# Inclusion Depth for Contour Ensembles

Chaves-de-Plaza, Nicolas F.; Mody, Prerak; Staring, Marius; van Egmond, Rene; Vilanova, Anna;
Hildebrandt, Klaus

**Important note**
To cite this publication, please use the final published version (if applicable).
Please check the document version above.

# Inclusion Depth for Contour Ensembles

Nicolas F. Chaves-de-Plaza ⓘ, Prerak Mody ⓘ, Marius Staring ⓘ, René van Egmond ⓘ, Anna Vilanova ⓘ,
and Klaus Hildebrandt ⓘ

*Abstract*—Ensembles of contours arise in various applications like simulation, computer-aided design, and semantic segmentation. Uncovering ensemble patterns and analyzing individual members is a challenging task that suffers from clutter. Ensemble statistical summarization can alleviate this issue by permitting analyzing ensembles' distributional components like the mean and median, confidence intervals, and outliers. Contour boxplots, powered by Contour Band Depth (CBD), are a popular non-parametric ensemble summarization method that benefits from CBD's generality, robustness, and theoretical properties. In this work, we introduce Inclusion Depth (ID), a new notion of contour depth with three defining characteristics. First, ID is a generalization of functional Half-Region Depth, which offers several theoretical guarantees. Second, ID relies on a simple principle: the inside/outside relationships between contours. This facilitates implementing ID and understanding its results. Third, the computational complexity of ID scales quadratically in the number of members of the ensemble, improving CBD's cubic complexity. This also in practice speeds up the computation enabling the use of ID for exploring large contour ensembles or in contexts requiring multiple depth evaluations like clustering. In a series of experiments on synthetic data and case studies with meteorological and segmentation data, we evaluate ID's performance and demonstrate its capabilities for the visual analysis of contour ensembles.

*Index Terms*—Uncertainty visualization, contours, ensemble summarization, depth statistics.

## I. INTRODUCTION

D IFFERENT applications in simulation, computer-aided design, and semantic segmentation have to deal with

Nicolas F. Chaves-de-Plaza is with TU Delft, 2628 Delft, CD, The Netherlands, and also with Holland PTC, 2629 Delft, JH, The Netherlands (e-mail: n.f.chavesdeplaza@tudelft.nl).

Prerak Mody is with the Department of Radiology, Leiden University Medical Center, 2333 Leiden, ZA, The Netherlands, and also with Holland PTC, 2629, JH Delft, The Netherlands (e-mail: p.p.mody@lumc.nl).

Marius Staring is with the Department of Radiology, Department of Radiation Oncology, Leiden University Medical Center, 2333 Leiden, ZA, The Netherlands (e-mail: m.staring@lumc.nl).

René van Egmond and Klaus Hildebrandt are with TU Delft, 2628, CD Delft, The Netherlands (e-mail: R.vanEgmond@tudelft.nl; k.a.hildebrandt@tudelft.nl).

Anna Vilanova is with TU Eindhoven, 5612 Eindhoven, AZ, The Netherlands (e-mail: a.vilanova@tue.nl).

This article has supplementary downloadable material available at https://doi.org/10.1109/TVCG.2024.3350076, provided by the authors.

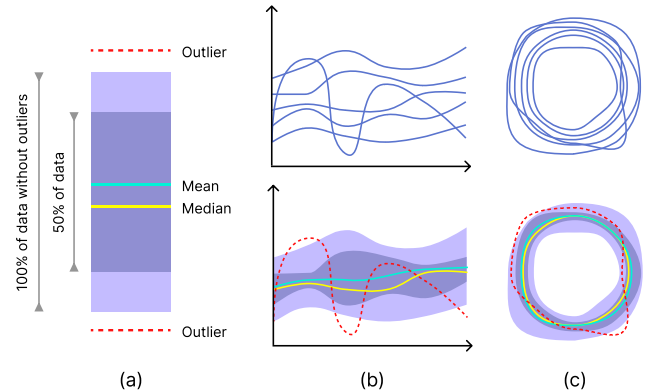Digital Object Identifier 10.1109/TVCG.2024.3350076

Fig. 1. Extension of the boxplot idiom (a) to the functional (b) and contour (c) data types.

ensembles of curves. Analyzing these ensembles permits understanding uncertainties in the results. We focus on ensembles of spatiotemporal scalar fields from which one can extract contours, closed and consistently-oriented curves. These appear in several domains. One example is meteorology, where analysts use ensembles of weather forecasts to analyze the predictions' variability under different initial conditions or changes in the computational model [1]. Another example is semantic segmentation, where ensembles are used to quantify the uncertainty that might come from the training data or the model [2]. In image-guided medical specialties, ensembles of segmentations are analyzed for planning the patients' treatments [3].

Visual inspection of the ensemble can facilitate its analysis and understanding. Spaghetti plots, which draw each contour in the ensemble using a different color, are a popular technique. They are attractive because they are accessible, represent all the data, and are simple to implement. Nevertheless, as the ensemble size increases, spaghetti plots become cluttered, potentially hiding interesting features of the ensemble. Motivated by these limitations, ensemble summarization methods have been proposed. They reduce information by extracting features of interest, such as representative members and contour variability, from the ensemble and visualize them using visual encodings based on lines and bands [4], [5], [6], [7].

A successful contour summarization technique is the contour boxplot (CBP) [7], which has been used in the fields of meteorology [8] and medicine [9], [10], [11]. As Fig. 1 illustrates, like traditional boxplots, CBPs depict four statistical features of an ensemble: the median, the trimmed mean, confidence intervals, and outliers. Underlying the CBP is the concept of statistical depth, which extends univariate order and rank statistics to

complex multivariate datasets by establishing a center-outward measure of centrality for the ensemble members [12].

In this paper, we propose an alternative notion of contour depth called Inclusion Depth (ID). ID contributes to the arsenal of depth-based contour analysis methods in three ways.

First, ID provides a novel statistical depth for ensembles of contours. It draws inspiration from Half-Region Depth (HRD) [13] and generalizes HRD from the class of functions to contours. This connection to HRD endows ID with theoretical properties and enables computational advantages analogous to those of HRD, also for ensembles of contours. In Section IV, we present the ID framework, detailing how it overcomes the challenges that extending the HRD approach to the contour case brings.

Second, ID leverages a simple principle that makes it accessible and facilitates the interpretation of the results. Specifically, ID leverages the inside/outside relationships between contours to estimate the ensemble's depth. To compute a contour's ID we compute how many other contours of the ensemble the contour contains and in how many other contours it is contained. Intuitively, a highly central contour has similar values for both quantities. An outlier might have an asymmetry of these quantities, if it's a magnitude outlier, or lower values for both, in the case of a shape outlier.

Third, the computation of ID scales better than the state-of-the-art Contour Band Depth (CBD), which was introduced jointly with the CBP idiom [7]. As Section III shows, for a $N$-contour ensemble, ID requires $\mathcal{O}(MN^2)$ operations while CBD needs at least $\mathcal{O}(MN^3)$, where $M$ is the contour size (i.e., resolution of the binary mask). In Section VI, we evaluate ID, empirically showing that performing only pairwise comparisons does not degrade ID's performance and yields depth scores qualitatively comparable to CBD's.

We further demonstrate the practical use of ID in Section VII by performing depth-based exploratory analysis of several real datasets from diverse domains like segmentation in radiotherapy and meteorological forecasting. Based on the results, we expect the faster but still performant ID will enable visual analysis of larger ensembles using depth-based visualizations like CBP, which allows both quantitative and qualitative interpretation of contour ensembles. Furthermore, it will bring applications that require multiple or/and fast depth evaluations like regression [14] and clustering [15] within reach.

## II. RELATED WORK

Our method fits in the context of uncertainty visualization. Ensembles permit quantifying predictive uncertainties due to changes in the initial conditions, the training data, or the model parameters [2]. Existing ensemble visualization techniques can be cataloged based on their data type, visualization method, and analytic task [16].

There are several alternatives for presenting a visual overview of contour ensembles. Spaghetti plots are a composition-after-visualization technique that plots each contour using a different color [17]. Although straightforward to implement and interpret, spaghetti plots become cluttered as the size of the ensemble grows, potentially hiding trends and interesting members. To address this issue, several ensemble summarization techniques have been proposed in recent years that aggregate contour data into salient features before visualizing it. Most available summarization techniques share a visual language that uses contour lines for the ensembles' representative members like the median, mean, and outliers, and bands for areas of interest like the ensemble's spread [18] and confidence intervals [7].

Available summarization techniques differ in the features they compute and the assumptions they make. Parametric model-based techniques assume a data distribution and use available models to derive statistical quantities. One approach fits a Gaussian distribution on the contours' PCA-reduced signed distance field (SDF) transform and uses it to derive a median and calculate bands [5], [19]. It is also possible to use a Gaussian model to describe each grid point and use this model together with iso-contour density and level-cross probability to extract the iso-contours' probability density [20]. Parametric techniques are conceptually attractive as they permit extracting information analytically [5]. Nevertheless, they impose assumptions on the data, like normality, which limits the applicability in practice. Our method is fully non-parametric, using a depth-induced ranking of the contours to detect outliers and derive quantities of interest like the median and robust mean.

