# Automated classification of single airborne particles from two-dimensional angle-resolved optical scattering (TAOS) patterns by non-linear filtering ☆

Giovanni Franco Crosta [a,*], Yong-Le Pan [b], Kevin B. Aptowicz [c], Caterina Casati [a], Ronald G. Pinnick [b], Richard K. Chang [d], Gorden W. Videen [b]

[a] Department of Environmental & Earth Sciences, University of Milan Bicocca 1, Piazza della Scienza, 20126 Milan, Italy
[b] U.S. Army Research Laboratory, Adelphi, MD 20783, USA
[c] Department of Physics, West Chester University, West Chester, PA 19383, USA
[d] Department of Applied Physics, Yale University, New Haven, CT 06511, USA

## ABSTRACT

Measurement of two-dimensional angle-resolved optical scattering (TAOS) patterns is an attractive technique for detecting and characterizing micron-sized airborne particles. In general, the interpretation of these patterns and the retrieval of the particle refractive index, shape or size alone, are difficult problems. By reformulating the problem in statistical learning terms, a solution is proposed herewith: rather than identifying airborne particles from their scattering patterns, TAOS patterns themselves are classified through a learning machine, where feature extraction interacts with multivariate statistical analysis.

Feature extraction relies on *spectrum enhancement*, which includes the discrete cosine FOURIER transform and non-linear operations. Multivariate statistical analysis includes computation of the principal components and supervised training, based on the maximization of a suitable figure of merit. All algorithms have been combined together to analyze TAOS patterns, organize feature vectors, design classification experiments, carry out supervised training, assign unknown patterns to classes, and fuse information from different training and recognition experiments. The algorithms have been tested on a data set with more than 3000 TAOS patterns. The parameters that control the algorithms at different stages have been allowed to vary within suitable bounds and are optimized to some extent.

Classification has been targeted at discriminating aerosolized *Bacillus subtilis* particles, a simulant of *anthrax*, from atmospheric aerosol particles and interfering particles, like diesel soot. By assuming that all training and recognition patterns come from the respective reference materials only, the most satisfactory classification result corresponds to 20% false negatives from *B. subtilis* particles and $< 11\%$ false positives from all other aerosol particles. The most effective operations have consisted of thresholding TAOS patterns in order to reject defective ones, and forming training sets from three or four pattern classes. The presented automated classification method may be adapted into a real-time operation technique, capable of detecting and characterizing micron-sized airborne particles.

© 2013 The Authors. Published by Elsevier Ltd. All rights reserved.

* Corresponding author. Tel.: +39 02 6448 2724, office: +39 02 6448 2754.
E-mail addresses: Giovanni_Crosta@uml.edu, crgvfa@yahoo.com (G.F. Crosta).

## 1. Introduction

The real-time detection, sizing and speciation of airborne particulate matter in the respirable size range (from 0.5 to 10 μm) are of great importance for civil, industrial and military purposes. These applications include monitoring of indoor or outdoor air, early warning of aerosolized biological threat, quality control of inhalable pharmaceuticals, characterization of pigment dispersions, in-flight detection of ice microcrystals and volcanic ash, and the physical and chemical analysis of tropospheric aerosol to predict cloud formation. In the most basic case, one would like to be able to detect single airborne particles in a prescribed size range, measure some physical and/or chemical properties in real time, and make decisions accordingly.

Experimental techniques based on single particle optical sensing (SPOS) have been known for a long time [1]. TAOS is a measure of the two-dimensional elastic-light-scattering patterns from single airborne particles, developed by a few research groups over the last two decades [e.g., 2–4]. A single moving particle is illuminated by a triggered laser beam, and the resulting scattering intensity pattern is stored for subsequent processing. Significant progress has been made on the instrumentation side, from particle sampling to collection optics, and pattern-capture devices. Even though TAOS patterns contain information related to particle size, refractive index, shape, surface roughness and other properties, analyzing and understanding the TAOS patterns is still a major open problem. In principle it could be dealt with as an inverse obstacle problem [5,6].

