use the deconvolved properties, although this choice has only a very marginal effect on the respective distributions, keeping the global trends virtually unchanged. Given the uncertainties on the distance estimates (which are of the order of ~30 per cent) and on the  $X_{CO}$  factor (of a factor 2), all these quantities have an uncertainty of at least a factor 2.

In the catalogue, we also provide the Heliocentric and Galactocentric coordinates of each cloud, determined as explained in online Appendix B.

# **4 DISTANCE DETERMINATION**

In order to compute the physical properties of clouds, we require knowledge of the distances. However, for a large survey such as SEDIGISM, there are very few existing direct measurements of the distances towards molecular clouds, and we mostly need to rely on estimates based on the kinematic distances (i.e. by assuming a Galactic rotation model, see Section 4.1), which rarely give a unique answer. Therefore, it is often required to search for ancillary indications to narrow down the distance assignment. In the following sections, we describe the computation of the kinematic distances, and how the problem of the kinematic distance ambiguities (KDA)

was solved. The kinematic distance solutions, along with their uncertainties, and our final decisions are listed in the catalogue. For each cloud we include two distance tags:  $d_{sol}$  specifies the type of distance solution, and  $d_{flag}$  specifies the method used to reach the final distance assignment. The numbering of  $d_{flag}$  reflects the order by which we check the different methods. Once a cloud gets a distance as per a given tag, we stop testing further methods. The flowchart depicting this decision process is shown in Fig. 6. These methods are all described in detail in Section 4.2, and summarized in Table 1.

### 4.1 Kinematic distances

To derive the kinematic distances of the clouds in our catalogue, we have used the Galactic rotation model of Reid et al. (2016), which has been constructed using maser parallax distance measurements. This model uses the revised values for  $R_0$  and  $V_0$  of 8.34 kpc and 240 km s<sup>-1</sup>, respectively. Besides these rotation curve parameters, this model also uses a Bayesian approach that can consider the source's proximity to spiral arms, displacement from the Galactic mid-plane and proximity to parallax sources to estimate the most likely distance. Since molecular clouds are not always confined to the spiral arms or associated with star formation we have relaxed those constraints.

Some clouds, however, have velocities that lie outside those allowed by the rotation model, and thus we are unable to assign them a distance. This is the case for 363 clouds, and they can be identified in the catalogue with the distance solution tag  $d_{sol} = NULL$  (and  $d_{flag} = -1$ ).<sup>8</sup> For the remaining clouds, if they lie outside the Solar circle, there is a unique kinematic distance solution. This is the case for 551 clouds, and these can be identified in the catalogue with the tag  $d_{sol} = NA$ , standing for *No Ambiguity* (and  $d_{flag} = 1$ ). When sources are located within the solar circle, there are two possible distance solutions, a *near* and a *far* one, which are equally spaced on either side of the tangent distance. Clouds that lie close to the tangent velocities (i.e. within 5 km s<sup>-1</sup>, to accommodate for uncertainties due to streaming motions, e.g. Brand & Blitz 1993; Wienen et al. 2015) were assigned the tangent distance, and given the tag  $d_{sol} = T$  (and  $d_{flag} = 2$ ). This is the case for 1080 clouds.

For sources with two possible distances, we performed an extensive cross-match with literature information, checked directly for H I self-absorption (HISA) in each cloud, and checked whether the cloud properties would make them statistically more likely to be at a specific distance solution, in order to solve the distance ambiguity. Upon completion of this procedure, clouds that were assigned a near distance were tagged in the catalogue with  $d_{sol} = N$  (corresponding to a total of 3679 clouds), while far distance clouds have  $d_{sol} = F$  (which amount to 4979 clouds). The full details on the procedure leading to our final distance decision are described in the following section.

Note that, despite our extensive effort in assigning distances to clouds, there are regions within our Galaxy for which we know that our kinematic distances are not reliable. We have therefore included a flag in the catalogue,  $d_{\text{reliable}}$ , which identifies clouds for which the distances are unreliable or non-existent ( $d_{\text{reliable}} = 0$ ) and those that have a reliable distance estimate ( $d_{\text{reliable}} = 1$ ). In particular, we have given a  $d_{\text{reliable}} = 0$  for clouds with a  $|v_{\text{lsr}}| < 10 \text{ km s}^{-1}$  with a near distance assignment. For those, the kinematic distance is too uncertain, since the  $v_{\text{lsr}}$  of the clouds is dominated by local motions,



**Figure 5.** Example of the medial axis for a molecular cloud in our sample (SDG316.766–0.020), which corresponds to an IRDC (SDC316.786–0.044 from Peretto & Fuller 2009), and the larger cloud often shortened to G316.75 (e.g. studied in Watkins et al. 2019). The grey-scale shows the <sup>13</sup>CO (2–1) integrated intensity, estimated using the voxels within the cloud's mask as defined by SCIMES, and the coloured pixels show the geometrical medial axis, colour-coded with the distance to the external cloud edge.

and therefore the distance assigned from a global rotation model has a distance uncertainty of the order of the distance value itself. We also assigned a  $d_{\text{reliable}} = 0$  to clouds for which we were not able to solve the distance ambiguity (i.e. clouds with a  $d_{\text{flag}} = 12$ , see Section 4.2.4). In addition, clouds towards the GC (and including most of the Galactic bar), i.e. within  $+353^{\circ} < \ell < 7^{\circ}$ , also have a very uncertain distance estimate (and are given a  $d_{\text{reliable}} = 0$ ), as the Galactic rotation model used for our kinematic distance assignment is not tailored to reproduce the complex dynamics of the gas in the centre of the Galaxy. The only exception being the clouds for which we have a maser parallax distance (as that is an exact measurement, independent of kinematic considerations), which are retained with a  $d_{\text{reliable}} = 1$ . The GC will be studied in more detail in future work, and we will then revise the catalogued distances for those clouds accordingly.

# 4.2 Solving the distance ambiguities

# 4.2.1 Maser parallaxes, dark clouds, IRDCs, and HiSA from literature

We performed a cross-match of our entire catalogue with literature information for any known robust indication of the distance of our clouds. We started by cross-matching our clouds with a compilation of known maser parallax measurements (Reid et al. 2009, 2014; Honma et al. 2012; Bobylev & Bajkova 2013; Wu et al. 2014). The matches were performed by checking if the position of the masers (in 3D) fell inside the mask of one of our SEDIGISM clouds. Sources with a known maser parallax measurement were assigned their maser parallax distance (instead of the kinematic distance). If there were more than one maser parallax measurement for a given cloud, then we take the average parallax distance. Clouds with a maser distance were given a  $d_{sol} = M$  and  $d_{flag} = 0$ , and this was the case for 11 clouds. The small number of SEDIGISM clouds with a maser parallax is due

<sup>&</sup>lt;sup>8</sup>Note that in the catalogue we assign these clouds a distance of -1, which effectively means that we have estimated their physical properties as if they were at 1 kpc distance, and that properties that have a linear dependence with distance will appear with negative sign.

 Table 1. Summary of the methods used to determine the distances of clouds, along with the number of clouds that had their distances assigned with each method.

$d_{\mathrm{flag}}$	Description	Number of clouds
-1	No distance information	363
0	Exact maser parallax distance	11
1	No distance ambiguity	551
2	Tangent distance	1080
3	Dark cloud (near distance)	77
4	IRDC (near distance)	751
5	Literature HISA (near distance)	91
6	Direct HISA measurement (near distance)	828
7	ATLASGAL source at near distance	252
8	Solomon distance to GP (near distance)	34
9	Size-linewidth scatter (near or far distance)	2263
10	ATLASGAL source at far distance	142
11	Extinction (near or far distance)	3178
12	Ambiguity not solved (defaulted to far)	1042

Cloud's (l,b,v)



**Figure 6.** Flowchart showing the distance assignment procedure adopted for the SEDIGISM clouds. The blue boxes highlight the methods used, and the grey boxes show the corresponding assigned distance and tag. The green and red arrows show the directions taken if a specific method succeeds or fails in providing a distance solution, respectively.

to the fact that most of the maser parallax catalogues cover Quadrants 1, 2, or 3, hence only very few maser parallax distances have been measured for sources in our longitude range, and of those, about half lie outside our latitude range.

We then did a cross-match with other literature catalogues [including dark clouds, infrared dark clouds (IRDCs), and HISA], using the clouds' centroid Galactic coordinates and velocity. For catalogues in which the major axes, minor axes and position angles are given, the match was done by checking if the centroid position of the SEDIGISM cloud falls in the elliptical footprint of the catalogued source. For catalogues that give no position angle, or provide only the beam size or a radius, we use the effective radius and the match is done by checking if the centroid of the SEDIGISM cloud falls

Using these criteria, we cross-matched our clouds with catalogued dark clouds with velocity information (Otrupcek, Hartley & Wang 2000), as well as with IRDCs, some with and some without velocity information (Simon et al. 2006; Du & Yang 2008; Jackson et al. 2008; Peretto & Fuller 2009; Chira et al. 2013; Liu; Wang & Xu 2013). Their extinction makes dark clouds and IRDCs appear in silhouette against a bright background (in the visible and in the IR, respectively). dark clouds typically reach high optical depths very quickly, and thus are typically tracing nearby clouds that absorb the stellar light from the Galactic disc. The IRDCs probe a higher column density regime, which means we observe deeper into the molecular clouds. Nevertheless, the concept is the same, in that we are more likely to see a cloud in extinction, if there is enough IR background to absorb against, thus placing such clouds preferably at their near distance solution (although this might not always be the case, e.g. Giannetti et al. 2015, found ~10 per cent of IRDCs to be located at the far distance). For our purpose, we have assumed any SEDIGISM sources which have a dark cloud, or an IRDC match to be at the near distance, and given a  $d_{\text{flag}} = 3$  or 4, respectively. Note that, in cases where the cross-match with IRDCs was only spatial (i.e. in the absence of available velocity information), we only consider the match to be reliable if there is a single SEDIGISM cloud associated with each IRDC: if an IRDC is in the same line of sight as multiple SEDIGISM clouds, more information - such as velocity information or a more detailed morphological match - would be needed in order to produce a robust association.

We also cross-matched our catalogue with known HISA (or H<sub>2</sub>CO absorption) features from the literature within the SEDIGISM coverage (Sewilo et al. 2004; Busfield et al. 2006; Pandian, Momjian & Goldsmith 2008; Anderson & Bania 2009; Urquhart et al. 2012; Anderson et al. 2015; Wienen et al. 2015). HISA occurs when cold HI gas in the foreground absorbs the warmer HI emission from background gas at the same velocity (e.g. Gibson et al. 2000). Therefore, the existence of HISA at a given velocity is often used as an indication that the cold gas that is absorbing is at a near distance<sup>9</sup> – as this makes it more likely to have background emission to absorb against, and that emission is less likely to be filled by other warmer HI emission along the line of sight between the observer and the cold cloud (e.g. Roman-Duval et al. 2009). SEDIGISM sources with a known HISA feature from the literature were assumed to be at the near distances, and given a  $d_{flag} = 5$ .

### 4.2.2 Direct and automated HISA determination

Given that many of the clouds in our catalogue do not have a counterpart with literature sources (given the improved sensitivity

<sup>9</sup>Note that not all near-distance clouds are expected to show a strong HISA feature, given the simultaneous requirement of (1) the existence of significant cold HI gas at the same velocities as the molecular cloud traced by <sup>13</sup>CO; (2) the existence of warm HI gas in the background, at the same velocity as the cloud; and (3) the non-existence of intervening warm HI gas between us and the cloud, that could fill in the cloud's intrinsic HISA. In addition, HII regions can also produce a direct HI absorption feature, even at the far distances, which can be mistakenly interpreted as an HISA feature. The HISA method for solving the KDA is therefore estimated to be ~80 per cent reliable (e.g. Anderson & Bania 2009).