The family of data-based non-parametric methods does not impose assumptions on the data distribution and, therefore, can describe the ensemble data on each point more accurately [21], [22]. Local summarization methods operate on the grid in which contours lie, computing point-wise statistics. Examples are contour probability plots, which extract bands by thresholding a scalar field of percentages [18], and EnConVis [4], which performs point-wise kernel density estimation, and then uses the per-point density to extract bands and representatives. Contour grid points are not independent of each other, so computing summaries based solely on point-wise estimates can fail to consider global characteristics of the contour data like the topological relationships between contours.

The method introduced in [6] uses a vector-to-closest-point representation along the contours boundary points to quantify their centrality based on the vector lengths and directions. This approach requires only comparisons between contours, making it more efficient than CBD. Nevertheless, it uses parametric statistical models that require parameter fitting to obtain the centrality estimates. Furthermore, it is unclear how the method performs under different ensemble distributions, which makes it hard to compare to existing contour depth methods like CBD.

## III. BACKGROUND: CONTOUR DEPTH AND BOXPLOTS

### A. Statistical Depth

Statistical depth provides a framework for extending concepts like the median, trimmed mean, and outliers, which depend on the points' ranks and orderings from the univariate to the multivariate case. Given a cloud of $N$ $d$-dimensional points $X \in \mathbb{R}^{N \times d}$, a depth function $D(z, X) : \mathbb{R}^d \to [0, 1]$ yields a center-outward measure of the centrality or depth of a point $z$ with respect to $X$. Intuitively, the farther away a point $z$ is from

the center of $X$, the lower its centrality. In practice, there are different methods for computing $D(z, X)$, which come with different guarantees in terms of the function's behavior like invariance to different geometric transformations of $X$ [12].

Statistical depth functions were originally devised to handle multivariate data. Nevertheless, their performance might decrease when $d \gg N$ due to the curse of dimensionality [23]. Furthermore, in some cases, data is more naturally represented as functions. In response to these observations, several definitions of depth that apply to functional data have been recently proposed [13], [23]. Two predominant functional depth methods are Band Depth (BD) [23] and Half-Region Depth (HRD) [13]. Inspired by the multivariate simplicial depth [24], BD computes a function's depth by comparing it to the bands formed by all other subsets of functions in the ensemble. Contour Band Depth, presented in the next subsection, generalizes BD's formulation and extends it to the case of contours.

Instead of forming bands, HRD looks at the proportion of functions lying on each side of the function of interest to determine its depth. The multivariate analog of HRD is Tukey's half-space depth [25]. HRD is more computationally efficient than BD, requiring only $N$ comparisons per function. Furthermore, it has been shown to yield comparable depths to BD [13]. The proposed Inclusion Depth generalizes HRD's formulation and extends it to the case of contours. In the following, we outline HRD.

Let $X = \{x_1, x_2, \ldots, x_N\}$ with $x_i : I \to \mathbb{R}$ be an ensemble of functions defined on the compact interval $I$. The graph of a function $x \in X$ can be defined as

$$G(x) = \{(t, x(t)), t \in I\}, \quad (1)$$

The epi and hypographs of $x$, which correspond to the regions above and below $G(x)$, can be defined as

$$hyp(x) = \{(t, y) \in I \times \mathbb{R} : y \leq x(t)\},$$
$$epi(x) = \{(t, y) \in I \times \mathbb{R} : y \geq x(t)\}. \quad (2)$$

The HRD of $x$ can be computed by evaluating the proportion of times $G(x)$ is contained in the epi and hypographs of other functions of the ensemble. Formally,

$$\text{HRD}(x|X) = \min\{\text{IN}_{hyp}(x), \text{IN}_{epi}(x)\}, \quad (3)$$

where

$$\text{IN}_{hyp}(x) = \frac{1}{N} \sum_{i=1}^{N} G(x) \subset hyp(x_i),$$

$$\text{IN}_{epi}(x) = \frac{1}{N} \sum_{i=1}^{N} G(x) \subset epi(x_i), \quad (4)$$

where $A \subset B$ is 1 if $A$ is contained in $B$ and 0 otherwise.

HRD in (3) attains its maximum value of 0.5 when $G(x)$ is contained in as many epi and hypographs of the other functions in the ensemble. The HRD satisfies several of the properties of a valid depth function [26]: linear invariance, maximality at the center, monotonically decreasing on rays, and upper-semicontinuity. Finally, a finite-dimensional version can

be obtained by drawing $d$ samples from $I$. When $d = 1$, the Half-Region Depth is equivalent to the Tukey depth.

### B. Contour Band Depth

Statistical depth allows for robust and model-free exploratory data analysis. Contour Band Depth (CBD) permits applying the depth methodology to contours [7]. Similarly to functional BD, CBD computes a contour's depth by determining how many bands formed by all other possible $J$-sized contour subsets (where $J \in \mathbb{Z}$ and $J \geq 2$) contain the contour. A contour is in a band if it contains the intersection of the band's contours and is contained in their union. To reduce the computational cost of verifying contour containment in $\sum_{i=2}^{N} \binom{N}{i}$ bands (where $N$ is the size of the contour ensemble), $J = 2$ is used. To alleviate the tendency of CBD with $J = 2$ to produce depth ties, a modified CBD (mCBD) was proposed [7]. Instead of strictly enforcing the containment property, mCBD considers the proportion of the contour that falls outside each band when computing its depth. CBD and mCBD compute an ensemble's depths in $\mathcal{O}(MN^3)$ time, where M is the contour size (i.e., binary mask resolution).

### C. Contour Boxplots

Boxplots offer a visualization of a dataset's summary statistics. Specifically, as Fig. 1 illustrates, a boxplot has four components. The gold and blue-colored lines represent the median and the trimmed mean, respectively. The trimmed mean is the average of the dataset with the outliers removed. Purple bands around the mean encode the interquartile range. Finally, outliers are shown using red dashed lines. As the middle and right side of Fig. 1 shows, the idea of boxplots can be extended to ensembles of functional [27] and contour [7] types through the concept of functional and contour depth. In these cases, the per-member depth values are used to compute the different statistics. The median is the member with the highest depth value and the interquartile ranges are bands formed by members whose depths fall in the specified ranges. Finally, the members with the lowest depths are flagged as outliers.

## IV. INCLUSION DEPTH

In this section, we introduce Inclusion Depth (ID). While ID can be defined for contours in $\mathbb{R}^2$ and $\mathbb{R}^3$, for the sake of simplicity, we consider the two-dimensional case.

Let $C = \{c_1, c_2, \ldots, c_N\}$ be an ensemble of contours, where a contour $c_i$ is a pair of a function $F_i : \Omega \to \mathbb{R}$ and an isovalue $q_i \in \mathbb{R}$. Here $\Omega$ is a compact domain in $\mathbb{R}^2$, such as a rectangle. A contour encloses a subset in the plane that we call the inside region

$$in(c_i) = \{p \in \Omega | F_i(p) < q_i\}. \quad (5)$$

ID is based on a simple principle. We evaluate for all pairs $c_i, c_j \in C$ whether or not $in(c_i)$ is contained in $in(c_j)$. Then, we form the fraction of contours of $C$ in which $in(c_i)$ is contained,

$$\text{IN}_{in}(c_i) = \frac{1}{N} \sum_{j=1}^{N} in(c_i) \subset in(c_j), \quad (6)$$
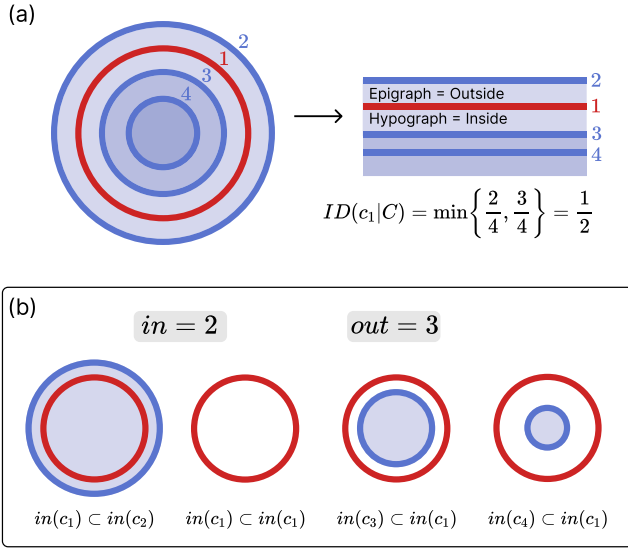
Fig. 2. Example of the ID computation for a 4-contour ensemble. In red is the contour for which we are currently estimating the depth. (b) shows the four comparisons that need to be performed to compute ID based on (8). Note that $c1 \subset c1$ (second column) counts for the inside and outside relationships.

and the fraction of contours of C that are contained $in(c_i)$,

$$\text{IN}_{out}(c_i) = \frac{1}{N} \sum_{j=1}^{N} in(c_j) \subset in(c_i). \qquad (7)$$

In the sums in (6) and (7), we interpret $in(c_i) \subset in(c_j)$ as the numerical value 1 if $in(c_i)$ is contained in $in(c_j)$ and as 0 otherwise. The ID is the minimum of the two fractions

$$\text{ID}(c_i|C) = \min\{\text{IN}_{in}(c_i), \text{IN}_{out}(c_i)\}. \qquad (8)$$

Fig. 2 illustrates the process of computing a contour's ID in a four-contour ensemble. As (a) depicts, ID is related to HRD. Specifically, the proof sketch in the appendix, available online, shows that if there is an invertible transform mapping the contours to graphs of functions, our definition of ID is the same as HRD. (b) presents the four comparisons required to compute $\text{IN}_{in}$ and $\text{IN}_{out}$ in (8). Note that the example uses perfectly nested simple single-loop contours for illustrative purposes. In practice, contours often have more complex shapes, are not necessarily nested, and can have multiple connected components.