Generally, two strategies can be pursued to analyze the observed TAOS patterns. Strategy 1 could be called model-driven analysis and understanding. The MAXWELL equations may be solved according to what is known about the particle, such as its size range, shape, refractive index, and orientation in the laboratory reference frame, in order to produce a scattering pattern. The latter is compared with experimental results [7,8]. If the comparison is satisfactory according to some figure of merit, then an attempt is made to infer the properties of the unknown scatterer from its TAOS data. In the simplest setting, this may consist of looking up a library of computed patterns. Such a computation can be carried out on simple shapes (e.g. spherically capped cylinders, prolate spheroids, pairs of tangent equal spheres) [9].

By contrast, strategy 2 is data-driven analysis. No *a priori* theoretical scattering model is assumed to be applicable. Experimental patterns are sorted by classes according to some criteria, i.e., by their origin (particulate material, range of sizes), or by their morphological properties (intensity peaks and valleys, statistical moments, etc.) [3]. In the latter case, if simulations are available, instead of requiring a perfect match between observation and calculation, a computed pattern is assimilated to an experimental class [10,11]. The discovery of hidden relations in the data set is handed over to a heuristic procedure, e.g., classification. A classifier is a "machine" which implements a learning algorithm, followed by a decision making algorithm, which assigns a datum to a class. Artificial learning includes the design, training and validation of a classifier. The ultimate goal is automated pattern classification. Herewith the scope is limited to the assignment of a new TAOS pattern to one of the known classes. Early implementations of such a classifier were described years ago [12,13] and shown to work well, although only applied to a relatively small number (tens) of patterns.

This paper reports the most recent progress dictated by strategy 2. The classifier described here relies on two interacting modules: feature extraction by the so-called spectrum enhancement algorithm, and linear classification by multivariate statistical analysis.

## 2. Measurement of TAOS patterns

TAOS instrumentation has been under development at the Department of Applied Physics, Yale University and at the US Army Research Laboratory (ARL) for approximately 15 years. From the early versions of the 1990s, the design and implementation of TAOS instrumentation has advanced to the point of being field-deployable and fit for outdoor or indoor air sampling [14]. The TAOS patterns from single individual atmospheric aerosol particles analyzed in this paper were measured by the same system as described previously [14].

A simplified schematic of the experimental arrangement is shown in Fig. 1(a), which is now briefly described. Atmospheric aerosols are drawn at a rate of 770 l/min through a metal duct into the laboratory and concentrated by a virtual impact concentrator (*Dycor* model *XMX*). The aerosol stream from the minor exit of the concentrator is drawn into Airtight chamber at 1 l/min for TAOS measurement. Particles in the aerosol jet are focused by Sheath nozzle and pass through the focal point $F_1$ of Ellipsoidal mirror (*Opti-Forms Inc.*, model *E64-3*) with major axis of 3.35 in. and eccentricity of 0.75. Once a particle arrives at $F_1$, a trigger signal is emitted by a detection subsystem, not shown in Fig. 1(a), and a pulse from the Nd:YAG 532 nm LASER (*Spectra Physics* model *X-30*, 30 ns pulse width) is emitted to illuminate the particle. The elastically scattered light originating from the particle at $F_1$ is reflected by Ellipsoidal mirror, passes through a quartz window (not shown) and is refocused to the right focal point, $F_2$. Iris at $F_2$ removes stray light. Rays from the virtual particle at $F_2$, which give the spatial distribution of scattered light, are projected by Lens into a plane pattern and recorded by an intensified charge-coupled device (ICCD, *Andor Technology Istar DH734-25F-03*). The latter is turned on by the same trigger signal as the main LASER. Downstream $F_1$, Prism deflects the main LASER beam to Beam dump, thus preventing damage to ICCD.