**Figure 7.** Three sketch examples of the automated HISA method, showing the original HI spectrum on the first column (in orange), with the cloud's velocity ranges denoted in blue, and the linear background fit done to the HI spectrum in black dotted line. The second column shows the cloud's background-subtracted HI spectra (in green), with the dotted dashed line representing the 0-emission level, and the vertical bar representing  $3\sigma_{\rm rms}$  of the HI emission. The two last columns show the criteria that we use to infer whether there is HISA in that particular sightline. Case 1 represents a line of sight with strong HISA, but the other two cases are not considered to have HISA. Case 3 shows an example of a false positive arising from criterion (i) alone, but which is mitigated by introducing criterion (ii).

and resolution of the SEDIGISM survey), we have also checked for the presence of HISA directly for each individual cloud. We have done so in an automated way, making use of HI 21 cm ATCA and Parkes data from both the Southern Galactic Plane Survey (McClure-Griffiths et al. 2005), and the ATCA HI Galactic Centre Survey (McClure-Griffiths et al. 2012). Both these data sets have a spatial resolution of 2 arcsec, a spectral resolution of 1 km s<sup>-1</sup> and an average noise level,  $\sigma_{\rm rms}^{\rm HI}$ , of ~1 K. We combined the data from these surveys into a single datacube, using the CONVERT<sup>10</sup> and KAPPA<sup>11</sup> packages from the STARLINK software (Currie et al. 2014), namely the WCSALIGN and WCSMOSAIC procedures. We then extracted sub-cubes covering the same spatial and velocity range of each of the SEDIGISM tiles, which we reprojected and resampled to the same pixel and channel sizes as the SEDIGISM data. Even though this procedure heavily oversamples the HI data, it facilitates the automated check of the HI emission in each of the SEDIGISM clouds, by directly using the assignment masks produced by SCIMES.

Our automated HISA procedure works in the following way: First, it selects all the voxels that belong to each cloud, and creates a projected 2D image of the cloud in the plane of the sky. This allows us to identify all lines of sight which belong to that cloud. Then, for each sightline, it determines the 'background' HI emission, by taking the HI emission one channel before, and one channel after the cloud's velocity range in that specific sightline, and fitting it with a linear function (see illustration on the first column of Fig. 7). We then subtract the HI emission inside the cloud from the fit of the background HI emission, on a channel by channel basis. This 'subtracted' datacube should have negative values whenever the HI content of a cloud is self-absorbing against the background HI emission (see the second column of Fig. 7). Therefore, we use these background-removed HI datacubes as the basis for our decision on whether a given sightline in a cloud has a significant HISA or not (see last two columns of Fig. 7).

To determine if a specific line of sight has HISA, we impose two conditions:

(i) The minimum intensity of the background-removed H I emission signal is lower than  $-3\sigma_{\rm rms}^{\rm HI}$ . This ensures that the self-absorption is significant, given the noise in the H I data.

(ii) The sum of the background-removed H I signal is negative, and has an absolute value larger than three times the cumulative noise, given by  $3\sqrt{N}\sigma_{\rm rms}^{\rm HI}$ , where N is the number of velocity channels across which the signal was summed up.

Step (ii) ensures that false positives are rejected. A false positive typically occurs when our simple background fit does not capture properly the variations of the HI background emission (e.g. by under- or overestimating the slope of the HI background emission), producing a signature similar to a p-Cygni profile, whose dip may be deeper than the HI noise – thus passing our criteria (i) (see Case 3 of Fig. 7). However, while a true self-absorbed profile would have negative emission throughout the entire cloud velocity range, resulting in the sum of the background-removed HI emission to be also negative (and significant), a false positive would have a sum that is within the noise of the HI data. We therefore use this criterion to remove potential false positives.

We then consider a cloud to have strong HISA only if the number of sightlines (i.e. 2D pixels) that satisfy condition (i) amount to at least one beam size in the HI data, and that satisfy condition (ii) amount to at least one SEDIGISM beam size. The results from this automated HISA determination are compiled in the catalogue under the *tag\_hisa* property, which is assigned a value of 1 for strong HISA, 0 if it is ambiguous (i.e. meeting only some of the criteria above), and -1 if there is no HISA. Clouds with a strong HISA from this method are taken to be at a near distance, and given a  $d_{flag} = 6$ .

### 4.2.3 ATLASGAL distances

The ATLASGAL survey (Schuller et al. 2009; Beuther et al. 2011) observed the dust continuum emission towards the inner Galactic

<sup>&</sup>lt;sup>10</sup>http://starlink.eao.hawaii.edu/docs/sun55.htx/sun55.html <sup>11</sup>http://starlink.eao.hawaii.edu/docs/sun95.htx/sun95.html

plane at 870  $\mu$ m, and produced a catalogue of 10163 compact sources<sup>12</sup> (CSC catalogue; Contreras et al. 2013; Urquhart et al. 2014c). In order to determine the distances to these clumps, there was a significant effort in assigning velocities to the continuum emission through a combination of extensive cross-match with molecular line data reported in the literature and dedicated follow-up observations (Wienen et al. 2012, Csengeri et al. 2016, Wienen et al. 2018, Urquhart et al. 2019). This was then combined with the Reid et al. (2016) Galactic rotation curve to calculate kinematic distances, and the distance ambiguities were resolved using the HISA method and using a friends-of-friends clustering algorithm to identify complexes. This successfully determined the distances to ~8000 ATLASGAL clumps (see Urquhart et al. 2018 for details).

Since all of the SEDIGISM survey is covered by ATLASGAL, we performed a cross-match between all clouds in our sample, to the ATLASGAL clumps with known  $v_{lsr}$  from Urquhart et al. (2018). This cross-match was done by considering the centroid positions and velocities of the ATLASGAL clumps, and placing them in the respective voxel in our 3D datacubes. We then checked whether that voxel falls within the mask of a SEDIGISM cloud (i.e. a perfect match), and otherwise estimate the distance to the nearest SEDIGISM cloud (in all three dimensions). We then consider ATLASGAL clumps that lie within one beam size of the edge of the nearest cloud, or within one  $\sigma_v$  of the cloud, to be a *partial match*. Out of the 5067 ATLASGAL sources within the SEDIGISM coverage, 4376 were matched as a perfect match to a SEDIGISM cloud, and 448 as being a partial match, leaving only 243 ATLASGAL clumps without a SEDIGISM counterpart. Most of these unmatched ATLASGAL clumps are either small clumps whose corresponding SEDIGISM emission did not satisfy our minimum size requirement, or they are in regions that form part of a smoother background that does not get assigned to a specific cloud (i.e. where the <sup>13</sup>CO emission does not have a local peak rising above the  $4\sigma_{\rm rms}$  requirement to be considered as independent peaks/leaves within the dendrogram). In total, these 4824 ATLASGAL clumps are contained within 1709 SEDIGISM clouds (i.e. ~16 per cent of SEDIGISM clouds).

Given that the distances to the ATLASGAL sample were estimated with the individual  $v_{lsr}$  of clumps (rather than that of the parent cloud), we do not use their distances directly. Instead, we are only interested in the type of distance solution determined for each ATLASGAL clump (near or far), in order to incorporate it in our distance assignment. In most cases, all ATLASGAL clumps within a given SEDIGISM cloud have a distance solution that agrees amongst them. However, there are a few cases where, within a SEDIGISM cloud, there are ATLASGAL clumps with both a near and far solution. In those cases, we define the 'global' ATLASGAL solution as being near, under the assumption that an indication for a near distance solution is more reliable than the absence of one (which is the most common reason for a far distance assignment). Note that, even though we had to do this step to provide a complete list of 'ATLASGAL distance solutions' for our SEDIGISM sample, none of the clouds for which the ATLASGAL distance solutions disagreed, actually took their final solution from ATLASGAL (instead they had their KDA lifted by other methods).

For SEDIGISM clouds with an ATLASGAL match, and for which the criteria in Sections 4.2.1 and 4.2.2 did not have an indication for a near distance, we check the distance solution from ATLASGAL. If that solution is *near*, then we adopt the near distance, and assign

### 4.2.4 Other distance indicators

In addition to the above methods, we also checked two often-used techniques that take the statistical distribution of the properties of molecular clouds into account. The first one is the method used by Solomon et al. (1987), which considers the physical distance of a cloud to the Galactic plane, should the cloud be assigned the far distance. If by taking the far distance the cloud is too far off the Galactic plane (i.e. >140 pc, which is the scale height of the gaseous Galactic disc, e.g. Solomon et al. 1987), then the near distance is favoured, and the cloud is given a  $d_{\text{flag}} = 8$ . Note that towards the far side of the Sagittarius and Scutum-Centaurus arm (around  $\ell \sim 290^{\circ}$ ), the Galactic mid-plane is known to be warped towards negative latitudes (e.g. Chen et al. 2019; Romero-Gómez et al. 2019). This implies that on the far-distance side, in the latitude range of  $300^{\circ} < \ell < 318^{\circ}$ , the Galactic plane descends below a latitude of  $-0.5^{\circ}$  (e.g. Reid et al. 2016), and therefore this area of the Galaxy is not well covered by our survey (since we cover a relatively narrow b range). Nevertheless, since the Galactic warp only becomes significant at Galactocentric distances of 8 kpc and beyond, any clouds in the longitude range of  $300^{\circ} < \ell < 318^{\circ}$  possibly following the warp are beyond the Solar circle, and should have unambiguous distances. Therefore, our criterion checking for the height above the Galactic plane is not affected by the existence of the Galactic warp.

The second method places each cloud on the size–linewidth relation ( $\sigma_v$ –R) (e.g. Larson 1981; Solomon et al. 1987), for both near and far distance solutions, and checks which solution provides the smaller distance to the empirical relation. We use this method to favour a given distance solution *only* if one solution is significantly closer to the empirical relation than the other solution (i.e. at least a factor 3 difference in log-space). More details on this method can be found in online Appendix D, and clouds that used this criteria were given a  $d_{\text{flag}} = 9$ .

Finally, we also used a method based on datacubes of the visual extinction in *K* band, as a function of distance (Marshall et al. 2006, in preparation, Elia et al., in preparation). From those cubes, the structures of significant extinction can be identified along each line of sight, by taking the distances at which the extinction has a significant jump. We then compare the extinction distances with the near and far kinematic distance estimates, by taking into account a 30 per cent uncertainty on the extinction distance as well as the kinematic distance uncertainties. We solve the KDA by taking the kinematic distance which has an extinction counterpart, if one exists. Clouds that used this criteria were given a  $d_{\text{flag}} = 11.^{13}$  In the future, this could potentially be expanded to also include *Gaia*-based 3D dust extinction maps (e.g. Lallement et al. 2019), although at the moment these only probe distances up to 3 kpc.

### 4.3 Revisiting previous distance estimates

In order to gauge how our distance estimates compare to the results from other surveys that covered the same area of the Galactic plane,

<sup>&</sup>lt;sup>12</sup>Note on nomenclature: we will refer to the ATLASGAL compact sources as 'clumps', as opposed to the larger scale SEDIGISM structures that we refer to as 'clouds'.

<sup>&</sup>lt;sup>13</sup>This method does have a few limitations, one being that it becomes less reliable for far distances, mainly as the extinction cubes have a pixel of 5 arcmin, and therefore roughly 10 times larger than the SEDIGISM beam size. Small clouds assigned a distance using this method should therefore be used with caution.