ID is more general than HRD, accommodating the different topologies that arise in higher dimensions. Fig. 3 shows examples of how ID deals with different cases. Note that, by subset operations, the definitions of $\text{IN}_{in}(c)$ and $\text{IN}_{out}(c)$ in (6) and (7) ensure that the two contours under comparison are nested. As the bottom right panel of Fig. 3 shows, when contours are not nested, the comparison will not add to the inside or outside counts, effectively reducing the depth of the contour under consideration. ID's results are invariant to homeomorphic transformations of the domain $\Omega$, a general class of transformations that includes affine transformations and Möbius transformations. Additionally, ID's results are invariant to the choice of inside and outside. We sketch proofs of these properties in the appendix,
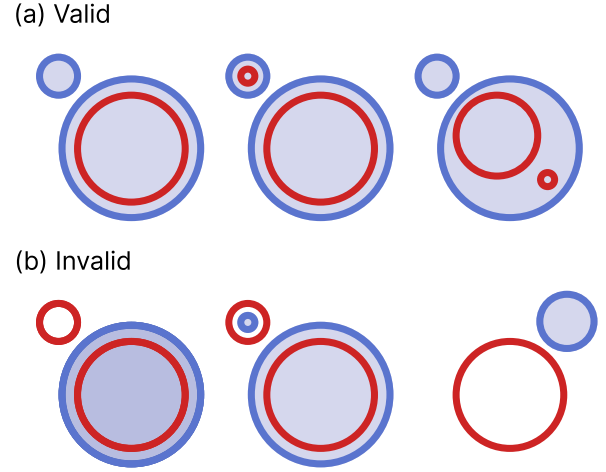


Fig. 3. Examples of how ID deals with different cases. If contours are nested (a), their relationship will add to the inside/outside counts. In other cases (b), the inside/outside counters will not increase, effectively reducing the overall depth.

---

**Algorithm 1:** Inclusion Depth (ID).

**Require:** $C, N$ ▷ Contour ensemble, number of contours
$\mathbf{d}^{\text{ID}} \leftarrow \{\}$ ▷ Inclusion depths
**for** $i = 1$ to $N$ **do**
  $num\_in \leftarrow 0; num\_out \leftarrow 0$ ▷ Inside/outside counts
  **for** $j = 1$ to $N$ **do**
    $num\_in \leftarrow num\_in + [in(c_i) \subset in(c_j)]$
    $num\_out \leftarrow num\_out + [in(c_j) \subset in(c_i)]$
  **end for**
  $\text{IN}_{in}(c_i) = num\_in/N$
  $\text{IN}_{out}(c_i) = num\_out/N$
  $\mathbf{d}^{\text{ID}} \leftarrow \mathbf{d}^{\text{ID}} \bigcup \min\{\text{IN}_{in}(c_i), \text{IN}_{out}(c_i)\}$
**end for**
**return** $\mathbf{d}^{\text{ID}}$

---

available online, and point the interested reader to the set theory literature, which further elaborates on properties of the $\subset$ operator like transitivity [28].

Algorithm 1 shows how to compute the ID of a contour ensemble. For computations, we assume $\Omega$ to be a rectangle, e.g., the bounding box of the ensemble of contours, and discretize the rectangle by a regular grid of size $M$. ID's scaling behavior depends mainly on the ensemble's size ($N$). Nevertheless, the grid size will also impact the algorithm's scaling behavior when performing the inside/outside comparisons. Therefore, ID has a computational complexity of $\mathcal{O}(MN^2)$, which is a significant improvement over the $\mathcal{O}(MN^3)$ complexity of CBD.

## V. EPSILON INCLUSION DEPTH

If the ensemble's contours are non-smooth and intersect, like the examples in the center of Fig. 4 show, inside/outside relationships will be ambiguous. In these cases, ID will produce ties and low-depth scores that reduce the method's practical utility. In this section, we present the Epsilon Inclusion Depth (eID) that relaxes the definitions of inside/outside in ID, reducing the effect of highly varying contours on the depth estimate.
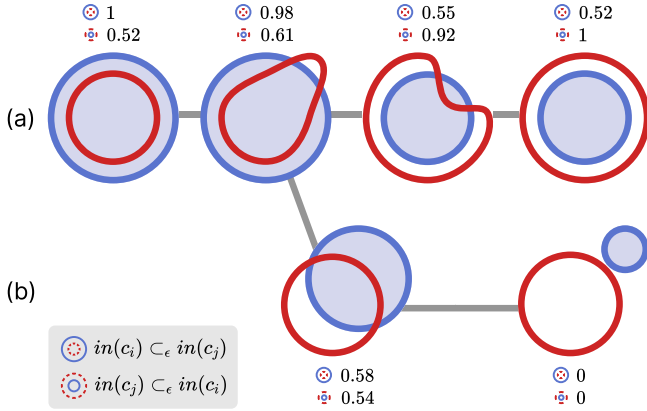
Fig. 4. Examples of computing the inside and outside relationships with the $\subset_\epsilon$ operator in (9) for different contour configurations. In red and blue are contour $\{c_i, c_j \in C\}$. The first row shows the transition of $c_i$ from being completely inside to completely outside of $c_j$. The second row shows the values that $\subset_\epsilon$ yields in ambiguous cases.

For this extension, we proceed analogously to HRD, for which modified HRD (mHRD) alleviates the problem that strongly varying functions pose for HRD by relaxing the requirement that the graph of a function must lie entirely in the epi or hypograph. mHRD determines the average proportion of the domain that a function's graph lies in the hypo and epigraphs of other functions [13]. This strategy is not directly applicable to the case of contours because of the lack of an independent variable. Therefore, we follow a strategy inspired by the modified Contour Band Depth in [7], which operates directly on the contours' domain and therefore does not require a dependent variable.

First, we define the epsilon subset operator $A \subset_\epsilon B$ for two sets $A, B \subset \mathbb{R}^2$. In contrast to the subset operator $\subset$, which returns either 0 or 1, $\subset_\epsilon$ yields a value in the interval [0,1]. It is defined as

$$A \subset_\epsilon B = 1 - \begin{cases} 0 & |A| = 0, \\ |A - B|/|A| & \text{otherwise,} \end{cases} \quad (9)$$

where $|A|$ denotes the area of $A$ and $A - B$ the set difference. Note that $A \subset_\epsilon B$ will be one if $B$ contains $A$. If a part of $A$ lies outside of $B$, $\subset_\epsilon$ will yield lower values.

The definition of eID is analogous to ID except that the $\subset$ operator is replaced by the $\subset_\epsilon$ operator. We consider the values

$$\text{IN}_{in}^\epsilon(c_i) = \frac{1}{N} \sum_{j=1}^{N} in(c_i) \subset_\epsilon in(c_j),$$

$$\text{IN}_{out}^\epsilon(c_i) = \frac{1}{N} \sum_{j=1}^{N} in(c_j) \subset_\epsilon in(c_i). \quad (10)$$

The eID is the minimum of the two values

$$\epsilon\text{ID}(c|C) = \min\{IN_{in}^\epsilon(c), IN_{out}^\epsilon(c)\}. \quad (11)$$

Fig. 4 shows how $\subset_\epsilon$ works across a variety of cases. As the extremes of the first row illustrate, when $in(c_i)$ (red) is completely inside or outside of $c_j$ (blue), the difference between $in(c_i) \subset_\epsilon in(c_j)$ and $in(c_j) \subset_\epsilon in(c_i)$ is the largest. When the

relationship between the contours is ambiguous, the second row of the figure shows that the difference shrinks. Also, the values of these quantities decrease, which has the effect of reducing the contribution of the $c_i/c_j$ comparison to the overall depth calculation. Finally, eID is invariant to area-preserving transformations. We sketch the proof of this property in the appendix, available online.

As the next sections show, eID provides meaningful results even when contours have many intersections. The implementation of eID only requires swapping $\subset$ for $\subset_\epsilon$ in Algorithm 1. eID maintains ID's computational complexity of $\mathcal{O}(MN^2)$.

## VI. EXPERIMENTS

In this section, we perform an extensive evaluation of the Inclusion Depth (ID) method using synthetic data. Specifically, we assess the scaling behavior of ID as the dataset's size increases and investigate the robustness of estimators derived with ID and the method's performance at identifying outliers. Before continuing with the experiments, we detail our experimental setup.

### A. Experimental Setup

In our experiments, we compare ID and eID with Contour Band Depth (CBD) and its modified version (mCBD). CBD's only parameter, the number of contours forming the band (J), is set to J=2 for all experiments.