The coordinate frame of the experiment is shown in Fig. 1(b). The origin is at $F_1$. The illuminating laser beam determines the $z$-axis, which coincides with the major axis of the ellipsoidal mirror. By taking into account the geometry of the system, every pixel of the $1024 \times 1024$ pixel ICCD detector corresponds to a unique pair of scattering angles, namely colatitude, $\theta$, and azimuth, $\varphi$. This arrangement collects scattered light within a large solid angle $\{75° \leq \theta \leq 135°\} \times \{0° \leq \varphi \leq 360°\}$, where $\theta = 0$ corresponds to forward scattering.
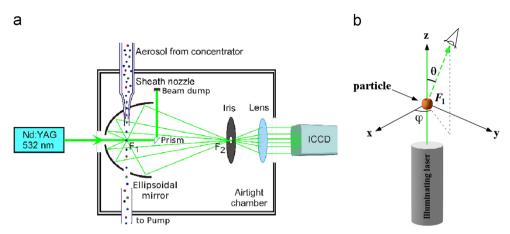
**Fig. 1.** (a) Simplified schematic of the experimental arrangement for TAOS pattern measurement from atmospheric aerosol particles and (b) the corresponding coordinate frame.
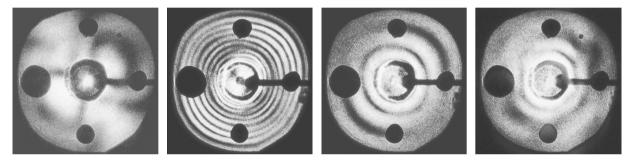


**Fig. 2.** Typical TAOS patterns, from (left to right) single *Bacillus subtilis* spore (class *Bq*), a dioctyl phthalate (DOP) droplet (*Dq*), a dried polystyrene latex (PSL) sphere (*Pq*), and a diesel soot aggregate (class *sq*). Patterns of such classes were used to train the classifier. Contrast enhanced by the Equalize command of GIMP for display purposes only.

Airborne particles of different materials can be artificially produced. For example, dried polystyrene latex (PSL) spheres and dioctyl phthalate (DOP) droplets were aerosolized by a nebulizer (*Royco Aerosol Generator*, model *256*) and used to calibrate the system. TAOS patterns from single *Bacillus subtilis* spore or aggregates of *B. subtilis* spores were also recorded.

Fig. 2 shows typical TAOS patterns of airborne particles from the classes of *B. subtilis* spores (class *Bq*), a dioctyl phthalate (DOP) droplet (*Dq*), a dried polystyrene latex (PSL) sphere (*Pq*), and a diesel soot aggregate of (*sq*), all used for classifier training, as explained in Sections 4.2 and 4.3. Left to right: pattern $\mathcal{N}$. *Bq*114 from a single *B. subtilis* spore with its typical bowtie; pattern $\mathcal{N}$. *Dq*010 from a 2.8 μm dioctyl phthalate (DOP) droplet, with narrow rings; pattern $\mathcal{N}$. *Pq*001 from a 1.03 μm dried polystyrene latex (PSL) sphere with broad rings; pattern $\mathcal{N}$. *sq*261 from a diesel soot aggregate, not showing any significant feature.

The five black round spots included in the pattern are images of the five holes drilled in the ellipsoid. The black horizontal stripe on the right of the pattern centre is the shadow of the prism holder (Fig. 1(a)).

As it appears from the shading symmetry of the scattering patterns of Fig. 2, the polarization plane of the incident wave is at 45°, clockwise with respect to the horizontal, encoded by the suffix −*q* in some class names. The outer boundary corresponds to $\theta = 75°$, the innermost circumference to $\theta = 135°$. These patterns were verified to be in good qualitative agreement with numerical simulations based on LORENZ–MIE theory [7,8].

The variety of patterns recorded from outdoor aerosol sampling is shown by Fig. 3. It would be attractive to divide such patterns into an adequate number of morphological classes, by introducing e.g., a representation *à la* DAVID MARR [15]. The task is beyond the scope of the present work. Instead, the eventual goal herewith is the discrimination of *Bq* patterns produced by single bacterial spores, from those produced by all other particles. The latter patterns include those from diesel soot aggregates (*sq*) and from general airborne particles (*K*0 to *K*5).

## 3. Classifiers

Machine learning and automated classification are not yet common in light-scattering analyses: this section provides some background of classifiers and their application to TAOS patterns.