Table 2.	Summary	of different	samples.
----------	---------	--------------	----------

Sample name	Description/conditions for selection	Number of sources
Full sample	Entire catalogue, with distances $(d_{\text{flag}} \neq -1)$	10 300
Science sample	$d_{\text{reliable}} = 1, Area > 3\Omega_{\text{beam}}, \text{edge} = 0$	6664
Distance limited sample	$d_{\text{reliable}} = 1, Area > 3\Omega_{\text{beam}}, \text{edge} = 0, 2.5 \text{ kpc} < d < 5 \text{ kpc}$	1743
Complete science sample	$d_{\text{reliable}} = 1, Area > 3\Omega_{\text{beam}}, \text{edge} = 0, M > 2.6 \times 10^3 \text{M}_{\odot}, R > 2.9 \text{pc}, d < 14.5 \text{kpc}, d_{\text{flag}} \neq 2$	1680

we have compared the results from our distance solutions to those of the ATLASGAL survey, as a reference (given that ATLASGAL had already performed a detailed comparison with other surveys, e.g. Urquhart et al. 2014b, 2018). Note however, that as per Section 4.2.3, only 1709 SEDIGISM clouds have an ATLASGAL counterpart (i.e. only  $\sim 16$  per cent of our sample), although this includes 95 per cent of all ATLASGAL clumps in our coverage (i.e. 4814 clumps). Of those, 1253 ATLASGAL clumps did not have an assigned distance, which we have now assigned.<sup>14</sup> For the ATLASGAL sources with a distance, the KDA solution between the two surveys agrees for 3080 ATLASGAL clumps. This leaves a total of 481 clumps (i.e. 13.5 percent of the ATLASGAL clumps with distances) with a distance solution that was revised by us, in most cases from a far distance to a near distance, by one of the other methods listed in Section 4.1. Most of these revisions were done using our HISA method (321 clumps,  $d_{\text{flag}} = 6$ ), followed by 60 clumps revised using the IRDC matches ( $d_{\text{flag}} = 4$ ), and 53 clumps with literature HISA ( $d_{\text{flag}} = 5$ ). A further 23 clumps were revised using the maser parallax measurements ( $d_{\text{flag}} = 0$ ), 7 clumps using the distance around the Larson relation ( $d_{\text{flag}} = 9$ ), and 2 clumps using the dark cloud association ( $d_{\text{flag}} = 3$ ). Finally, 2 clumps were re-assigned a near distance for being in the same complex as other ATLASGAL sources with a near distance ( $d_{\text{flag}} = 7$ ), 1 clump was revised as having a nonambiguous solution ( $d_{\text{flag}} = 1$ ), and 12 clumps had their distances revised to a tangent distance ( $d_{\text{flag}} = 2$ ), although for these cases the change from near or far solutions into the assumed tangent distance is within the uncertainties.

With the large survey coverage, and improved resolution and sensitivity of the SEDIGISM survey compared to other spectroscopic surveys covering the same Galactic longitudes (e.g. the MopraCO survey, Burton et al. 2013; the ThrUMMS, Barnes et al. 2015; the Dame et al. 2001 survey), here we present the most extensive sample of molecular clouds towards the inner Galaxy yet, with 10663 clouds in total. With our comprehensive effort to combine different independent methods to determine the distance solutions for each SEDIGISM cloud, we have been able to assign distances to 10 300 clouds, 7993 of which have well-characterized (reliable) distance assignments.

# 5 GLOBAL PROPERTIES OF THE SEDIGISM SAMPLE

For our analysis of the statistical properties of the SEDIGISM molecular clouds, we have excluded any clouds whose projected footprint size is smaller than 3 beams (i.e. any clouds that are barely resolved). We also excluded clouds with an unreliable distance  $(d_{\text{reliable}} = 0)$ , and those that are incomplete because they touch a

survey coverage edge (edge = 1). With these criteria, we select a total of 6664 clouds for our analysis, which we will refer to as our 'science sample'. In addition, we will refer to the science sample above the completeness limits (as per online Appendix C) as our 'complete science sample' (which also exclude clouds at a tangent distance – see Section 6 for more details). Table 2 summarizes the specific details of the several samples that we use in the paper.

### 5.1 Distribution of individual properties

Fig. 8 shows the distributions of a number of different properties, namely the total mass (M), the velocity dispersion ( $\sigma_n$ ), the medial axis length ( $length_{MA}$ ), the average surface density ( $\Sigma$ ), the virial parameter ( $\alpha_{vir}$ ), and the aspect ratio from the medial axis ( $AR_{MA}$ ). The histograms correspond to the full science sample (in light grey), from which we highlight the subset of clouds with an ATLASGAL counterpart (in dark grey), and from those, also clouds with a signpost of high-mass star formation (HMSF, in red), as per Urquhart et al. (2014b). These signposts of HMSF include the existence of methanol masers (Urguhart et al. 2013a, 2015, which used the masers from the Methanol Multibeam Survey, Caswell et al. 2010; Green et al. 2012); HII regions (Urguhart et al. 2013b, which combined information from the CORNISH survey, Hoare et al. 2012; Purcell et al. 2013, and the GLIMPSE survey, Benjamin et al. 2003); or massive young stellar objects, YSOs (Urquhart et al. 2014b, which matched ATLASGAL sources with YSOs and HII regions identified by the Red MSX Source survey, Lumsden et al. 2013; Urguhart et al. 2014a). In total, we have 435 SEDIGISM clouds within the full sample (330 in the science sample, i.e.  $\sim$ 4 per cent of clouds) that have signposts of active HMSF (similar to the fraction of high-mass star-forming clouds found by Barnes et al. 2011). We note, however, that for this work, we did not cross-match our SEDIGISM clouds with HMSF tracers directly: our sample of HMSF clouds is purely a subsample of the ATLASGAL sources, and so any HMSF signposts outside that are not accounted for. This will be explored in future work. We also computed the main statistics (i.e. the median, lower and upper quartiles, skewness, and kurtosis) of these distributions, plus that of the equivalent radius (R), which we compile in Table 3. These distributions, however, could potentially be affected by our different completeness at different distances within our science sample. In order to check how this might affect the global results, we have also computed the histograms using a distance limited sample (with 2.5 kpc < d <5.0 kpc), shown in online Appendix E (online Fig. E1). The statistics for the distance limited sample are also compiled in Table 3, showing that they follow broadly the same trends as the science sample.

Noticeably, the median values in Table 3 and the histograms from Fig. 8 show that clouds with an ATLASGAL counterpart tend to be at the higher end of the distributions of mass, velocity dispersion, size, aspect ratio, and surface density, as compared to the science sample. This is even more so for clouds with an HMSF signpost (whose median values are again higher than those of the ATLASGAL subsample). The increase in the median values of those properties as we go from the science sample to the HMSF sub-sample range from a

<sup>&</sup>lt;sup>14</sup>Note that most of these are towards the central Galaxy, for which the kinematic distances are less reliable. If we consider only the sample we use for science (as per Section 5), the number of ATLASGAL clumps that so far did not have a distance assigned and for which we are able to assign a reliable distance is 308.



**Figure 8.** Histograms of global properties: Mass (top-left), velocity dispersion (top-centre), medial axis length (top-right), average surface density (bottom-left), virial parameter (bottom-centre), and aspect ratio from the medial axis (bottom-right). The histograms are for the science sample (light grey), clouds that have an ATLASGAL counterpart (dark grey), and clouds that have an HMSF signpost (red). The normalization of all histograms was made with respect to the total number of clouds in the science sample. The vertical dashed line on the mass histogram shows our mass completeness limit (see online Appendix C), and the dashed lines on the virial parameter histogram represent an  $\alpha_{vir} = 1$  and 2.

modest increase of a factor 2 (e.g. for the aspect ratio and velocity dispersion) up to an order of magnitude increase (for the mass). The only exception to this trend is the virial parameter, for which the median values (and the quartiles) are similar between all three subsets.

Interestingly, while the science sample typically has a distribution with a significant tail (i.e. with high kurtosis values), as we move from the full sample to the ATLASGAL sub-sample and then to clouds with an HMSF signpost, the shape of the distribution of all properties (except for the aspect ratio) becomes progressively flatter (smaller values of kurtosis) and symmetric (smaller values of skewness) with HMSF clouds occupying nearly the same parameter space as clouds with an ATLASGAL counterpart but without HMSF signpost. This is rather interesting as it suggests that there is no single 'global' property of clouds that is sufficient to determine, on its own and unambiguously, their potential to host high-mass star formation, and perhaps a complex combination of several conditions is needed. It is worth noting that some global properties like magnetic fields are, of course, not considered here. In Section 6.5, we will investigate if the ability to form high-mass stars might instead be influenced by the Galactic environment.

### 5.2 Scaling relations

Fig. 9 shows two of the most common scaling relations in the literature: the size–linewidth relation in the top panels, where the dashed line represents the Larson relation,  $\sigma_v^2 \propto R$  (Larson 1981; Heyer et al. 1998); and the Heyer relation,  $\sigma_v^2/R \propto \Sigma$  (Heyer et al. 2009), in the lower panels, where the solid black line shows  $\alpha_{vir} = 1$  as defined in Section 3.2, and the dashed lines correspond to a  $\alpha_{vir} =$ 

1 when including the contribution of external pressure ( $P_{ext} = 1, 10$  and 100  $M_{\odot}$  pc<sup>-3</sup> km<sup>2</sup> s<sup>-2</sup>). On the left-hand panels, we show our SEDIGISM science sample in grey-scale, and the subset of clouds with an ATLASGAL counterpart in blue, and those with a signpost of HMSF in red. From these, we can see that although our SEDIGISM clouds do show some correlation on both plots, neither of these follow the scaling relations proposed by previous works.

The right-hand side panels show a compilation of literature catalogues of molecular clouds in green colour scale, including both Galactic studies (Oka et al. 2001; Heyer et al. 2009; Roman-Duval et al. 2010; Barnes et al. 2016; Rice et al. 2016; Miville-Deschênes et al. 2017; Colombo et al. 2019; Rigby et al. 2019) and extragalactic studies (Rosolowsky & Blitz 2005; Bolatto et al. 2008; Wong et al. 2011; Gratier et al. 2012; Wei, Keto & Ho 2012; Donovan Meyer et al. 2013; Colombo et al. 2014; Leroy et al. 2015; Utomo et al. 2015; Faesi, Lada & Forbrich 2016; Freeman et al. 2017; Pan & Kuno 2017; Schruba et al. 2017; Tosaki et al. 2017). On those, we overplot the loci of the distribution of our science sample as the black ellipse, produced from a principal component analysis<sup>15</sup> (PCA; Pearson 1901), similar to Colombo et al. (2019). The ellipse

<sup>15</sup>The PCA analysis (Pearson 1901) can be useful to identify the directions of maximal and minimal variance of data with large intrinsic scatter, thus equivalent to finding the direction and scatter of the underlying scaling relation (which are typically estimated using a linear regression fit). As we are simply interested in using the PCA as a representation of the loci of the distributions, we did not take into account the uncertainties in the measured quantities for this analysis.

**Table 3.** Statistics of some of the physical properties of the SEDIGISM clouds, namely the mass (M), velocity dispersion ( $\sigma_v$ ), equivalent radius ( $R_{eq}$ ), medial axis length (*length*<sub>MA</sub>), medial axis aspect ratio ( $AR_{MA}$ ), surface density ( $\Sigma$ ), and virial parameter ( $\alpha_{vir}$ ), for the entire science sample, and for a distance-limited sample (to minimize distance-biased results). Within these samples we also list the statistics for the subsets of clouds with an ATLASGAL counterpart or with a HMSF signpost. Q25 and Q75 represent the lower (25 per cent) and upper (75 per cent) quartiles of the distributions.