We implemented the CBD and ID methods and the experiments in a Python-based framework.[1] Contour depth methods receive as input a list of binary Numpy [29] arrays of size $M = 300 \times 300$ corresponding to a discretization of (5). These binary masks can be obtained, for example, as the output of a segmentation algorithm or by thresholding scalar fields using an iso-value. While acceleration through parallelization is possible, our focus in implementation prioritizes asymptotic algorithmic scaling over specific optimizations.

Similar to [7], we use synthetic ensembles of circular shapes contaminated with outliers to assess the methods' performance. We extend the experiments of contour depth by considering different types of outliers separately, following the experimental paradigm used to evaluate the functional Half-Region Depth [13]. The first row of Fig. 5 showcases the different outliers we consider (orange contours). We expect the depth scores that CBD and ID yield to be lower for contours that deviate from the ensemble's main trend and higher for those that conform to it.

To generate ensembles of contours contaminated with outliers, we define a stochastic model from which we can sample shapes. The model results from a mixture of a base model $r_0$ and a second model $r_1$, which depends on the outlier type under consideration. For both $r_0$ and $r_1$, we use stochastic processes indexed by the shape's angle, yielding angle-correlated values for the shape's radius. We define the base model $r_0$ as

$$r_0(\theta) = f_0(\theta) + \epsilon_0(\theta), \quad (12)$$

---

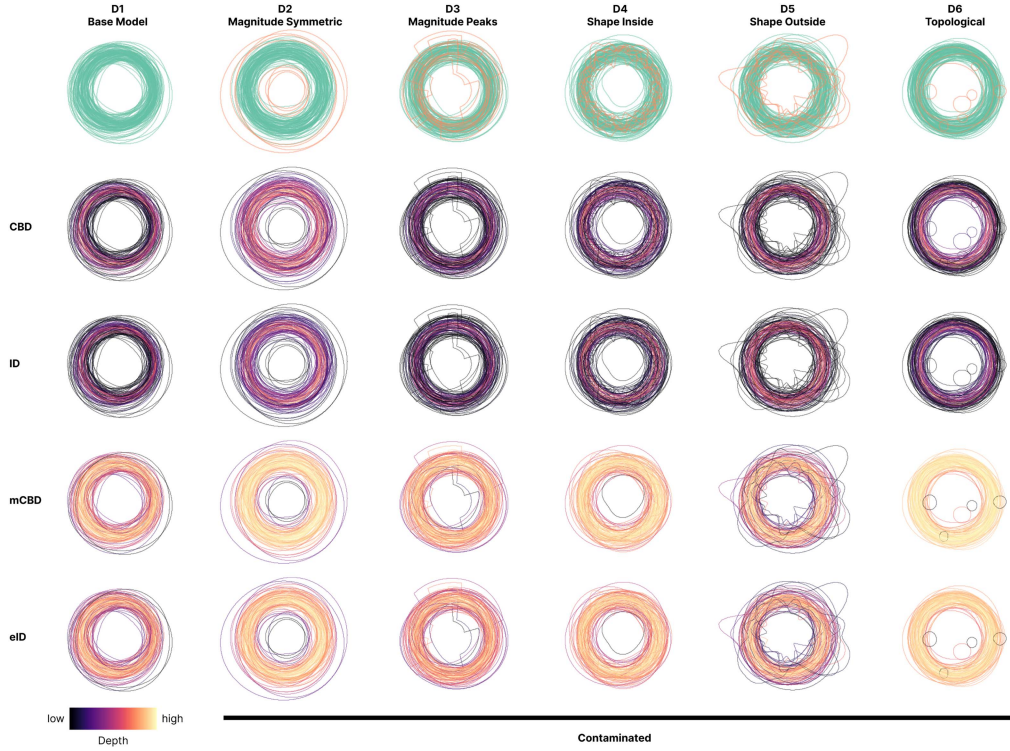[1] Code can be found at https://graphics.tudelft.nl/inclusion-depth-paper

Fig. 5. First row presents an overview of the synthetic datasets we used in the experiments, with the outliers highlighted in orange. The last four rows plot the ensembles assigning the lines' colors based on the depths each method yielded. Darker and brighter colors denote lower and higher depth values, respectively. The color scale was scaled based on the min and max depth value per dataset/depth method combination to facilitate the comparison of the depth-induced rankings across methods.

where $\theta \in \mathbb{R}^1 00$ is a vector containing 100 equally spaced samples of the interval $[0, 2\pi]$ and $f_0(\theta) = 0.5$ is the mean radius function.

To add randomness to the mean shape, we use Gaussian Processes (GP), defined by a mean and an exponentiated quadratic kernel

$$k_{mid}(\theta_i, \theta_j) = \sigma_{mid}^2 \exp\left(-\frac{(g(\theta_i) - g(\theta_j))^2}{2l_{mid}^2}\right), \qquad (13)$$

where $\theta_i, \theta_j \in \theta$, $g : \mathbb{R} \to \mathbb{R}$ is a function that transforms the domain and $mid$ can be zero or one depending on whether we refer to $r_0$ or $r_1$ in (14).

We define $\epsilon_0(\theta)$ in (12) as the sum of two zero-mean GPs with $g = \sin$ and $g = \cos$ in (13), respectively. Using these periodic functions ensures that the start and end of the $\theta$ interval are mapped to the same radius. The kernel's parameters $\sigma_0$ and $l_0$ define the shape of the contour by affecting the amplitude and the frequency of the angle-correlated noise. We set $\sigma_0 = 0.003$ and $l_0 = 0.9$.

To obtain a binary mask from the zero-centered shape defined by the polar coordinates $(\theta, r(\theta))$, we convert them to Cartesian coordinates using $y = r \sin(\theta)$ and $x = r \cos(\theta)$, and rasterize the resulting closed polygon in a square grid with the target size $M$ with scikit-image's polygon2mask. The panel in the upper left corner of Fig. 5 shows a $N = 100$ ensemble generated by sampling the base model $r_0$ (D1).

For the experiments, we define five datasets of contour ensembles (D2-D6 in Fig. 5) based on the three types of outliers we describe next. In all cases, we obtain an outlier-contaminated ensemble by sampling from the mixture

$$r(\theta) = r_0(\theta) + \rho r_1(\theta), \qquad (14)$$

where $\rho \sim Bern(0.1)$ introduces an outlier with a probability of 0.1 and $r_1$ is defined analogously to $r_0$ in (12). In the following, we describe the different outlier types.

First, we consider magnitude outliers in which we alter the shape's mean radius. We define the auxiliary random variable $sign = 2\gamma - 1$ where $\gamma \sim Bern(0.5)$. $sign$ indicates whether the magnitude contamination corresponds to shrinking (-1) or enlarging (1) the shape. The first dataset with magnitude outliers is the Symmetric Magnitude Contamination (D2) for which $f_1(\theta) = 0.3 \cdot sign$. We define a second dataset with magnitude outliers which we call Peaks Magnitude Contamination (D3). Instead of changing the magnitude of the shape's radius, in D3 we only contaminate a subinterval $(\theta_l, \theta_r)$ of $\theta$ where $\theta_l < \theta_r$ and both $\theta_l$ and $\theta_r$ are uniformly distributed random variables. Specifically, for D3, we define $f_1$ as

$$f_1(\theta) = \begin{cases} sign \cdot inc & \theta_l \leq \theta \leq \theta_r \\ 0 & \text{otherwise} \end{cases},$$

where $inc = 0.3$, and $\theta_l$ and $\theta_r$ are defined for every $\theta_i \in \theta$.

The second type of outlier we consider is shape outliers. To obtain shape outliers, instead of altering the mean radius of the
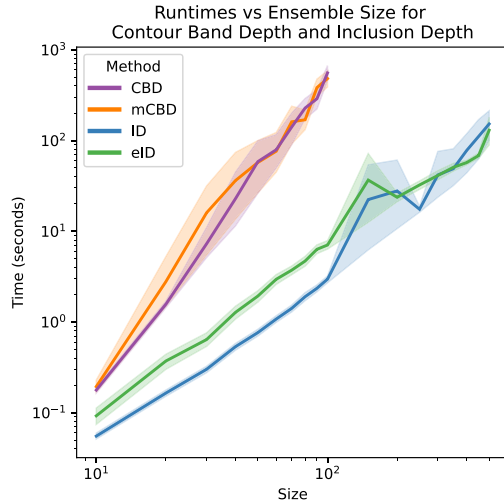
Fig. 6. Comparison of mean runtimes across datasets and replications of CBD, mCBD, ID and eID. Both $x$ and $y$-axis use logarithmic scales and shaded area denotes the 95 percent confidence interval.