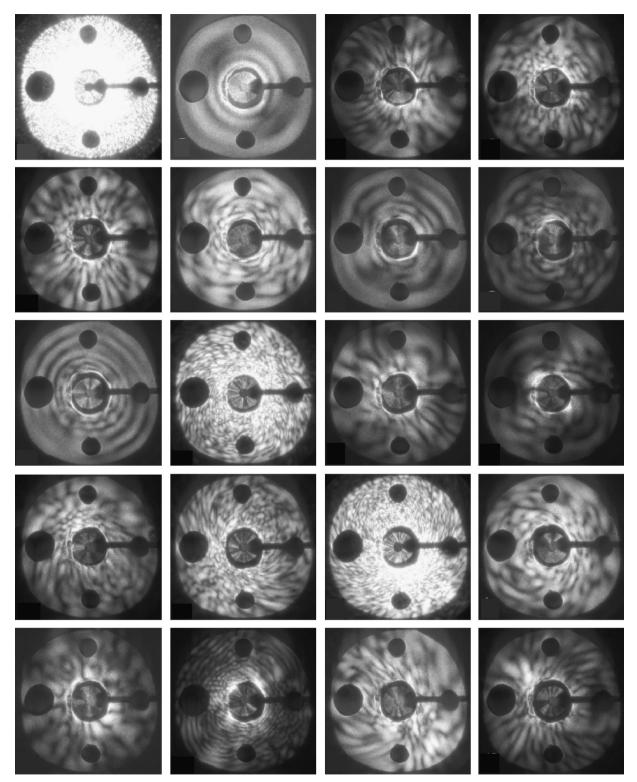**Fig. 3.** The morphological variety of TAOS patterns from outdoor sampling (classes $K0$ to $K5$).

### 3.1. Outline of machine learning

Quoting V.N. Vᴀᴘɴɪᴋ [16], statistical learning is formalised by the minimization of a suitably defined risk functional on the basis of empirical data. The components of statistical learning are: a generator, a supervisor, a learning machine and a loss function.

**Table 1**
List of the backward *TAOS* pattern data sets.

| Material | Pol. (deg) | Directory | Label | From | to | Role |
|----------|-----------|-----------|-------|------|------|------|
| Background | 45 | *qBackground_100* | *Gq* | 001 | 100 | Thresholding |
| *B. subtilis* spores | 45 | *qDugway-BG_098* | *Bq* | 103 | 200 | Reference |
| Dioctyl phthal. | 45 | *qDioctylpht_200* | *Dq* | 001 | 200 | Reference |
| PSL sphere | 45 | *qPS1-1034 nm_039* | *Pq* | 001 | 039 | Reference |
| PSL sphere | 45 | *qPS2-1034 nm_147* | *Qq* | 001 | 147 | None |
| Atmospheric aerosols | | | | | | |
| Set A | 45 | *2004Okt06_Ambient-Pol45-A* | *K0* to *K2* | 0001 | 2500 | Recognition |
| Set B | 45 | *2004Okt06_Ambient-Pol45-B* | *K2* to *K5* | 2501 | 5000 | Recognition |
| Set C | 45 | *q2004Okt06K5993* | *K5* | 5001 | 5993 | Recognition |
| Diesel soot | 45 | *qDieselSoot_288* | *sq* | 001 | 288 | Interfering a.k.a. confounder |

Detailed presentations of the statistical [16] and mathematical theory of learning [17] are available. In the TAOS-specific context, the generator is the source of empirical data, the set $\{g\}$ of TAOS patterns, each of them occurring with unknown probability density $f_G[g]$. Given a pattern $x$, the supervisor returns an attribute, for example the pattern's class of belonging $B$ $[g]$, with unknown conditional probability $f_S[B|g]$. Classes of belonging are provided by Table 1. Machine learning is formalized by the synthesis of a function $A[.]$, which, evaluated on a pattern, returns the assigned class $A[g]$. In turn, as it will be explained below, $A[.]$ is the result of feature extraction (Section 3.3), followed by linear classification (Section 3.4). The goal of learning is to synthesize a machine which minimizes losses, or errors. In this context, the goal has been restated as the maximization of correct class assignments i.e., the occurrences of $A[g]=B[g]$, according to the figure of merit defined by Eq. (23) below. Maximization makes sense, because the class assignment function $A[.]$ depends on a finite number of control parameters, which form the vector $\overrightarrow{\psi}$ of Eq. (11), and because $\overrightarrow{\psi}$ is allowed to vary within prescribed bounds in the finite set $\Psi_{ad}$ (Eq. (15)). Once the machine has learnt, it can assign any pattern to a class. In other words, the learning machine here is a classifier. Since $B[g]$ is known, the training procedure implemented herewith is an instance of supervised learning.