	Science sample					Distance limited sample (2.5 kpc $< d < 5.0$ kpc)				
Sub-set	Median	Q25	Q75	Skewness	Kurtosis	Median	Q25	Q75	Skewness	Kurtosis
$M (\times 10^3 \text{ M}_{\odot})$										
Science	1.25	0.40	3.59	52.7	3521.1	0.43	0.13	2.04	7.6	94.7
With ATLASGAL	5.13	1.69	13.83	24.1	690.6	3.74	1.20	10.52	4.3	32.4
With HMSF	12.00	3.56	27.24	13.8	220.9	10.31	3.35	23.08	3.1	17.0
$\sigma_v (\mathrm{kms^{-1}})$										
Science	0.76	0.55	1.07	6.9	162.0	0.77	0.51	1.25	2.6	15.5
With ATLASGAL	1.29	0.97	1.81	7.5	126.7	1.35	0.99	1.93	2.5	14.1
With HMSF	1.66	1.25	2.20	7.7	95.4	1.69	1.28	2.29	2.4	11.6
$R_{\rm eq}$ (pc)										
Science	2.31	1.39	3.64	3.4	37.2	1.34	0.80	2.62	1.9	7.2
With ATLASGAL	3.56	2.16	5.65	1.5	6.8	3.07	1.78	4.53	1.0	3.8
With HMSF	4.82	2.79	6.92	1.4	6.0	4.10	2.66	5.81	0.6	3.1
$length_{MA}$ (pc)										
Science	7.51	4.21	13.51	3.8	42.3	4.82	2.54	10.83	2.1	8.8
With ATLASGAL	13.52	7.73	23.04	1.6	6.4	12.61	6.32	21.07	1.2	4.5
With HMSF	18.88	10.73	29.74	1.2	4.5	16.95	10.82	26.28	0.8	3.4
AR <sub>MA</sub>										
Science	4.9	3.4	7.1	1.6	7.7	5.6	3.8	8.3	1.6	7.0
With ATLASGAL	6.6	4.5	9.6	1.4	6.3	7.6	5.1	10.9	1.4	6.0
With HMSF	7.6	5.3	10.8	1.6	7.5	9.0	6.2	11.7	1.6	7.7
$\Sigma (M_{\odot} pc^{-2})$										
Science	73.0	58.0	99.7	5.1	70.7	75.4	57.6	112.7	3.0	19.3
With ATLASGAL	128.2	98.3	170.2	4.2	42.3	139.9	103.9	190.8	2.2	12.1
With HMSF	158.1	120.4	221.1	3.7	27.6	186.0	137.1	252.9	1.9	8.5
$\alpha_{\rm vir}$										
Science	1.25	0.79	2.10	8.4	128.2	1.85	1.23	3.08	4.3	33.0
With ATLASGAL	1.36	0.81	2.58	7.5	82.5	1.79	1.05	2.98	2.9	16.0
With HMSF	1.28	0.76	2.62	7.0	61.3	1.49	0.94	2.78	2.5	10.3

contours in the right-hand panels of Fig. 9 correspond to a  $2\sigma$  level, i.e. it contains  $\sim$  95 per cent of the data points, while the central point corresponds to the mean. The remaining 5 per cent of data points are overplotted as circles.

We have also performed this PCA analysis for the cloud catalogues from the fiducial sample of the COHRS survey [in <sup>12</sup>CO (3-2), Colombo et al. 2019], and from the CHIMPS survey [in <sup>13</sup>CO (3-2), Rigby et al. 2019], which we plot in Fig. 9 as yellow and purple ellipses, respectively. Although both of these surveys have a slightly higher spatial resolution than SEDIGISM (17 arcsec versus 28 arcsec), they both cover the first quadrant, making them highly complementary to the SEDIGISM survey. In fact, the native resolution of CHIMPS was smoothed to 27 arcsec for their source extraction and derivation of cloud properties that we use here, thus making it very similar to that of the SEDIGISM survey. For completeness, we summarize the directions of major variance from the PCA analysis for these three surveys in Table 4, which can be compared to the expected slopes from the literature. Note, however, that even though the slopes from the PCA analysis can be suggestive of a correlation, in all the cases we performed the PCA here, the major and minor axis are similar (within a maximum of a factor 3 difference), which indicates that these are not tight correlations.

The clouds from the COHRS survey were extracted using the same method as us (SCIMES) but, because it uses  $^{12}\mathrm{CO}$  (3–2), it

typically traces larger clouds, with larger velocity dispersions (partly due to the fact that <sup>12</sup>CO traces more diffuse gas than <sup>13</sup>CO, but also due to line broadening from optical depth effects, and from a coarser spectral resolution of  $1 \text{ km s}^{-1}$ ).<sup>16</sup> The CHIMPS survey coverage overlaps with COHRS, but it uses the optically thinner <sup>13</sup>CO (3–2). Even though the clouds from CHIMPS were extracted using Fellwalker (Berry 2015), which segments the emission into their individual peaks (hence not allowing for the grouping of several peaks into complexes), and their line tracer is not the same as ours (using a higher energy transition of <sup>13</sup>CO), the properties of the CHIMPS clouds agree remarkably well with those of our SEDIGISM sample. There is only a small shift in the sizes of the SEDIGISM clouds towards larger values (as we can see in the topright panel of Fig. 9) and, although the CHIMPS sample spans to lower average surface densities than the SEDIGISM sample (as we can see on the lower-right panel), both samples have clouds reaching similar values towards the high surface density end. These differences can be easily understood as a consequence of (1) the cloud

<sup>&</sup>lt;sup>16</sup>A comprehensive comparison of the COHRS cloud population with other surveys can be found in Colombo et al. (2019), namely their fig. 13, which can be used to compare with the relative position of the SEDIGISM cloud catalogue.



**Figure 9.** Top row: size–linewidth relation ( $\sigma_v$  versus  $R_{eq}$ ), where the dashed line represents the Larson relation. Bottom row: scaling relation between  $\sigma_v^2/R$  and gas surface density  $\Sigma$ , where the lines correspond to  $\alpha_{vir} = 1$ : the solid line is without external pressure, and the dashed lines are when including external pressure (from top down, at a constant  $P_{ext} = 100$ , 10, and 1 M<sub>☉</sub> pc<sup>-3</sup> km<sup>2</sup> s<sup>-2</sup>). The left-hand panels show these relations for the SEDIGISM sample alone, where the grey-scale represents the density of points for the entire science sample, the blue circles show the clouds with an ATLASGAL counterpart, and the red circles show the clouds that have an HMSF signpost. The panels on the right show, in green, the density of points from a compilation of literature catalogues which include both Galactic and extragalactic studies (see the text for full list of references). Our SEDIGISM sample is represented by the black ellipse (from a PCA analysis, and where the ellipse contour contains 95 per cent of the data) and black points (which show the remaining 5 per cent of clouds). Similarly, we also show the PCA ellipses for the fiducial sample of the COHRS survey in orange (Colombo et al. 2019), and the CHIMPS survey in purple (Rigby et al. 2019), both of which are high-resolution surveys towards the first Galactic quadrant – complementary to SEDIGISM. For reference, the dashed grey boxes on the right-hand panels show the plotting range of the corresponding left-hand panel.

segmentation used by the CHIMPS survey, breaks up the emission more, thus extracting smaller (and less dense) clouds, whilst the grouping of individual clumps into larger cloud complexes achieved by our usage of SCIMES for the SEDIGISM segmentation will tend to incorporate such small diffuse clumps into larger complexes and (2) the <sup>13</sup>CO (3–2) transition used in CHIMPS has a higher critical density and will typically trace warmer gas than the brighter <sup>13</sup>CO (2–1) transition of SEDIGISM, which will mean that the CHIMPS clouds will typically be able to trace less mass for a given brightness temperature.

Most interestingly, these plots show that the choice of tracer and the specific limitations of the surveys change our global view of the properties of molecular clouds. Looking at the <sup>12</sup>CO emission from the COHRS survey, we could argue that these clouds are in a pressureconfined regime (i.e. lying above the  $\alpha_{vir} = 1$  line when external pressure is not included, but could be consistent with being virialized if a moderate external pressure is at play). However, looking at the same clouds with an optically thinner tracer (i.e. with CHIMPS) changes our perception of their energy balance, with clouds moving closer to a more gravitationally bound regime, or else requiring only

**Table 4.** Slopes ( $\alpha$  and b) recovered from a PCA analysis on the scaling relations, where  $\sigma_v \propto R^{\alpha}$ , and  $(\sigma_v^2/R) \propto \Sigma^b$ . The mean values of each pair or quantities (i.e. the centres of the ellipses in Fig. 9) are noted with the supperscript *m*.

Sample	α	$[\sigma_v^m, R^m]$	b	$[(\sigma_v^2/R)^m,\Sigma^m]$
SEDIGISM	0.52	[2.19, 0.79]	3.91	[78.3, 0.29]
CHIMPS	0.39	[1.43, 0.86]	2.15	[48.3, 0.52]
COHRS	0.27	[4.48, 2.33]	14.79	[79.1, 1.22]
Expected	0.5 <sup><i>a</i></sup>		$1.0^{b}$	

<sup>a</sup>Larson (1981).

<sup>b</sup>Heyer et al. (2009).

a very weak external pressure to be virialized. This points out a rather important issue: although molecular clouds are highly hierarchical, they are part of a continuous medium that smoothly blends into the diffuse warm neutral medium, with no hard boundary. We know that the ISM is not composed of a discrete set of entities, and yet this discretization is (and has been) a crucial step in our understanding of the cold molecular medium. What we use to define them thus changes what we actually trace. Simple measures of the energy balance of clouds at any one single level are incapable of providing a complete picture of the true physics that describe and regulate the evolution of clouds. Instead, we need to move into trying to put a sequence together for the general trend of the change in molecular cloud properties with tracer density (which could even perhaps be used as a proxy for time). Studies looking into the evolution of these global properties, within molecular clouds - i.e. as we move inside the internal hierarchy of clouds - are necessary for taking our understanding of the physics inside molecular clouds to the next level. This is one of the key advantages of using a dendrogram-based segmentation of the ISM that we shall explore in future work.

# 6 GALACTIC DISTRIBUTION OF THE MOST EXTREME CLOUDS

Using the longitude  $(\ell)$  and distance (d) of the clouds in our catalogue, we can estimate their Galactocentric coordinates, which we use to plot our clouds on a 'top-down' view of the Galaxy. These are shown in Fig. 10, overlaid on an artist's impression of the Milky Way [by NASA/JPL-Caltech/R. Hurt (SSC/Caltech)]. The main known gaseous spiral arms are labelled in the bottom-left panel. The top-left panel of Fig. 10 shows our full SEDIGISM catalogue with distances, the top-right panel shows the distribution of our science sample, and the bottom-right panel shows the science sample colour-coded depending on whether the clouds have an ATLASGAL counterpart (blue), or an HMSF signpost (red). Using this top-down Galactic distribution of clouds in the science sample, we estimate a typical mass surface density of gas associated with clouds to be of the order of 1  $\times$   $10^5\,M_{\odot}\,kpc^{-2}$  (and ranging from  ${\sim}4.4$   $\times$   $10^2$  to  $1.3 \times 10^6 \,\mathrm{M_{\odot} \, kpc^{-2}}$ ). Note that the values for the average and minimum mass surface densities are only lower limits, as they are likely affected by our completeness limits. On the bottom-left panel of Fig. 10, we show our complete science sample, i.e. clouds within the science sample that lie above our mass and radius completeness limit (as detailed in online Appendix C), and are located within a Heliocentric distance of 14.5 kpc (the distance used to determine our completeness limit). The complete science sample also excludes clouds with a tangent distance. For those clouds, although the physical properties are reliable (since the near and far distances are relatively close together), their Galactic position falls into a single line at the

tangent distance, which introduces some biases for the statistical tests we will be performing with this sample (see online Appendix F for more details). Our complete science sample consists of 1680 clouds.