TABLE I
MEAN AND STANDARD DEVIATION OF THE PREPROCESSING (T1), DEPTH CALCULATION LOOP (T2) AND FULL (T3=T1+T2) TIMES IN SECONDS FOR D1 WITH $N = 100$

| Method | t1 (secs) | t2 (secs) | t3 (secs) |
|--------|-----------|-----------|-----------|
| CBD | $6.75 \pm 1.77$ | $612.31 \pm 351.40$ | $619.06 \pm 351.14$ |
| mCBD | $6.48 \pm 1.46$ | $697.02 \pm 328.91$ | $703.50 \pm 328.49$ |
| ID | $0.00 \pm 0.00$ | $2.31 \pm 0.37$ | $2.31 \pm 0.37$ |
| eID | $0.00 \pm 0.00$ | $7.37 \pm 3.98$ | $7.37 \pm 3.98$ |

circular shape, we modify the parameters of the covariance matrix of $\epsilon_1$ which define the amplitude ($\sigma_1$) and the frequency ($l_1$) of the noise along the shape's boundary. Specifically, increasing $\sigma_1$ leads to higher amplitude while increasing $l_1$ increases the number of peaks. For the Shape Inside (D4) dataset, we keep $\sigma_1 = 0.003$ but decrease the frequency to $l_1 = 0.01$ to ensure that the shape varies while staying within the ensemble's envelope. For the Shape Outside (D5) dataset, we set $\sigma_1 = 0.009$ and $l_1 = 0.04$, which results in highly varying shapes that spill outside the bounds defined by the normal members of the ensemble. We expect D4 outliers to be more challenging to detect than D5 ones, given that they fall inside the ensemble's envelope.

The final type of outlier we consider are topological outliers which correspond to contours that have holes or disconnected components not present in other members of the ensemble. To create the Different Topologies dataset (D6), we randomly downscale $r_1$ using a uniform distribution between 0.1 and 0.2 for the scaling factor. Note that we use the same parameters for $r_1$ as for $r_0$. After determining the $(x, y)$ coordinates of the shrank shape, we translate them to a random location that lies either inside or outside (with equal probability) of the mean circular shape defined by $r_0$.

For the experiments, we consider several ensemble sizes $N \in \{i * 10 : 1 \le i \le i_{\max}$, where $i_{\max} = 10$ for CBD and $i_{\max} = 30$ for ID. We compute 10 realizations of each dataset/size/depth method combination to establish the results' statistical significance. We ran all the experiments presented in this section on a Mac Book Pro (2022) with an M1 Pro processor (without GPU acceleration) and 32 GB RAM.

### B. Experiment 1: Scaling Behavior

Fig. 6 depicts the time in seconds that each depth method takes for ensembles of different sizes. For each size, we compute the mean and standard deviation across replications and datasets (D1-6). The first thing to note is that we only ran CBD methods

until $N = 100$. After this point, the CBD method took too long to compute. In contrast, we considered ensembles up to size $N = 300$ for ID. The figure shows how ID and eID, with a computational complexity of $\mathcal{O}(MN^2)$, scale more favorably than CBD methods, which are $\mathcal{O}(MN^3)$.

In addition to the aggregated runtime, we investigated the time the preprocessing and depth calculation loop portions of each method take. Table I shows this information for D1 with $N = 100$. As the table shows, all methods spend most of their time in the depth calculation loop (t2). CBD methods take, on average, an order of magnitude more time than ID methods. The large standard deviations of CBD methods' timings are caused by outlier timings that arose likely due to other processes in the machine interfering with the experiment's process. Within each method family, the modified version takes more time because they require more operations than the strict versions. Finally, CBD methods have a larger preprocessing time (t1) than ID methods, which do not require preprocessing. This is specific to our implementation, which precomputes CBD's bands before starting the depth calculation loop.

### C. Experiment 2: Outlier Detection

Depths can be used to perform robust statistical analysis by removing outliers, which are contours with low depth. For the second experiment, we evaluate ID's performance in identifying outliers in D2-D6 in Fig. 5. Specifically, given a set of outliers $\mathcal{O}_m$ for a method $m$ and a reference set $\mathcal{O}_r$, we compute the percentage of correctly identified outliers with respect to the reference set as

$$PO_{m,r} = \begin{cases} 0 & \text{if } |O_r| = 0 \\ \frac{|O_m \cap O_r|}{|O_r|} & \text{otherwise,} \end{cases} \tag{15}$$

where $|\cdot|$ denotes the number of outliers in the set.

For a method $m$, we define its set of outliers $O_m$ as the $\lceil N\alpha \rceil$ members with the lowest depths, where $\lceil \cdot \rceil$ is the ceiling operator. For the results we report next, we used $\alpha = 0.3$. We compare the outliers of each depth method identified against the ground truth (GT) outliers, which we define as the reference set $O_r$. Table II shows the mean and the standard deviation of the percentage of the outliers each method detected with respect to the GT ones for D2-D6 with $N = 100$.

As the table indicates, except for D2, strict depth methods are more effective at identifying outliers. This result agrees with the functional depth literature, which shows that strict depth methods have a higher sensitivity to outliers [23]. The

| Dataset | CBD (%) | mCBD (%) | ID (%) | eID (%) |
|---|---|---|---|---|
| D2 | $76.16 \pm 13.31$ | $98.12 \pm 4.22$ | $90.08 \pm 7.68$ | $98.12 \pm 4.22$ |
| D3 | $77.54 \pm 14.67$ | $58.94 \pm 13.19$ | $71.46 \pm 14.48$ | $49.89 \pm 14.20$ |
| D4 | $88.14 \pm 15.59$ | $17.43 \pm 14.74$ | $85.07 \pm 14.34$ | $8.06 \pm 7.68$ |
| D5 | $85.21 \pm 16.92$ | $69.27 \pm 8.94$ | $83.37 \pm 18.90$ | $54.54 \pm 7.78$ |
| D6 | $66.11 \pm 9.31$ | $68.19 \pm 16.23$ | $81.46 \pm 13.33$ | $66.52 \pm 18.90$ |

most challenging dataset for mCBD and eID was D4, with inside-shape outliers. Although both methods performed poorly, mCBD did a better job, which potentially indicates that the extra comparisons of CBD endow the method with a higher sensitivity for detecting shape outliers.

As the table indicates, no strict method consistently outperforms the other. ID performed better for the dataset with symmetric magnitude contamination (D2) and topological outliers (D6). In the other cases, CBD achieved better scores. Similarly, except for D4, the performance of modified depth methods was comparable across datasets. These results show how, in practice, the choice of method will depend on the type of data at hand. In agreement with previous literature in band depths [7], [23], the strength of CBD lies in identifying outliers like those in D4, which have a significantly different shape but fall within the ensemble's band envelope.

Finally, we also compare the methods' outlier detection performance qualitatively. The four bottom rows of Fig. 5 present the spaghetti plots with lines colored according to the depths that different methods yield. The figure evidences the similarities between CBD and ID, and mCBD and eID. As expected, CBD and ID methods assign lower depth values to contours that deviate from the ensemble's main trend. CBD and ID produce a wider range of depth values, demonstrated by the color gradient which contains black and bright yellow lines. In contrast, mCBD and eID yield mostly high-depth scores with some contours receiving lower ones. Graphically, this translates to overall brighter color gradients. Despite this visual change, it is possible to observe that the depth-induced rankings of the contours are similar between strict and modified versions.

### D. Experiment 3: Estimator's Robustness

Depth values permit generalizing uni-variate order and rank statistics to the multivariate case. For this experiment, we are interested in the quality of the trimmed mean, which is one of the robust statistics that the contour boxplot visualization uses. To compute the $\alpha$-trimmed mean ($M_m^\alpha$) of an ensemble of contours we average binary masks of the top $N - \lceil N\alpha \rceil$ contours, depthwise, and extract a new contour from the resulting scalar field using 0.5 as iso-value. Specifically, we compute the $\alpha$-trimmed mean contour for method $m$ using the expression

$$M_m^\alpha = \frac{\sum_{i=1}^{N - \lceil N\alpha \rceil} in(c_i)}{N - \lceil N\alpha \rceil}, \qquad (16)$$

| Dataset | $M_S$ | $M_{CBD}^\alpha$ | $M_{mCBD}^\alpha$ | $M_{ID}^\alpha$ | $M_{eID}^\alpha$ |
|---|---|---|---|---|---|
| D1 | $1.42 \pm 0.06$ | $1.17 \pm 0.10$ | $1.13 \pm 0.04$ | $1.15 \pm 0.08$ | $1.12 \pm 0.05$ |
| D2 | $1.77 \pm 0.08$ | $1.47 \pm 0.12$ | $1.32 \pm 0.14$ | $1.37 \pm 0.14$ | $1.31 \pm 0.12$ |
| D3 | $1.51 \pm 0.08$ | $1.26 \pm 0.11$ | $1.20 \pm 0.12$ | $1.24 \pm 0.10$ | $1.18 \pm 0.11$ |
| D4 | $1.46 \pm 0.08$ | $1.24 \pm 0.09$ | $1.14 \pm 0.05$ | $1.22 \pm 0.08$ | $1.13 \pm 0.05$ |
| D5 | $1.50 \pm 0.08$ | $1.24 \pm 0.10$ | $1.17 \pm 0.07$ | $1.24 \pm 0.10$ | $1.17 \pm 0.07$ |
| D6 | $1.60 \pm 0.16$ | $1.48 \pm 0.23$ | $1.17 \pm 0.08$ | $1.24 \pm 0.14$ | $1.15 \pm 0.06$ |

We compute the average $MSE$ across replications for $N = 100$ and include also the standard deviation of the estimates. We multiply both the mean and std by $\times 10^2$ to reduce clutter.

where $in(c_1), \ldots, in(c_{N - \lceil N\alpha \rceil})$ are the binary masks of the inside regions associated with the $N - \lceil N\alpha \rceil$ contours with the highest depth, according to method $m$. In addition to each method's trimmed mean, we also consider the sample mean ($M_S$), which we compute per dataset/replication combination by using (16) without trimming the ensemble. $M_S$ represents a worst-case scenario in which outliers were not removed. For the experiments in this section, we set $\alpha = 0.3$.