### 3.2. How is a classifier designed

*Samples, classes as subsets, classes as attributes*: A sample, $\mathfrak{S} = \{g\}$, is a finite set of elements. Hereinafter, $g$ stands for a TAOS pattern, or a region-of-interest thereof. A class is a subset $\mathfrak{C} \subset \mathfrak{S}$ of elements having at least one property in common, summarized by the attribute $C$, either because of origin or as a consequence of judgment.

In supervised training, the class of belonging of $g$ is defined by the attribute $B[g]$, regarded as the value returned by the function $B[.]$ evaluated at $g$. In turn, $B[.]$ takes values in the finite set $\mathfrak{B} = \{B_1, B_2, ..., B_T\}$. For each $g$, $B[g] \in \mathfrak{B}$ is known to the supervisor, not to the machine. The latter, given $g$, computes the assigned class $A[g] \in \mathfrak{B}$ i.e., the value of $A[.]$ at $g$ shall also be in $\mathfrak{B}$, but may differ from $B[g]$.

*What shall be submitted to the classifier*?: The choice affects the representation of a pattern for classification. Herewith, TAOS patterns are represented by suitably derived features. An algorithm, a.k.a. the feature extraction module, has been used for this purpose. Features form a real-valued vector (Section 3.3) that encodes morphology-related information and depends on the above mentioned control parameters.

*What is the structure of the classifier*?: Two interacting modules characterize the classifier, the feature extraction and the linear classification modules.

*How are features extracted*?: In this work, features have been obtained by means of the spectrum enhancement algorithm [18,19]. Basically, the spatial derivatives of integer or fractional order of the image are evaluated and non-linear transformations are applied in the Fourier domain (a.k.a. the reciprocal domain). Details of the algorithm are provided in Section 3.3.

*How is the classifier made to work*?: In general, a classifier has to be trained, then validated and eventually applied.

*How does supervised training occur*?: From each class of belonging, a few patterns are selected and a training set is formed. One starts with an *n*-tuple of control parameters. The corresponding features are extracted from all patterns. The linear classifier performs multivariate statistical analysis of feature vectors and forms a *P*-dimensional linear space $\mathcal{Z}$ accordingly (Eqs. (19) and (20)). The assignment of a pattern to a class is based on computation of distances in $\mathcal{Z}$ (Eq. (22)). A pattern $g$ of known class $B[g]$ is assigned class $A[g]$. Classification is rated by a figure of merit such as $F[\mathbf{M}]$ of Eq. (23) below. The *n*-tuple is changed and class assignment is repeated. In other words, training consists of interaction between the feature extraction and the linear classification modules, driven by $F[\mathbf{M}]$.

*How is reliability attained*?: During the validation of supervised training, a class is assigned to other patterns, the *B*-class of which is known, but which has not affected the formation of $\mathcal{Z}$. Class assignment is rated by another $F[\mathbf{M}]$, and the result is taken into account. Furthermore, different training sets are formed (Eq. (24)) and the overall classification is eventually rated (Eq. (26)).

*How is the classifier eventually applied*?: The $n$-tuple of parameters that perform best with all training sets under test is accepted for assigning a class to previously unused patterns, the $B$-class of which is either known or unknown. This is the recognition stage. Since training by each set creates its own linear space $\mathcal{Z}$, and class assignment occurs there independently, one needs a rule to assemble the recognition results into a unique index for each recognized pattern: this leads to information fusion (Section 3.4.8 below).

A detailed description of the feature extraction and linear classification modules follows next.