We caution that showing clouds with this top-down perspective, although suggestive, can be misleading - indeed we know that the uncertainties on the distances can amount to  $\sim 1 \text{ kpc}$ , particularly when streaming motions around spiral arms can be important, and this can easily displace clouds across entire spiral arms. In addition, the exact position and strength of these arms is still quite uncertain (e.g. Taylor & Cordes 1993; Reid et al. 2014; Vallée 2017). In fact, the very existence of four strong spiral arms is still subject of debate, especially as studies in the optical/near-IR (e.g. Drimmel 2000; Siebert et al. 2011, 2012; Gaia Collaboration 2018), suggest that we only have two main stellar spiral arms - which could indicate that the four spiral arms that we see in the gas, are not as well defined as this figure depicts, and are perhaps more flocculent in nature. This idea is also supported by our relatively low values of molecular gas mass surface densities, which place the Milky Way at the bottom of the distribution of the values retrieved for a sample of 15 nearby spiral galaxies Sun et al. (2018), whose typical molecular gas mass surface densities are of the order of  $10^6 - 10^8 \text{ M}_{\odot} \text{ kpc}^{-2}$ . Hence, these top-down perspective plots are used here merely as a first look at the Galactic distribution of clouds. A more detailed study of arm/interarm dependence requires using a model of the spiral pattern, and is most accurately done in the  $\ell bv$  space, which is beyond the scope of this paper.

In order to look for effects that could depend on the Galactic environment, without the need to assume any specific spiral arm model, we have examined the spatial distribution of clouds with extreme properties (i.e. clouds that form the tails of a distribution). and compared those to the global Galactic distribution of clouds. The idea behind this exercise is a purely statistical one, which will test whether the most extreme clouds follow the same spatial distribution as the global population of clouds, or whether they show significant deviations from it. As an attempt to take this analysis a step further, we can make the loose assumption that the spiral arms should preferentially be represented by the crowded regions of the global population, while the inter-arm regions would be preferentially associated with the least crowded places. This assumption is purely qualitative (due to the uncertainties in the distances), and we make no attempt to effectively associate clouds with spiral arms or interarm regions. For our purpose, we use the complete science sample as our global cloud population (bottom-left panel of Fig. 10), from which we selected a number of sub-samples that comprise the most extreme clouds. This selection was made by taking the most extreme 100 clouds of each distribution (corresponding to the top or bottom 6 per cent), and the specific selection criterion is indicated at the top of each panel in Figs 11 and 12.

The comparison between the sub-samples and the global cloud population was done by performing the Pearson's  $\chi^2$  statistical test, which tests whether the frequency distribution of certain events observed in a sample is consistent with a particular theoretical distribution. The full details of the  $\chi^2$  statistical test that we performed are explained in online Appendix F. In brief, for our purpose, we used the 2D Galactic distribution of clouds in the complete science sample as our theoretical distribution. In practice, we built a normalized 2D histogram with the spatial distribution of clouds in the complete science sample (using their Galactocentric coordinates), using a spatial bin of  $0.3 \times 0.3$  kpc – this map represents the probability of an observation falling in a specific spatial bin (see left-hand panels of online Figs F1–F3). We then compute the  $\chi^2$  statistics using the observed 2D distribution of each sub-sample (shown in the



**Figure 10.** Top down view of the Galaxy, with the deprojected position of SEDIGISM clouds overplotted on an artistic impression of the Milky Way [NASA/JPL-Caltech/R. Hurt (SSC/Caltech)]. The position of the Sun is marked with a '+' in all panels. The top-left panel shows the density plot of the entire catalogue, and the top-right panel shows the science sample. The bottom-left panel shows the Galactic distribution of the clouds in the complete science sample (i.e. above our completeness limit, and excluding clouds with a tangent distance assignment). For these three panels, the colour scale and the size of the symbols is related to the local density of clouds (more crowded areas are shown in white, and with larger symbols). The bottom-right panel shows all the sources in the science sample colour-coded depending on whether they have an ATLASGAL counterpart (in blue), an HMSF signpost (in red), or neither (in black).

central panels of online Figs F1–F3), and the observed  $\chi^2$ -values are compared to the values obtained from a pure random draw of clouds from the theoretical distribution (i.e. effectively obtaining a *p*-value, which we call  $p_{\rm rnd}$ , see online Fig. F4). Given the statistical fluctuations, as well as the uncertainties in the distributions, and binning effects (neither of which are taken into account for this exercise), the exact  $\chi^2$  values and  $p_{\rm rnd}$  that we derive should not be taken at face value. Instead, they are more useful for a relative comparison of the sub-samples, as an indication for which subsamples are most different to the global cloud population. The results from our  $\chi^2$  statistical test are summarized in online Table F1. We describe all of the studied tails of distributions in Sections 6.1–6.5.

### 6.1 The most massive molecular cloud complexes

Some observations of nearby spiral galaxies (e.g. Koda et al. 2009, Colombo et al. 2014), as well as some galaxy-scale numerical models (e.g. Dobbs et al. 2008, Fujimoto et al. 2014, Duarte-Cabral & Dobbs 2016, Pettitt et al. 2018) – both of which benefit from a more straightforward association of clouds to spiral arms – have suggested



Figure 11. Top down view of the Galaxy as in Fig. 10, showing the SEDIGISM clouds of the complete science sample that are part of the top 10 per cent of clouds in terms of Mass (top-left), surface density (top-right), medial axis length (bottom-left), and aspect ratio from the medial axis (bottom-right). The specific condition that corresponds to this cut-off is indicated at the top of each panel. Clouds are colour-coded depending on whether they have an ATLASGAL counterpart (in blue), an HMSF signpost (in red), or neither (in black).

that the most massive clouds are preferentially located along spiral arms. This is widely accepted and understood in the context of spiral arms being able to concentrate more material, and thus able to form larger and more massive giant molecular clouds. Similarly to the argument for encountering the most massive clouds in the arms, with a higher concentration of material in the spiral arms we ought to expect the highest surface density clouds to lie in the spiral arms as well. In this spirit, we have plotted the Galactic distribution of the 100 most massive clouds in our sample, and the 100 clouds with the highest surface density, in the top-left and top-right panels of Fig. 11, respectively.

Our  $\chi^2$  tests comparing these two distributions to the global cloud population, give us  $\chi^2$  values of 670 and 638 (which corresponds to a  $p_{rnd}$  of 0.05 and 0.16), for the extreme mass and surface

density clouds respectively. This suggests that the distribution of high surface density clouds still follows the original distribution of clouds, implying that such clouds might be found in crowded areas (or spiral arms), simply from statistics. The distribution of the most massive clouds, however, is less consistent with a pure random draw of clouds from the parent distribution. If the disparities between the two distributions were caused by having more high-mass clouds in the spiral arms than what is statistically expected, then we should see an excess of high-mass clouds in the most crowded areas of the global distribution. However, considering the spatial distribution of these clouds on Fig. 11 (top panels), and the relative difference between the predicted and measured counts shown in online Fig. F1 (top and middle rows), it is not obvious that this is the case, with clouds having both an excess and lack of counts in different crowded



**Figure 12.** Same as Fig. 11, showing the SEDIGISM clouds of the complete science sample that are part of the top 10 per cent of clouds with a high virial parameter (top-left), and high velocity dispersion (top-right). The lower panels show the bottom 10 per cent of clouds in the same properties: with a low virial parameter (bottom-left) and a low velocity dispersion (bottom-right). The specific condition that corresponds to this cut-off is indicated at the top of each panel. Clouds are colour-coded depending on whether they have an ATLASGAL counterpart (in blue), an HMSF signpost (in red), or neither (black).

areas. The specific regions where the most high-mass clouds are found to be in excess or lacking, are not particularly striking in terms of their environment, leaving our interpretation inconclusive.

### 6.2 The most elongated clouds

A subject of increasing interest in the SF community is the origin and properties of the most elongated clouds. While some numerical and observational studies suggest that extremely long filamentary clouds would be formed as the result of the Galactic shear in the inter-arm regions (e.g. Kim & Ostriker 2002; Shetty & Ostriker 2006; Ragan et al. 2014; Duarte-Cabral & Dobbs 2016, 2017), other studies suggest that at least some of these might trace the 'spines' of the spiral arms (e.g. Goodman et al. 2014; Wang et al. 2015; Zucker, Battersby & Goodman 2015).

We have thus looked at the Galactic distribution of the 100 longest clouds in the SEDIGISM sample, as well as the 100 clouds with the largest aspect ratio. These are shown in the bottom panels of Fig. 11 (left and right, respectively). Our  $\chi^2$  tests for these two distributions, give us a  $\chi^2$  value of 715 and 671 (corresponding to a  $p_{rnd}$  of 0.005 and 0.04) for the extreme length and aspect ratio clouds respectively. This suggests that the Galactic distribution of both these sub-samples are different from the global cloud population (although this is most evident for the sample of largest clouds). However, neither of them seem to show any clear preference for crowded or non-crowded areas (see also online Fig. F1 bottom panel and online Fig. F2 top panel).

This analysis has a few caveats, though. The first one is that there are more of these elongated clouds located at the near distance, than there are at the far distance. This could be linked to resolution limitations which will result in more distant filaments appearing less elongated. The second caveat is the fact that both of these quantities are purely the projected ones (the length and aspect ratio on the plane of the sky). If long filamentary clouds are indeed shaped by the shear from the Galactic differential rotation, we do not expect them to be randomly orientated. Therefore, this projection is likely to affect our ability to select the truly elongated structures, in specific parts of the Galaxy, being particularly critical in lines of sight where we expect the clouds' elongations to be roughly along our line of sight. The third caveat is the fact that even the longest molecular filaments in our Galaxy (such as the  $\sim 100 \text{ pc}$  long Nessie filament, Jackson et al. 2010), do not appear in our segmentation as a single entity - they are instead composed of several (smaller) filamentary sections. Finally, the relative lack of large (and massive) clouds nearby (with d < 2.5 kpc), can also point at a possible bias from the cloud segmentation, in which we might still be more likely to break the most nearby clouds into smaller substructures (even though we use a clustering algorithm designed to minimize this effect). All of these effects could result in the underestimation of both the length and aspect ratio of clouds, and thus the 100 most extreme clouds we take for this analysis might not correspond to the most extreme cases in physical space.