A robust trimmed mean is one not affected by outliers. In other words, the trimmed mean contour should be close to the population's average shape. Therefore, to evaluate the depth methods' estimators, we compare them against the binary mask of $f_0$ in (12), which we denote $M_P$. To compare the trimmed means with $M_P$ we compute the mean squared error (MSE) between the masks

$$MSE(M_m^\alpha, M_P^\alpha) = \frac{\sum_{r=0}^{rows} \sum_{c=0}^{cols} [M_m^\alpha(r, c) - M_P^\alpha(r, c)]^2}{rows \times cols}, \qquad (17)$$

where $M_m^\alpha(r, c)$ is the value of the binary array of the trimmed mean $M_i^\alpha$ under consideration at the given row and column.

Table III presents the mean and the standard deviation of the MSE for D1-D6 with the ensemble size $N = 100$. Both CBD and ID methods yield lower average MSE when compared to the sample mean $M_S$. This shows that removing outliers, only considering the most central contours, leads to more robust estimators closer to the population mean $M_P$. In most cases, the mean MSE of $M_\alpha^{CBD}$ is higher than that of $M_\alpha^{ID}$. The same observation holds for the modified versions, which suggests that the outliers ID methods remove contribute more to deviating the trimmed mean from the population estimate. Finally, modified depth methods obtain lower MSE than their strict counterparts. Considering that strict methods performed better at identifying outliers, this result suggests that other contours besides artificially introduced outliers might contribute more towards making the mean estimates less robust. These results show that both CBD and ID methods yield robust mean estimates that are closer to the population estimate than $M_S$.

## VII. VISUAL COMPARISON ON REAL DATA

The previous results demonstrated ID's robustness and more favorable scaling behavior compared to CBD using synthetic data. We now illustrate the use of ID with medical image semantic segmentation and meteorological forecasting datasets.

The contours in these real datasets tend to cross over a lot. Therefore, we focus the analysis on eID, which yields more visually meaningful results in these cases. Unless stated otherwise, we used the same setup for the depth computation methods and ran the analyses in the same machine as in the experiments with synthetic data.

### A. Medical Image Segmentation Ensembles

*Data:* In image-guided medical specialties, clinicians use three-dimensional images of the patient's anatomy to plan the treatment. A core step of the treatment planning process is to segment anatomies of interest like malignancies and the organs-at-risk. With the advent of deep learning-based auto-contouring technologies, this step has been largely automated [30]. Nevertheless, clinicians still need to perform a quality assessment of the segmentations, which requires understanding the uncertainty in the predictions.

We consider the computerized tomography (CT) of a patient with head and neck cancer treated at HollandPTC between 2018 and 2020. The IRB approved the research protocol for the use of patient data in research, all patients signed an informed consent form. For the analysis, we focus on the brain stem and the parotid gland because these structures are not always clearly visible in CT, which can increase inter-clinician variability. In these cases, a visual statistical summary can help clinicians understand the range of predictions. We used a collection of 3D segmentation models based on the popular UNet architecture [31] to generate an ensemble of segmentation predictions of the right parotid gland. Specifically, we trained 30 models on different subsets of the training split of the dataset of the Head and Neck Auto Segmentation MICCAI Challenge [32], a technique known as bootstrapping in the machine learning community. The MICCAI dataset contains CT scans of patients with head and neck cancer with ground truth segmentations of nine organs at risk. To further augment the ensemble size, and the variability of the predictions, we trained each model using different learnable weight initializations. Using the resulting models to segment the parotid gland yields an ensemble of 120 scalar maps of per-voxel softmax probabilities. We extracted the contour ensemble that CBD and ID receive as input by thresholding these arrays with an iso-value of 0.8. For the results below, we computed the depths of the ensemble of contours in 2D $540 \times 540$ pixels slices of the right parotid gland and brain stem segmentation volumes.

*Analysis:* The top row of Fig. 7 visualizes the raw ensemble of contours of the brain stem and parotid gland using spaghetti plots. The variability in the contours of the two structures differs due to changes in visibility in the CT scans. The brain stem shows significantly more variability than the parotid gland, especially on the upper left side, where several contour lines go out of the way of the main shape. The overplotting in the spaghetti plots creates high-density areas that hint at the location of shape representatives. Conversely, one would expect that contours in less dense areas that deviate from the ensemble's main trend correspond to outliers.

Visual statistical summaries remove the need from presenting all ensemble members while still conveying relevant statistical
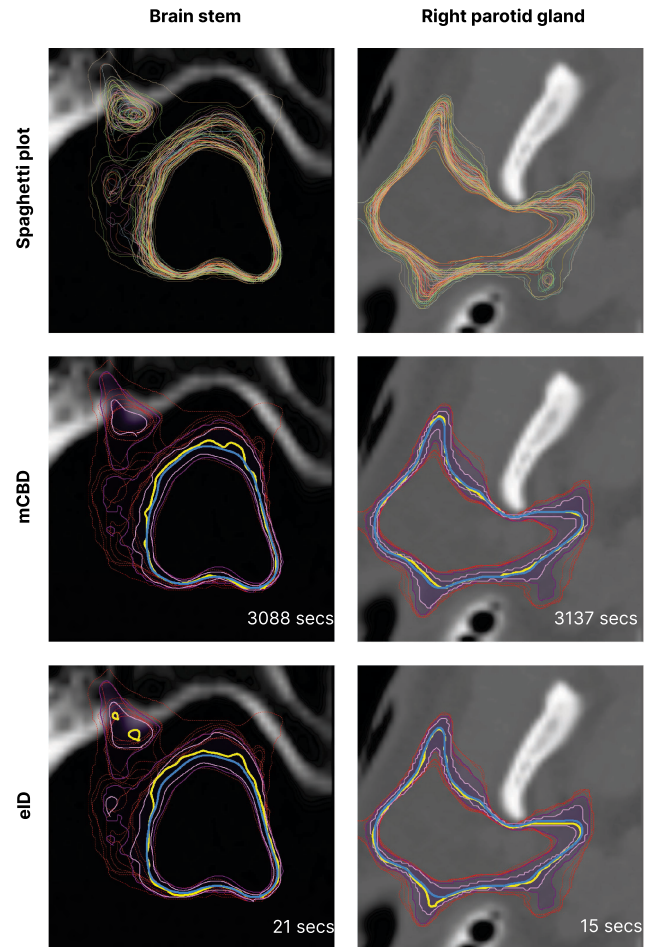


Fig. 7. Contour boxplots that provide a statistical summary of an ensemble of contours of a slice of the brain stem (top row) and right parotid gland (bottom row) of a head-and-neck cancer patient. We generated the contour boxplots using the depths obtained from the mCBD and eID. The yellow and blue lines correspond to the median and mean, respectively. Two bands are depicted in shades of purple as formed by members with the top 50% and 100% depths, not considering outliers, which are shown using dashed red lines. The timings indicate the duration of the contour depth computation process.

features like the representative contours and the ensemble's variability. For each anatomical structure, Fig. 7 presents contour boxplots generated with depths from the mCBD and eID, using $\alpha = 0.1$ for the trimming. The first thing to note is the different runtimes. For a $N = 100$ ensemble, mCBD took more than twenty minutes to compute the depths. In contrast, it took eID seconds. These results show that ID can support larger datasets without requiring special hardware, which increases its practical value.