## 3.3. Feature extraction through enhancement of the spatial frequency spectrum

Let $\Omega$ denote a square of sidelength $L/2$, hereinafter called a "tile." The grayscale image $g[.]$, at which the enhanced spectrum algorithm works, is a (scalar valued) function supported in $\Omega$. The domain in which a *TAOS* pattern is supported is not a square. Therefore, a region of interest (Section 4.3) has to be cut out of a raw pattern and transformed into a tile. The function supported in the region of interest will be denoted by $g[.]$. It represents the whole pattern, as far as classification is concerned.

By spectrum enhancement [19], $\sigma\eta$ for short, one understands any sequence of linear and non-linear operations in the spatial frequency domain: from the power spectral density of the pattern a reference function, which plays the role of a model, is subtracted and the difference function is further processed.

Let $\mathcal{Q}\Omega$ denote a square of sidelength $L$, and $\mathcal{T}$ represent the surface of the torus obtained by deforming $\mathcal{Q}\Omega$ and affixing its opposite sides together. Let $\vec{x} \equiv \{x_1, x_2\} \in \mathcal{Q}\Omega$ be the position vector in the direct domain and $\vec{u} \equiv \{u_1, u_2\}$ be the spatial frequency vector in the reciprocal domain. Let $\mathcal{Q}g[\vec{x}]$ represent a scalar function that is continuous on $\mathcal{T}$. One way of obtaining such a $\mathcal{Q}g[\vec{x}]$ from $g[.]$ supported in $\Omega$ is the application of the twofold $\mathcal{Q}[.] = \text{flop}[\text{flip}[.]]$ reflection. Next let $\mathcal{Q}\Omega$ be discretized by a square grid of step-length $\ell$. The FOURIER transform $G[\vec{u}]$ of $\mathcal{Q}g[\vec{x}]$ is supported at grid nodes in the square

$$0 \leq |u_1|, |u_2| \leq u_{max} = \frac{L}{2\ell} - 1 \text{ cycles/image} \tag{1}$$

and is distribution-valued. In particular $G[\mathbf{O}] = a_{0,0}\delta[\mathbf{O}]$, where obviously $|a_{0,0}| > 0$ for any non-degenerate image. As a consequence of the continuity of $\mathcal{Q}g[\vec{x}]$ on $\mathcal{T}$, the graph of $G[\vec{u}]$ exhibits no cross artefact [20, Chapter 4].

Let $\vec{u}$ be represented in polar coordinates $u \equiv \{u, \vartheta\}$, where $u = |\vec{u}|$ is wavenumber and $\vartheta$ (not to be confused with colatitude, $\theta$) the polar angle such that $0 \leq \vartheta \leq 2\pi$. The power spectral density is denoted by $|G[\vec{u}]|^2$. If $\Theta := [\vartheta_L, \vartheta_H]$, with $\vartheta_L < \vartheta_H$, stands for an arc of radius $u$, then its length, $|\Theta|$, is $|\Theta| = (\vartheta_H - \vartheta_L)u \leq \pi u$.

**Definition** (*Arc-averaged spectral density*). The *normalized, arc-averaged spectral density profile* is the function $s[.]$ of $u$ alone defined in $0 \leq u \leq u_{max}$ (cycles/image) according to

$$s[u] = \frac{1}{|\Theta|} \int_\Theta 10 \, \text{Log}_{10}\left[\frac{|G[\vec{u}]|^2}{|a_{0,0}|^2}\right] u \, d\vartheta, \tag{2}$$

where the integral reduces to a finite sum over the grid nodes belonging to $\Theta$.

Let $m[u]$ be a *model spectral density*. For example, one can choose the continuous function parameterized by $p$ and defined by

$$m^{(p)}[u] := 0, 0 \leq u \leq 1, \quad m^{(p)}[u] := -10 \, \text{Log}_{10}[u^p], \quad u \geq 1 \text{ cycles/image}, \tag{3}$$

where $p$ ($> 0$) is the *model exponent*.