### 6.3 The most dynamically active clouds

Given the complex global Galactic dynamics, we would expect to see, at least at first order, some link between the most dynamic places in the Galaxy, with the kinetic properties of clouds. In this sense, we have isolated the 100 clouds with the highest virial parameter, and highest velocity dispersion. Their Galactic distribution is shown in Fig. 12, top panels. Our  $\chi^2$  test for clouds with a large virial parameter gives us a  $\chi^2$  value of 699 (corresponding to a  $p_{rnd}$  of 0.01), indicating that they differ from a random statistical subset of the global cloud population, in terms of their Galactic placement (see also online Fig. F2 middle panel). On the other hand, clouds with a large velocity dispersion have a higher  $\chi^2$  value of 747 (which corresponds to a much smaller  $p_{rnd}$  of 0.001), making this distribution less like the global cloud population. Most of the differences in the statistics of this sub-sample comes from an excess of high-velocity dispersion clouds relatively nearby (see top-right panel of Fig. 12, and bottom panel of online Fig. F2), which then also propagates (although less severely) into clouds with high virial parameters also being mostly nearby. We believe that these trends could be partly due to observational biases (see online Fig. C2, and the discussion in online Appendix C).

Interestingly, these dynamically active clouds typically make up two types of populations. The first is most closely associated with crowded regions (potentially associated with the near Sagittarius, Scutum and Norma spiral arms), which is where we expect more frequent cloud–cloud interactions, in line with the results from numerical simulations of spiral galaxies (e.g. Duarte-Cabral & Dobbs 2017; Pettitt et al. 2018). This population of clouds is also actively forming high-mass stars. The larger values of velocity dispersion and virial parameters could thus be also an indication of larger internal motions of clouds, perhaps partly driven by their active gravitational contraction, or by internal feedback from the forming stars, or both.

The second population of clouds are devoid of HMSF signposts, and some even lacking an ATLASGAL counterpart (i.e. less dense). Most of these are also at large distances, which could suffer from a completeness effect in the ATLASGAL and HMSF tracers. Alternatively, this second population could represent clouds relatively close to the Galactic bar, and/or in the streams of gas feeding the GC region – all regions prone to experiencing a significant shear driven by the global Galactic dynamics. This dichotomy (of clouds in the two extremes of their SF history sharing the same integrated dynamical properties) highlights the caveats of performing a standard virial analysis and deriving any conclusions therefrom alone.

#### 6.4 The most dynamically quiescent clouds

On the opposite extreme of the dynamical status of molecular clouds, we have also explored the location of the clouds that are relatively quiet (which we refer to as the most 'dynamically quiescent' clouds), which include clouds with a low virial parameter, or a low velocity dispersion. These types of clouds are often not subject of much attention (mostly as they typically lie close to survey limitations in terms of spectral resolution). Nevertheless, some recent numerical work by Pettitt et al. (2018) has suggested that, in grand-design spiral galaxies, while clouds with high virial parameter are most often associated with spiral arms, clouds with low virial parameters have a weaker correspondence with the spiral arms, with many inter-arm clouds being remnants of large arm complexes or simply formed *in situ* from small overdensities in filaments and arm spurs.

To investigate these dynamically quiescent clouds in SEDIGISM, we have selected the 100 clouds in the complete science sample with the lowest virial parameter, and the lowest velocity dispersion. Their Galactic distribution is shown in the bottom panels of Fig. 12. Our  $\chi^2$  tests give us  $\chi^2$  values of 675 and 669 (corresponding to a  $p_{rnd}$  of 0.04 and 0.05) for the low virial parameter and low velocity dispersion respectively. This suggests that the Galactic distribution of the most dynamically quiescent clouds is only mildly different to that of the global cloud population. Their distribution in Fig. 12 (see also online Fig. F3 top and middle row) suggests that they are not found in very crowded areas (possibly favouring inter-arm locations).

Clouds with a low virial parameter are often interpreted to be gravitationally bound (i.e. where gravity dominates over turbulence). However, these clouds are not necessarily collapsing – indeed if they were, the collapse itself would increase the virial parameter again (e.g. Kauffmann et al. 2013). Our results show that these dynamically quiescent clouds are mostly devoid of HMSF or even high column densities (which would result in an ATLASGAL counterpart), perhaps indicating that their evolution is not regulated by their own gravity but by interaction with the Galactic potential, the large-scale shear motions and perhaps also by large-scale magnetic fields.

We caution however, that even though a handful of dynamically quiescent clouds are relatively nearby, most of them are at d >8.0 kpc. In terms of absolute numbers, the science sample does contain nearby low velocity dispersion clouds, but most of those are below the size and/or mass threshold used to build the complete science sample. The usage of a completeness limit for the whole SEDIGISM sample (and especially one largely above the resolution element) was an attempt to remove any bias from the resolution and distance. However, our intrinsic observational limitations may still be responsible for at least part of this signature, as we can see that the average measured velocity dispersion of the complete science sample has a correlation with distance (see online Fig. C2, and the respective discussion in online Appendix C). Furthermore, at the far distances our sample may also not be complete in terms of the detection of an ATLASGAL counterpart or HMSF signposts, potentially biasing the interpretation above.

Nevertheless, these type of clouds could potentially be interesting to follow up with the goal to investigate whether this tentative trend does hold up, with a more in-depth analysis, considering the survey limitations and a detailed modelling of the spiral pattern of the Galaxy.

#### 6.5 The high-mass star-forming clouds

One of the questions we wanted to address here is whether the Galactic distribution of clouds that host ongoing high-mass star formation is uniform, or whether they are preferably located in spiral arms as our preliminary study of the SEDIGISM science verification field suggested (Schuller et al. 2017). In particular, if high-mass star-forming clouds are tracing the arms, we are also interested in exploring whether that is purely due to a statistical sampling (as suggested by e.g. Elmegreen & Elmegreen 1986; Moore et al. 2012; Eden et al. 2013); or whether there is an excess of high-mass star-forming regions in the crowded spiral arms, suggestive of SF triggering from the passage of a spiral wave (e.g. Lin & Shu 1964; Roberts 1969; Toomre 1977; Martínez-García, González-Lópezlira & Bruzual-A 2009).

Fig. 10 (bottom-right panel), shows the distribution of all clouds with an HMSF signpost in our science sample (in red). The  $\chi^2$ statistical test, performed using only the clouds in the complete science sample (from which only 211 clouds have an HMSF signpost) gives a  $\chi^2$  value of 735, which translates into a  $p_{rnd}$  of 0.001. This indicates that the distribution of clouds with an HMSF signpost does not mimic the global distribution of clouds. Upon closer inspection of Figs 10 and F3 (available online), it becomes clear, however, that most of the deviations from the global distribution of clouds do not arise from crowded or non-crowded areas, but rather shows a distance effect. Indeed, most of the clouds with signs of on-going high-mass star formation are located relatively close to us. The extremely high density of points there (compared to elsewhere in the Galaxy), is likely to be a simple consequence of completeness in the HMSF signposts (namely HII regions and massive YSOs).

Interestingly, if we look at the higher mass clouds or the higher surface density clouds (Fig. 11 top panels), not all of these host highmass star formation. This is true even if we just consider the most nearby clouds, where we should be less affected by completeness issues in terms of HMSF signposts. As we have seen in Section 5, there does not seem to be a unique global property of a molecular cloud that defines the ability of a cloud to form high-mass stars - and the same applies for the Galactic environment. Perhaps to isolate clouds with a potential to form massive stars, we need to use a combination of conditions that need to be satisfied, or even just the most extreme conditions within a cloud (rather than the integrated properties). Applying a single global threshold law (such as a gas surface density threshold or mass-radius threshold, e.g. Krumholz & McKee 2008; Kauffmann & Pillai 2010; Baldeschi et al. 2017) to define the potential to form massive stars, is probably not a single unique descriptor. Fig. 13 highlights this issue, where we can see clouds with and without high-mass star formation that have the same mass and radius. In this figure, we also show as a dashed line, the empirical relation for high-mass star formation inferred by Kauffmann & Pillai (2010), and confirmed by other works (e.g. Kauffmann et al. 2010a, b; Urguhart et al. 2018). Note that the plotted line is the Kauffmann & Pillai (2010) original threshold scaled up so



**Figure 13.** Mass–radius relation for the SEDIGISM clouds, where the greyscale represents the density of points for the entire science sample, the blue circles show the clouds with an ATLASGAL counterpart, and the red circles show the clouds that have an HMSF signpost. The dashed line shows the empirical relation from Kauffmann & Pillai (2010), where clouds above this line are expected to form high-mass stars. The plotted threshold is at  $M[M\odot] = 1053$  (R[pc])<sup>1.33</sup>, which is scaled up from the original HMSF threshold from Kauffmann & Pillai (2010), to account for the different opacity law used (see online Appendix G).

as to be consistent with our adopted opacity law (see online Appendix G for more details).

Although this empirical relation was determined for clumps, rather than for clouds (as we use here), the bulk of the parameter space that we probe is similar to that in Kauffmann & Pillai (2010): their sizes range from <0.1 to 10 pc (compared to our range of 0.3 to ~30 pc), and their masses range from 1 to >  $10^4 M_{\odot}$  (compared to our range of 10 to > $10^5 M_{\odot}$ ).

If we use that relation directly with our clouds, we would miss some true positives (107 out of 330 clouds with an HMSF signpost lie below the empirical threshold, i.e. missing  $\sim$ 33 per cent of all clouds that we know are actively forming massive stars), as well as potentially provide a significant number of false positives (455 out of a total of 678 clouds above the HMSF threshold do not have a detected HMSF signpost, i.e.  $\sim$ 70 per cent of clouds above the empirical line for HMSF). Since completeness limits could play a role in the non-detection of the signposts for HMSF, we estimate that the number of false negatives (i.e. missed true positives) is a lower limit, while the number of false positives is an upper limit.

The detection of potential false positives was not ruled out by Kauffmann & Pillai (2010): indeed they note that their threshold appears to capture a necessary condition for HMSF, but not a sufficient one. Alternatively, it could also be that part of the clouds above the HMSF threshold line but for which we have no detected HMSF (i.e. the false positives), are in fact clouds that simply have not done so yet, because of the potential large latency periods prior to star formation. In that sense, the trends in properties going from the science sample to clouds with an ATLASGAL counterpart and then clouds with an HMSF signpost (from Section 5) could be an indication of the cloud evolution towards HMSF during this latency period (with clouds progressively building up their mass, becoming larger, denser, and more dynamically active – with larger

velocity dispersions), even if this remains a stochastic process for each individual cloud (e.g. Barnes et al. 2018).

More intriguing, however, are the missed true positives. These clouds lie below the empirical line supposedly representing the threshold below which HMSF would not occur, and yet they have tracers of ongoing HMSF. Nevertheless, it is possible that the material probed by Kauffmann & Pillai (2010) is intrinsically tracing higher density material than what we do, which could shift the exact position of the cloud sample with respect to the empirical line for HMSF, thus potentially making this relation inappropriate for usage with our sample. An indication that this might indeed be the case, is the fact that the subsample of SEDIGISM clouds with an HMSF signpost that we present here, is purely a subsample of the ATLASGAL clumps, which seem to confirm the Kauffmann & Pillai (2010) relation on clump scales (e.g. Urquhart et al. 2018). This highlights a potential caveat of using such relations blindly, as perhaps they are not applicable on cloud scales, when the density profiles become shallower, and the more diffuse material contributes to increasing the sizes of the clouds, whilst providing only moderate increase to the enclosed mass. A hierarchical study of this transition within clouds would be required to understand where this relation might break.

### 7 SUMMARY AND CONCLUSIONS

The SEDIGISM survey has covered  $\sim$ 84 deg<sup>2</sup> of the inner Galaxy with <sup>13</sup>CO (2–1). From the contiguous portion of the survey (i.e. excluding the W43 field), we extracted the entire molecular cloud population with a large dynamic range in spatial scales, using the SCIMES algorithm. We determined the distances to the clouds, using the kinematic distances, and a number of methods to solve the distance ambiguities (including masers, IRDC, dark clouds, HISA, distance to the Larson's size–linewidth relation, distance to the Galactic plane, and extinction distances). The full catalogue that we release contains 10 663 molecular clouds, 10 300 of which with measurements of physical properties.