In terms of the boxplot's statistical features, we start by analyzing the median, depicted as a yellow line. In both cases, the median that mCBD and eID yield is not the same contour. Nevertheless, the contours' shapes are visually similar. When we inspected the depth-induced rankings of the contours, we noticed both medians obtained high depth with both methods, but their ranks varied, which resulted in a different contour being displayed. For instance, for the parotid gland, mCBD's median had the 8th highest depth according to eID. Similarly, eIDs'
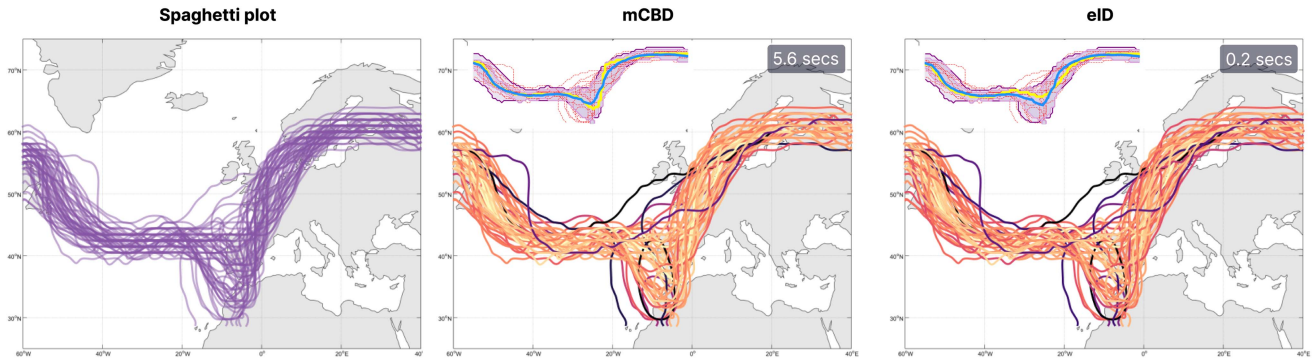
Fig. 8.    mCBD and eID depths for an ensemble of 500 hPa geopotential height contour lines. The inset of each method presents the corresponding contour boxplot with the $N \times 20\% = 10$ contours with the lowest depth set as outliers. We used an opacity of 0.1 for the spaghetti plot. For the depth line plots, we scaled the color scale based on the min and max depth each method yielded. The timings indicate the duration of the contour depth computation process.

median was ranked 8th according to mCBD. We observed a similar trend with the brain stem. The similarity of the rankings induced by mCBD and eID depths can be observed by comparing the method's trimmed means (blue lines) and the outliers (dashed red lines). The figure shows how the means both methods yield are very similar. While the MSE between mCBD's and eID's median contours is 0.013 (brain stem) and 0.011 (parotid gland) pixels, it is 0.0005 (brain stem) and 0 (parotid gland) pixels for the means. In the appendix, available online, we provide a juxtaposed view of the medians and $\alpha$-trimmed means to facilitate visual inspection. The $\alpha$-trimmed means are similar because mCBD's and eID's inlier and outlier sets largely agree. For the right parotid gland, for example, both sets perfectly intersect. Similar to the medians, the boxplots' confidence bands vary across methods. The reason is that although mCBD's and eID's inlier and outlier sets agree, the ordering of the contours is not the same (Pearson correlation coefficient between mCBD and eID depths of 0.96 for the brain stem and 0.98 for the parotid gland). Therefore, the shape of the bands might differ if, for example, a contour that fell in the 50% band for mCBD does not for eID.

## B. Meteorological Forecasting

*Data*: A common use case for contour statistical models is to analyze meteorological forecast data. In this work, we consider data from the European Centre for Medium-Range Weather Forecasts (ECMWF). Specifically, the ECMWF Ensemble Prediction System (EPS) provides ensembles of predictions for different variables like precipitation, temperature, and pressure. The forecasts include $N = 50$ perturbed members and a control run. We analyze the same data as in [5], which is the forecast from 00:00 UTC 15 October 2012. More details about this type of data can be found at [33]. The region under consideration encompasses $101 \times 41 \times 62$ grid points, which corresponds to latitude, longitude, and geopotential height dimensions. For the analysis, we consider 2D fields, corresponding slices of the region where the geopotential height is $500hPa$. To obtain contours from this field, we threshold them using an iso-value

of 5,600 m. The left-most panel of Fig. 8, depicts the extracted contours laid over the geographical region they span.

*Analysis:* The spaghetti plot in the first column of Fig. 8 permits assessing the extent of the ensemble and suggests trends in high-density areas. Furthermore, it permits identifying contour portions lying outside of the ensemble's envelope as potential outliers. The second and third panels of Fig. 8 color the lines using the depth that mCBD and eID assigned to each contour. Darker and brighter colors represent lower and higher values, respectively. Visual inspection reveals the similarity between the methods' outputs. Contours that are deep within the ensemble's predominant shape are brighter, signaling higher depth. In contrast, contours with portions deviating from the ensemble's representative trend are darker. Additionally, the color gradients in both mCBD and eID line plots are similar. Inspecting the depth scores closer revealed a large agreement between the methods' inlier and outlier sets, which had 97% and 90% overlap ($\alpha=0.2$). Nevertheless, similarly to the case of segmentation data, the specific depth values vary (correlation coefficient of 0.96), altering the depth-induced rankings and leading to different medians (yellow) and bands being displayed. The $\alpha$-trimmed means (blue), with an MSE of 0.0021 pixels (compared to the medians' MSE of 0.037 pixels), and outliers (red) evidence the agreement of robust estimates based on mCBD and eID depth scores.

## VIII. DISCUSSION AND CONCLUSION

In this paper, we presented Inclusion Depth (ID), a new depth notion applicable to contour ensembles. The concept of statistical depth permits extending order and rank-based statistics to the multivariate case. Depth-induced orderings allow summarizing the ensemble members in terms of their median, trimmed mean, and confidence bands, and obtaining robust estimators by removing outliers.

ID provides theoretical guarantees on the depth estimates, derived from its relationship with Half-Region Depth. Additionally, based on the simple principle of assessing contours inside/outside relationships, ID is accessible and its results interpretable. Using synthetic data we demonstrated ID's more

favorable $\mathcal{O}(MN^2)$ scaling, compared CBD's $\mathcal{O}(MN^3)$ [7]. The experiments showed that ID and eID are successful at identifying a wide range of outliers and yield robust estimators of the ensemble's mean, comparable to CBD's. These robust estimators enable extending robust statistical theory and analysis to contours. Finally, by applying ID to real datasets, we further demonstrated the method's practical value to analyze contour ensembles when paired with visualizations like contour boxplots.

In the literature, it has been noted that CBD can be accelerated in several ways. First, CBD's outer loop is highly parallelizable, so it could significantly profit from GPU acceleration. In this paper, we did not focus on improvements that could be added on top of the methods. Rather, we propose an alternative depth notion that is asymptotically faster than CBD. Similarly to CBD, ID has a highly parallelizable loop, so this improvement would also benefit ID. Second, in terms of algorithmic improvements, [34] proposes a faster way to compute functional Band Depth. Contours, with the different possible topologies, are not straightforward to adapt to this methodology. Therefore, it remains future work to verify whether these optimizations are possible. Same as with parallelization, it holds that such an improvement would likely benefit both CBD and ID.

The experiments with synthetic data showed that ID and eID detect outliers with comparable performance to CBD across several outlier types. Nevertheless, there is still room for improvement. Particularly in the case of eID, which performed weakly at identifying shape outliers with a magnitude similar to other ensemble members. Improving outlying detection performance constitutes future work. We anticipate that introducing information about the contour's metric structure, similar to [6], could help in cases where inside/outside relationships do not suffice. Second, the eID can assign low non-zero depth scores to outlying contours. mCBD uses an automatic thresholding method that optimizes the ensemble's mean depth to set outliers' depth to zero. This procedure removes the need to find a threshold for the trimming operations via trial and error, like in eID's case. To reduce users' burden, we will investigate options to integrate an automated thresholding procedure similar to mCBD's in our framework.

The improved computational complexity of ID brings within reach the usage of depth-based order and rank statistics for larger datasets in interactive settings. In domains like computer-aided design, simulation, and medical image segmentation, it is common to deal with three-dimensional objects [11]. Our method is quite general and can be applied to three-dimensional contours with ease. Second, currently unimodal distribution is assumed, however, when studying contour's ensembles it is common to first identify the main modes of variation [4], [19], [35]. CBD could make this identification more robust to certain types of outliers [15] but at the cost of reduced interactivity. Using ID instead would permit performing real-time interactive depth-based clustering on larger contour ensembles. Finally, the interactivity that ID unlocks calls for reimagining contour boxplots for interactive scenarios. For instance, it could be possible to change parameters or weights in the depth function and see them reflected in the contour boxplot in real time.

## REFERENCES

[1] D. B. Stephenson and F. J. Dolas-Reyes, "Statistical methods for interpreting Monte Carlo ensemble forecasts," *Tellus A: Dyn. Meteorol. Oceanogr.*, vol. 52, no. 3, pp. 300–322, 2000. [Online]. Available: https://doi.org/10.3402/tellusa.v52i3.12267

[2] M. Abdar et al., "A review of uncertainty quantification in deep learning: Techniques, applications and challenges," *Inf. Fusion*, vol. 76, pp. 243–297, 2021. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S1566253521001081

[3] F. Renard, S. Guedria, N. D. Palma, and N. Vuillerme, "Variability and reproducibility in deep learning for medical image segmentation," *Sci. Rep.*, vol. 10, no. 1, 2020, Art. no. 13724. [Online]. Available: https://doi.org/10.1038/s41598--020-69920-0

[4] M. Zhang, Q. Li, L. Chen, X. Yuan, and J.-H. Yong, "EnConVis: A unified framework for ensemble contour visualization," *IEEE Trans. Vis. Comput. Graphics*, vol. 29, no. 4, pp. 2067–2079, Apr. 2023.

[5] F. Ferstl, M. Kanzler, M. Rautenhaus, and R. Westermann, "Visual analysis of spatial variability and global correlations in ensembles of iso-contours," *Comput. Graph. Forum*, vol. 35, no. 3, pp. 221–230, 2016. [Online]. Available: https://onlinelibrary.wiley.com/doi/abs/10.1111/cgf.12898

[6] I. Demir, M. Jarema, and R. Westermann, "Visualizing the central tendency of ensembles of shapes," in *Proc. SIGGRAPH ASIA Symp. Vis.*, New York, NY, USA, 2016, Art. no. 3. [Online]. Available: https://doi.org/10.1145/3002151.3002165

[7] R. T. Whitaker, M. Mirzargar, and R. M. Kirby, "Contour boxplots: A method for characterizing uncertainty in feature sets from simulation ensembles," *IEEE Trans. Vis. Comput. Graphics*, vol. 19, no. 12, pp. 2713–2722, Dec. 2013.

[8] M. Rautenhaus et al., "Visualization in meteorology–A survey of techniques and tools for data analysis tasks," *IEEE Trans. Vis. Comput. Graphics*, vol. 24, no. 12, pp. 3268–3296, Dec. 2018.

[9] M. Mirzargar and R. T. Whitaker, "Representative consensus from limited-size ensembles," *Comput. Graph. Forum*, vol. 37, no. 3, pp. 13–22, 2018. [Online]. Available: https://onlinelibrary.wiley.com/doi/abs/10.1111/cgf.13397

[10] P. Voglreiter et al., "RFA guardian: Comprehensive simulation of radiofrequency ablation treatment of liver tumors," *Sci. Rep.*, vol. 8, no. 1, 2018, Art. no. 787. [Online]. Available: https://doi.org/10.1038/s41598--017-18899-2

[11] M. Raj, M. Mirzargar, J. S. Preston, R. M. Kirby, and R. T. Whitaker, "Evaluating shape alignment via ensemble visualization," *IEEE Comput. Graph. Appl.*, vol. 36, no. 3, pp. 60–71, May/Jun. 2016.

[12] K. Mosler, *Depth Statistics*. Berlin, Germany:Springer, 2013, pp. 17–34. [Online]. Available: https://doi.org/10.1007/978--3-642-35494-6_2

[13] S. López-Pintado and J. Romo, "A half-region depth for functional data," *Comput. Statist. Data Anal.*, vol. 55, no. 4, pp. 1679–1695, 2011. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S0167947310004123

[14] C. Gao, "Robust regression via mutivariate regression depth," *Bernoulli*, vol. 26, no. 2, pp. 1139–1170, 2020. [Online]. Available: https://doi.org/10.3150/19-BEJ1144

[15] R. Jörnsten, "Clustering and classification based on the L1 data depth," *J. Multivariate Anal.*, vol. 90, no. 1, pp. 67–89, 2004. [Online]. Available: https://www.sciencedirect.com/science/article/pii/S0047259X04000272

[16] J. Wang, S. Hazarika, C. Li, and H.-W. Shen, "Visualization and visual analysis of ensemble data: A survey," *IEEE Trans. Vis. Comput. Graphics*, vol. 25, no. 9, pp. 2853–2872, Sep. 2019.

[17] J. Sanyal, S. Zhang, J. Dyer, A. Mercer, P. Amburn, and R. Moorhead, "Noodles: A tool for visualization of numerical weather model ensemble uncertainty," *IEEE Trans. Vis. Comput. Graphics*, vol. 16, no. 6, pp. 1421–1430, Nov./Dec. 2010.

[18] A. Kumpf, B. Tost, M. Baumgart, M. Riemer, R. Westermann, and M. Rautenhaus, "Visualizing confidence in cluster-based ensemble weather forecast analyses," *IEEE Trans. Vis. Comput. Graphics*, vol. 24, no. 1, pp. 109–119, Jan. 2018.

[19] F. Ferstl, K. Bürger, and R. Westermann, "Streamline variability plots for characterizing the uncertainty in vector field ensembles," *IEEE Trans. Vis. Comput. Graphics*, vol. 22, no. 1, pp. 767–776, Jan. 2016.

[20] K. Pothkow and H.-C. Hege, "Positional uncertainty of isocontours: Condition analysis and probabilistic measures," *IEEE Trans. Vis. Comput. Graphics*, vol. 17, no. 10, pp. 1393–1406, Oct. 2011.

[21] T. Athawale, E. Sakhaee, and A. Entezari, "Isosurface visualization of data with nonparametric models for uncertainty," *IEEE Trans. Vis. Comput. Graphics*, vol. 22, no. 1, pp. 777–786, Jan. 2016.

[22] K. Pöthkow and H.-C. Hege, "Nonparametric models for uncertainty visualization," *Comput. Graph. Forum*, vol. 32, no. 3pt2, pp. 131–140, 2013.

[23] S. López-Pintado and J. Romo, "On the concept of depth for functional data," *J. Amer. Stat. Assoc.*, vol. 104, no. 486, pp. 718–734, 2009. [Online]. Available: https://doi.org/10.1198/jasa.2009.0108

[24] R. Y. Liu, "On a notion of data depth based on random simplices," *Ann. Statist.*, vol. 18, no. 1, pp. 405–414, 1990. [Online]. Available: http://www.jstor.org/stable/2241550

[25] J. W. Tukey, "Mathematics and the picturing of data," in *Proc. Int. Congr. Math.*, Vancouver, 1975, pp. 523–531. [Online]. Available: https://cir.nii.ac.jp/crid/1573950399770196096

[26] R. Serfling and Y. Zuo, "General notions of statistical depth function," *Ann. Statist.*, vol. 28, no. 2, pp. 461–482, 2000. [Online]. Available: https://doi.org/10.1214/aos/1016218226

[27] Y. Sun and M. G. Genton, "Functional boxplots," *J. Comput. Graphical Statist.*, vol. 20, no. 2, pp. 316–334, 2011. [Online]. Available: https://doi.org/10.1198/jcgs.2011.09224

[28] T. Jech, *Set Theory*, 3rd ed., Berlin, Germany: Springer, 2003.

[29] C. R. Harris et al., "Array programming with NumPy," *Nature*, vol. 585, no. 7825, pp. 357–362, Sep. 2020. [Online]. Available: https://doi.org/10.1038/s41586--020-2649-2

[30] E. Montagnon et al., "Deep learning workflow in radiology: A primer," *Insights Imag.*, vol. 11, no. 1, pp. 1–15, 2020. [Online]. Available: https://doi.org/10.1186/s13244--019-0832-5

[31] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *Proc. Med. Image Comput. Comput.-Assisted Intervention*, N. Navab, J. Hornegger, W. M. Wells, and A. F. Frangi, Eds. Cham:Springer International Publishing, 2015, pp. 234–241.

[32] P. F. Raudaschl et al., "Evaluation of segmentation methods on head and neck CT: Auto-segmentation challenge 2015," *Med. Phys.*, vol. 44, no. 5, pp. 2020–2036, 2017. [Online]. Available: https://aapm.onlinelibrary.wiley.com/doi/abs/10.1002/mp.12197

[33] M. Leutbecher and T. Palmer, "Ensemble forecasting," *J. Comput. Phys.*, vol. 227, no. 7, pp. 3515–3539, 2008.

[34] Y. Sun, M. G. Genton, and D. W. Nychka, "Exact fast computation of band depth for large functional datasets: How quickly can one million curves be ranked?," *Stat*, vol. 1, no. 1, pp. 68–74, 2012.

[35] B. Ma and A. Entezari, "An interactive framework for visualization of weather forecast ensembles," *IEEE Trans. Vis. Comput. Graphics*, vol. 25, no. 1, pp. 1091–1101, Jan. 2019.

**Marius Staring** is a professor and vice director of LKEB with the Leiden University Medical Center. He and his team develop generic machine-learning approaches for automated image analysis and apply these in the clinical and life sciences. He is an associate editor of *IEEE Transactions on Medical Imaging*, and a member of program committees of MICCAI, IEEE ISBI, SPIE MI, and WBIR. Open-sourcing his methods has been a common theme in his career, exemplified by the image registration package Elastix, see https://elastix.lumc.nl/.

**René van Egmond** is an associate professor of cognitive ergonomics with the Faculty of Industrial Design Engineering, Delft University of Technology, Delft, The Netherlands. His expertise lies in the fields of Product Sound Design & Perception and Informational Ergonomics. His research is focused on understanding how people process information streams in complex environments and how people process this information and deal with this complexity.

**Anna Vilanova** is a full professor in visual analytics with the Department of Mathematics and Computer Science, Eindhoven University of Technology. She is leading a research group on the subject of visual analytics focusing on high dimensional data, explainable AI, and medical visualization for complex imaging data. She has been a member of international program committees, chair, and editor of conferences & journals in visualization. She has been an elected member of the EUROGRAPHICS executive committee and vice president of EUROGRAPHICS.

**Nicolas F. Chaves-de-Plaza** is currently working toward the PhD degree with the Department of Intelligent Systems, Delft University of Technology, Delft, The Netherlands. He works on developing visualization and interaction tools to support clinician-driven segmentation of 3D medical images in the context of Adaptive Proton Therapy.

**Prerak Mody** is currently working toward the PhD degree with the Division of Image Processing (Dutch abbreviation LKEB), Leiden University Medical Center. His research focuses on clinically-applicable Bayesian and interactive deep learning techniques.

**Klaus Hildebrandt** is a tenured assistant professor with the Computer Graphics and Visualization Group, Department of Intelligent Systems, Delft University of Technology, Delft, The Netherlands. His research interests include visual computing, geometric data processing, physical simulation, and computational and discrete differential geometry.