In this paper, we have explored some of the global properties of clouds using a sub-sample of the full catalogue (i.e. our 'science sample'), consisting of 6664 well-resolved sources and for which the distance estimates are reliable. In particular, we compare the scaling relations retrieved from SEDIGISM to those of other surveys, including Galactic and extragalactic work. We find that the locus of the SEDIGISM clouds is similar to that of other surveys, but that the specific scaling relations vary widely between surveys – even between those that cover the same area in the Galaxy, just with different tracers. The intrinsic scatter in these relations is very large, making all the correlations rather unconstrained.

We also explored the properties of clouds with and without tracers of high-mass star formation, and we find that for most distributions (mass, size, surface density, velocity dispersion), the median values of the distributions is higher for clouds with an HMSF signpost, potentially indicative of an evolutionary sequence. However, the distributions become progressively flatter, with the clouds with HMSF spanning a wide range of values for all properties we looked at. These results suggest that there is no single global property of a cloud that is able to define their ability to form massive stars, and the usage of a simple threshold to isolate clouds forming highmass stars is not complete (providing both false negatives and false positives).

Finally, we have looked into potential links between the Galactic environment of clouds and their properties, by looking at the Galactic distribution of the most extreme clouds. For that purpose, we have isolated the most extreme 100 clouds in each distribution (i.e. clouds that make up the tails of the distributions), and compared their Galactic distribution to that of the cloud population above our completeness limits (i.e. our complete science sample), using a  $\chi^2$ statistical test. This provides a means to determine whether extreme clouds follow a Galactic distribution that differs significantly from the global cloud population. We find that, for most properties, the Galactic distribution of the most extreme molecular clouds is is only marginally different to that of the global cloud population. The Galactic distribution of the largest clouds, the most turbulent clouds and the high-mass star-forming clouds are those that deviate most significantly from the global cloud population. We also find that the least dynamically active clouds (with low velocity dispersion or low virial parameter) are situated further afield, mostly in the least populated areas, and therefore could hint at those being mostly in inter-arm regions. However, we find that part of these trends might be due to completeness limits (e.g. in case of the HMSF tracers), and intrinsic survey limitations, which result in a trend of decreasing velocity dispersion with distance, hampering our ability to make any firm conclusions from this data alone.

In future work, we shall follow up some of these tentative trends using distance-limited samples, with the incorporation of detailed models of the spiral arms, and with more complete cross-match with signposts of HMSF (e.g. by comparing with the Hi-GAL samples, and their L/M ratio as an indicator for more embedded HMSF and their respective evolutionary stage) to mitigate some of the observational biases that are potentially at play in the work presented here.

# ACKNOWLEDGEMENTS

ADC acknowledges the support from the Royal Society, through a University Research Fellowship (URF/R1/191609). ADC and AJR acknowledge the support from the UK STFC consolidated grant ST/N000706/1. DC acknowledges support by the Deutsche Forschungsgemeinschaft, DFG, through project number SFB956C. LB and RF acknowledge support from CONICYT grant Basal AFB-170002. HB acknowledges support from the European Research Council under the Horizon 2020 Framework Program via the ERC consolidator grant CSF-648505. HB furthermore thanks for financial help from the DFG via the SFB881 'The Milky Way System' (subproject B1). CLD acknowledges funding from the European Research Council for the FP7 ERC consolidator grant project ICYBOB, grant number 818940. MW acknowledges funding from the European Union's Horizon 2020 research and innovation programme under the Marie Sklodowska-Curie grant agreement no. 796461. SB and NS acknowledge support from the Agence National de Recherche (ANR/France) and the Deutsche Forschungsgemeinschaft (DFG/Germany) through the project GENESIS (ANR-16-CE92-0035-01/DFG1591/2-1). The STARLINK software (Currie et al. 2014) is currently supported by the East Asian Observatory. This publication is based on data acquired with the Atacama Pathfinder Experiment (APEX) under programmes 092.F-9315 and 193.C-0584. APEX is a collaboration among the Max-Planck-Institut fur Radioastronomie, the European Southern Observatory, and the Onsala Space Observatory.

# DATA AVAILABILITY

With this paper, we release the complete catalogue of SEDIGISM molecular clouds as extracted using the SCIMES code, alongside the masks of each molecular cloud in the catalogue as fits files, and additional ancillary tables in http://sedigism.mpifr-bonn.mpg.de.

# REFERENCES

- Anderson L. D., Bania T. M., 2009, ApJ, 690, 706
- Anderson L. D., Armentrout W. P., Johnstone B. M., Bania T. M., Balser D. S., Wenger T. V., Cunningham V., 2015, ApJS, 221, 26
- Baldeschi A. et al., 2017, MNRAS, 466, 3682
- Barnes P. J. et al., 2011, ApJS, 196, 12
- Barnes P. J., Muller E., Indermuehle B., O'Dougherty S. N., Lowe V., Cunningham M., Hernandez A. K., Fuller G. A., 2015, ApJ, 812, 6
- Barnes P. J., Hernandez A. K., O'Dougherty S. N., Schap William J. I., Muller E., 2016, ApJ, 831, 67
- Barnes P. J., Hernandez A. K., Muller E., Pitts R. L., 2018, ApJ, 866, 19
- Battersby C. et al., 2011, A&A, 535, A128
- Benjamin R. A. et al., 2003, PASP, 115, 953
- Berry D. S., 2015, Astron. Comput., 10, 22
- Bertoldi F., McKee C. F., 1992, ApJ, 395, 140
- Beuther H., Kainulainen J., Henning T., Plume R., Heitsch F., 2011, A&A, 533, A17
- Bobylev V. V., Bajkova A. T., 2013, Astron. Lett., 39, 809
- Bolatto A. D., Leroy A. K., Rosolowsky E., Walter F., Blitz L., 2008, ApJ, 686, 948
- Brand J., Blitz L., 1993, A&A, 275, 67
- Burton M. G. et al., 2013, Publ. Astron. Soc. Aust., 30, e044
- Busfield A. L., Purcell C. R., Hoare M. G., Lumsden S. L., Moore T. J. T., Oudmaijer R. D., 2006, MNRAS, 366, 1096
- Caswell J. L. et al., 2010, MNRAS, 404, 1029
- Chen X., Wang S., Deng L., de Grijs R., Liu C., Tian H., 2019, Nat. Astron., 3, 320
- Chira R.-A., Beuther H., Linz H., Schuller F., Walmsley C. M., Menten K. M., Bronfman L., 2013, A&A, 552, A40
- Colombo D. et al., 2014, ApJ, 784, 3
- Colombo D., Rosolowsky E., Ginsburg A., Duarte-Cabral A., Hughes A., 2015, MNRAS, 454, 2067
- Colombo D. et al., 2019, MNRAS, 483, 4291
- Contreras Y. et al., 2013, A&A, 549, A45
- Csengeri T. et al., 2016, A&A, 586, A149
- Currie M. J., Berry D. S., Jenness T., Gibb A. G., Bell G. S., Draper P. W., 2014, in Manset N., Forshay P., eds, ASP Conf. Ser. Vol. 485, Astronomical Data Analysis Software and Systems XXIII. Astron. Soc. Pac., San Francisco, p. 391
- Dame T. M., Hartmann D., Thaddeus P., 2001, ApJ, 547, 792
- Dempsey J. T., Thomas H. S., Currie M. J., 2013, ApJS, 209, 8
- Dobbs C. L., Glover S. C. O., Clark P. C., Klessen R. S., 2008, MNRAS, 389, 1097
- Donovan Meyer J. et al., 2013, ApJ, 772, 107
- Drimmel R., 2000, A&A, 358, L13
- Du F., Yang J., 2008, ApJ, 686, 384
- Duarte-Cabral A., Dobbs C. L., 2016, MNRAS, 458, 3667
- Duarte-Cabral A., Dobbs C. L., 2017, MNRAS, 470, 4261
- Duarte-Cabral A., Acreman D. M., Dobbs C. L., Mottram J. C., Gibson S. J., Brunt C. M., Douglas K. A., 2015, MNRAS, 447, 2144
- Eden D. J., Moore T. J. T., Morgan L. K., Thompson M. A., Urquhart J. S., 2013, MNRAS, 431, 1587
- Elia D. et al., 2013, ApJ, 772, 45
- Ellsworth-Bowers T. P. et al., 2013, ApJ, 770, 39
- Elmegreen B. G., Elmegreen D. M., 1986, ApJ, 311, 554
- Faesi C. M., Lada C. J., Forbrich J., 2016, ApJ, 821, 125
- Freeman P., Rosolowsky E., Kruijssen J. M. D., Bastian N., Adamo A., 2017, MNRAS, 468, 1769
- Fujimoto Y., Tasker E. J., Wakayama M., Habe A., 2014, MNRAS, 439, 936 Gaia Collaboration, 2018, A&A, 616, A11
- Giannetti A., Wyrowski F., Leurini S., Urquhart J., Csengeri T., Menten K. M., Bronfman L., van der Tak F. F. S., 2015, A&A, 580, L7
- Gibson S. J., Taylor A. R., Higgs L. A., Dewdney P. E., 2000, ApJ, 540, 851
- Ginsburg A. et al., 2013, ApJS, 208, 14
- Goodman A. A. et al., 2014, ApJ, 797, 53
- Gratier P. et al., 2012, A&A, 542, A108
- Green J. A. et al., 2012, MNRAS, 420, 3108

- Güsten R., Nyman L. Å., Schilke P., Menten K., Cesarsky C., Booth R., 2006, A&A, 454, L13
- Heyer M. H., Brunt C., Snell R. L., Howe J. E., Schloerb F. P., Carpenter J. M., 1998, ApJS, 115, 241
- Heyer M. H., Carpenter J. M., Snell R. L., 2001, ApJ, 551, 852
- Heyer M., Krawczyk C., Duval J., Jackson J. M., 2009, ApJ, 699, 1092
- Hildebrand R. H., 1983, Q. J. R. Astron. Soc., 24, 267
- Hoare M. G. et al., 2012, PASP, 124, 939
- Honma M. et al., 2012, PASJ, 64, 136
- Jackson J. M., Finn S. C., Rathborne J. M., Chambers E. T., Simon R., 2008, ApJ, 680, 349
- Jackson J. M., Finn S. C., Chambers E. T., Rathborne J. M., Simon R., 2010, ApJ, 719, L185
- Kauffmann J., Pillai T., 2010, ApJ, 723, L7
- Kauffmann J., Bertoldi F., Bourke T. L., Evans N. J. I., Lee C. W., 2008, A&A, 487, 993
- Kauffmann J., Pillai T., Shetty R., Myers P. C., Goodman A. A., 2010a, ApJ, 712, 1137
- Kauffmann J., Pillai T., Shetty R., Myers P. C., Goodman A. A., 2010b, ApJ, 716, 433
- Kauffmann J., Pillai T., Goldsmith P. F., 2013, ApJ, 779, 185
- Kim W.-T., Ostriker E. C., 2002, ApJ, 570, 132
- Koda J. et al., 2009, ApJ, 700, L132
- Krumholz M. R., McKee C. F., 2008, Nature, 451, 1082
- Lallement R., Babusiaux C., Vergely J. L., Katz D., Arenou F., Valette B., Hottier C., Capitanio L., 2019, A&A, 625, A135
- Larson R. B., 1981, MNRAS, 194, 809
- Leroy A. K. et al., 2015, ApJ, 801, 25
- Lin C. C., Shu F. H., 1964, ApJ, 140, 646
- Liu X.-L., Wang J.-J., Xu J.-L., 2013, MNRAS, 431, 27
- Lumsden S. L., Hoare M. G., Urquhart J. S., Oudmaijer R. D., Davies B., Mottram J. C., Cooper H. D. B., Moore T. J. T., 2013, ApJS, 208, 11
- McClure-Griffiths N. M., Dickey J. M., Gaensler B. M., Green A. J., Haverkorn M., Strasser S., 2005, ApJS, 158, 178
- McClure-Griffiths N. M., Dickey J. M., Gaensler B. M., Green A. J., Green J. A., Haverkorn M., 2012, ApJS, 199, 12
- Marshall D. J., Robin A. C., Reylé C., Schultheis M., Picaud S., 2006, A&A, 453, 635
- Martínez-García E. E., González-Lópezlira R. A., Bruzual-A G., 2009, ApJ, 694, 512
- Mattern M. et al., 2018, A&A, 619, A166
- Miville-Deschênes M.-A., Murray N., Lee E. J., 2017, ApJ, 834, 57
- Molinari S. et al., 2010, A&A, 518, L100
- Moore T. J. T., Urquhart J. S., Morgan L. K., Thompson M. A., 2012, MNRAS, 426, 701
- Oka T., Hasegawa T., Sato F., Tsuboi M., Miyazaki A., Sugimoto M., 2001, ApJ, 562, 348
- Ossenkopf V., Henning T., 1994, A&A, 291, 943
- Otrupcek R. E., Hartley M., Wang J.-S., 2000, Publ. Astron. Soc. Aust., 17, 92
- Pan H.-A., Kuno N., 2017, ApJ, 839, 133
- Pandian J. D., Momjian E., Goldsmith P. F., 2008, A&A, 486, 191
- Pearson K., 1901, Phil. Mag., 2, 559
- Peretto N., Fuller G. A., 2009, A&A, 505, 405
- Pettitt A. R., Dobbs C. L., Acreman D. M., Price D. J., 2014, MNRAS, 444, 919
- Pettitt A. R., Dobbs C. L., Acreman D. M., Bate M. R., 2015, MNRAS, 449, 3911
- Pettitt A. R., Egusa F., Dobbs C. L., Tasker E. J., Fujimoto Y., Habe A., 2018, MNRAS, 480, 3356
- Purcell C. R. et al., 2013, ApJS, 205, 1
- Ragan S. E., Henning T., Tackenberg J., Beuther H., Johnston K. G., Kainulainen J., Linz H., 2014, A&A, 568, A73
- Reid M. J. et al., 2009, ApJ, 700, 137
- Reid M. J. et al., 2014, ApJ, 783, 130
- Reid M. J., Dame T. M., Menten K. M., Brunthaler A., 2016, ApJ, 823, 77
- Rice T. S., Goodman A. A., Bergin E. A., Beaumont C., Dame T. M., 2016, ApJ, 822, 52

- Rigby A. J. et al., 2016, MNRAS, 456, 2885
- Rigby A. J. et al., 2019, A&A, 632, A58
- Roberts W. W., 1969, ApJ, 158, 123
- Roman-Duval J., Jackson J. M., Heyer M., Johnson A., Rathborne J., Shah R., Simon R., 2009, ApJ, 699, 1153
- Roman-Duval J., Jackson J. M., Heyer M., Rathborne J., Simon R., 2010, ApJ, 723, 492
- Romero-Gómez M., Mateu C., Aguilar L., Figueras F., Castro-Ginard A., 2019, A&A, 627, A150
- Rosolowsky E., Blitz L., 2005, ApJ, 623, 826
- Rosolowsky E. W., Pineda J. E., Kauffmann J., Goodman A. A., 2008, ApJ, 679, 1338
- Rosolowsky E. et al., 2010, ApJS, 188, 123
- Schruba A. et al., 2017, ApJ, 835, 278
- Schuller F., Menten K. M., Contreras Y., Wyrowski F., Schilke e. a., 2009, A&A, 504, 415
- Schuller F. et al., 2017, A&A, 601, A124
- Schuller F. et al. 2020, MNRAS, in press
- Scoville N. Z., Solomon P. M., 1975, ApJ, 199, L105
- Sewilo M., Watson C., Araya E., Churchwell E., Hofner P., Kurtz S., 2004, ApJS, 154, 553
- Shetty R., Ostriker E. C., 2006, ApJ, 647, 997
- Siebert A. et al., 2011, MNRAS, 412, 2026
- Siebert A. et al., 2012, MNRAS, 425, 2335
- Simon R., Rathborne J. M., Shah R. Y., Jackson J. M., Chambers E. T., 2006, ApJ, 653, 1325
- Solomon P. M., Rivolo A. R., Barrett J., Yahil A., 1987, ApJ, 319, 730
- Stark A. A., Lee Y., 2006, ApJ, 641, L113
- Stutzki J., Guesten R., 1990, ApJ, 356, 513
- Sun J. et al., 2018, ApJ, 860, 172
- Taylor J. H., Cordes J. M., 1993, ApJ, 411, 674
- Toomre A., 1977, ARA&A, 15, 437
- Tosaki T. et al., 2017, PASJ, 69, 18
- Traficante A., Fuller G. A., Smith R. J., Billot N., Duarte-Cabral A., Peretto N., Molinari S., Pineda J. E., 2018a, MNRAS, 473, 4975
- Traficante A., Lee Y. N., Hennebelle P., Molinari S., Kauffmann J., Pillai T., 2018b, A&A, 619, L7
- Urquhart J. S. et al., 2012, MNRAS, 420, 1656
- Urquhart J. S. et al., 2013a, MNRAS, 431, 1752
- Urquhart J. S. et al., 2013b, MNRAS, 435, 400
- Urquhart J. S., Figura C. C., Moore T. J. T., Hoare M. G., Lumsden S. L., Mottram J. C., Thompson M. A., Oudmaijer R. D., 2014a, MNRAS, 437, 1791
- Urquhart J. S. et al., 2014b, MNRAS, 443, 1555
- Urquhart J. S. et al., 2014c, A&A, 568, A41
- Urquhart J. S. et al., 2015, MNRAS, 446, 3461
- Urguhart J. S. et al., 2018, MNRAS, 473, 1059
- Urquhart J. S. et al., 2019, MNRAS, 484, 4444
- Utomo D., Blitz L., Davis T., Rosolowsky E., Bureau M., Cappellari M., Sarzi M., 2015, ApJ, 803, 16
- Vallée J. P., 2014, ApJS, 215, 1
- Vallée J. P., 2017, Astron. Rev., 13, 113
- Wang K., Testi L., Ginsburg A., Walmsley C. M., Molinari S., Schisano E., 2015, MNRAS, 450, 4043
- Watkins E. J., Peretto N., Marsh K., Fuller G. A., 2019, A&A, 628, A21
- Wei L. H., Keto E., Ho L. C., 2012, ApJ, 750, 136
- Wienen M., Wyrowski F., Schuller F., Menten K. M., Walmsley C. M., Bronfman L., Motte F., 2012, A&A, 544, A146
- Wienen M. et al., 2015, A&A, 579, A91
- Wienen M., Wyrowski F., Menten K. M., Urquhart J. S., Walmsley C. M., Csengeri T., Koribalski B. S., Schuller F., 2018, A&A, 609, A125
- Williams J. P., de Geus E. J., Blitz L., 1994, ApJ, 428, 693
- Wilson C. D. et al., 2011, MNRAS, 410, 1409
- Wong T. et al., 2011, ApJS, 197, 16
- Wu Y. W. et al., 2014, A&A, 566, A17
- Zucker C., Battersby C., Goodman A., 2015, ApJ, 815, 23

# SUPPORTING INFORMATION

Supplementary data are available at MNRAS online.

### SEDIGISM\_Cloud\_Catalogue\_supplementary

Please note: Oxford University Press is not responsible for the content or functionality of any supporting materials supplied by the authors. Any queries (other than missing material) should be directed to the corresponding author for the article.

<sup>1</sup>School of Physics & Astronomy, Cardiff University, Queen's building, The parade, Cardiff CF24 3AA, UK

<sup>2</sup>Max-Planck-Institut für Radioastronomie (MPIfR), Auf dem Hügel 69, D-53121 Bonn, Germany

<sup>3</sup>School of Physical Sciences, University of Kent, Ingram Building, Canterbury, Kent CT2 7NH, UK

<sup>4</sup>Department of Astronomy, University of Florida, 211 Bryant Space Sciences Center, Gainesville, FL 32611-2055, USA

<sup>5</sup>Laboratoire d'Astrophysique de Marseille, CNRS, Aix Marseille Université, UMR 7326, F-13388 Marseille, France

<sup>6</sup>Department of Physics & Astronomy, West Virginia University, P. O. Box 6315, Morgantown, WV 26506, USA

 <sup>7</sup>Space Science Institute, 4765 Walnut St. Suite B, Boulder, CO 80301, USA
 <sup>8</sup>School of Science and Technology, University of New England, NSW 2351, Australia

<sup>9</sup>Observatorio Astrofisico di Arcetri, Largo Enrico Fermi 5, 1-50125 Firenze, Italy

<sup>10</sup>Max-Planck-Institut für Astronomie, Königstuhl 17, D-69117 Heidelberg, Germany

<sup>11</sup>Laboratoire d'astrophysique de Bordeaux, CNRS, Univ. Bordeaux, B18N, allée Geoffroy Saint-Hilaire, F-33615 Pessac, France

<sup>12</sup>Departamento de Astronomía, Universidad de Chile, Casilla 36-D, Santiago, Chile

<sup>13</sup>Department of Physics & Astronomy, University of Exeter, Stocker Road, Exeter EX4 4QL, UK

<sup>14</sup>Astrophysics Research Institute, Liverpool John Moores University, 146 Brownlow Hill, Liverpool L3 5RF, UK

<sup>15</sup>Haystack Observatory, Massachusetts Institute of Technology, 99 Millstone Road, Westford, MA 01886, USA

<sup>16</sup>Korea Astronomy & Space Science Institute, 776 Daedeokdae-ro, 34055 Daejeon, Republic of Korea

<sup>17</sup>Department of Physics, Faculty of Science, Hokkaido University, Sapporo 060-0810, Japan

<sup>18</sup>Istituto di Astrofisica e Planetologia Spaziali, INAF, via Fosso del Cavaliere 100, I-00133 Roma, Italy

<sup>19</sup>I. Physikalisches Institut, Universität zu Köln, Zülpicher Str 77, D-50937 Köln, Germany

<sup>20</sup>European Southern Observatory, Alonso de Cordova 3107, Casilla 19001, Santiago 19, Chile

<sup>21</sup>Astronomy Department, University of Wisconsin, 475 North Charter St, Madison, WI 53706, USA

<sup>22</sup>Department of Space, Earth and Environment, Chalmers University of Technology Onsala Space Observatory, SE-439 92 Onsala, Sweden

<sup>23</sup>INAF – Osservatorio Astronomico di Cagliari, Via della Scienza 5, I-09047 Selargius (CA), Italy

<sup>24</sup>ipag, Univ. Grenoble Alpes, CNRS, F-38000 Grenoble, France

<sup>25</sup>School of Engineering, Macquarie University, NSW 2109, Australia

<sup>26</sup>Kavli Institute for Astronomy and Astrophysics, Peking University, 5 Yiheyuan Road, Haidian District, Beijing 100871, People's Republic of China

This paper has been typeset from a TeX/LATeX file prepared by the author.