**Definition** (*Log-enhanced spectrum*). The $m^{(p)}$-*log enhanced spectrum* $h^{(p)}[u]$ is defined by

$$h^{(p)}[u] := s[u] - m^{(p)}[u], \quad 0 \leq u \leq u_{\max}. \tag{4}$$

The above operations are algebraic sums of logarithms. Since $s[0] = 0$ and $m^{(p)}[0] = 0$, $\forall p$, then $h^{(p)}[.]$ complies with $h^{(p)}[0] = 0$. Intuitively, the function $h^{(p)}[.]$ represents deviations of $s[.]$ from the model $m^{(p)}$. The values of $L$, $u_{\max}$, $|\Theta|$ and $p$ are among the control parameters of the classifier.

**Definition** (*Enhanced spectrum*). The $m^{(p)} -$ (plain) *enhanced power spectral density* is defined by

$$H^{(p)}[\vec{u}] := |\vec{u}|^p \frac{|G[\vec{u}]|^2}{|a_{0,0}|^2} + \delta[\vec{u}], \tag{5}$$

where $|\vec{u}|^p := (u_1^2 + u_2^2)^{p/2}$.

The main result pertaining to $\sigma\eta$ consists of the following.

**Theorem** (*Spectrum enhancement and spatial differentiation of integer order [19]*). *Assume the image is not degenerate and that all partial derivatives of g[.] up to a suitable order exist as tempered distributions:*

(a) *If the model exponent satisfies $p/2 = N(> 0)$, integer, then $H^{(p)}[\vec{u}]$ has the representation*

$$H^{(p)}[\vec{u}] = \frac{1}{|a_{0,0}|^2} \sum_{n=0}^{N} \binom{N}{n} \left| \mathcal{F}\left[ \frac{\partial^N \mathcal{Q}g}{\partial^{(N-n)}x_1 \partial^n x_2} \right] \right|^2 + \delta[\vec{u}]. \tag{6}$$

(b) *Let $p/2 = N + \lambda$ such that $N + 1 \in \mathbf{N}$, $0 < \lambda < 1$. Then*

$$H^{(p)}[\vec{u}] = \frac{(u_1^2 + u_2^2)^\lambda}{|a_{0,0}|^2} \sum_{n=0}^{N} \binom{N}{n} \left| \mathcal{F}\left[ \frac{\partial^N \mathcal{Q}g}{\partial^{(N-n)}x_1 \partial^n x_2} \right] \right|^2 + \delta[\vec{u}]. \tag{7}$$

(c) *In either case, if all* Fourier *coefficients satisfy*

$$|a_{l,m}| \geq \epsilon > 0 \tag{8}$$

*then the relation between $H^{(p)}[.]$ and $h^{(p)}[.]$ is*

$$h^{(p)}[u] = \frac{10}{|\Theta|} \int_{\Theta} \text{Log}_{10}[H^{(p)}[\vec{u}]] u \; d\vartheta. \tag{9}$$

*The relation between spectrum enhancement and the spatial differentiation of fractional order when $p/2$ is not an integer also can be shown* [19, Theorem 2].

Unlike linear regression of the Log $|G|$ vs. Log$[u]$ plot used in fractal analysis [20], $\sigma\eta$ serves to separate structure from texture as suggested e.g., by the Osher–Rudin paradigm [21]. $\sigma\eta$ places the emphasis on the low-frequency behavior of $h^{(p)}[.]$ which corresponds to image structure.

A polynomial $q[.]$ of suitable degree $d$ is obtained from $h^{(p)}[.]$ via singular value decomposition [22]. The support of $q[.]$ is the interval $0 \leq u \leq u_{\max}$, where discrete wavenumbers have non-negative integer values (Eq. (1)). After interpolation, one may want to restrict the support to

$$(0 \leq)u_L \leq u \leq u_H(\leq u_{\max}). \tag{10}$$

The variables and parameters on which $q[.]$ depends are arranged as a 12-tuple:

$$\vec{\psi} := \{r_B, r_E, \varphi_B, \varphi_E, \gamma, \zeta, \vartheta_L, \vartheta_H, p, d, u_L, u_H\}, \tag{11}$$

where $r_B, r_E, \varphi_B$ and $\varphi_E$ define the region of interest (Section 4.3), $\zeta$ is the average intensity threshold (Eq. (34)), $\gamma$ is the median filter mask size, in pixels, applied to $g[.]$ in the spatial domain (Section 4.3), $\vartheta_L$ and $\vartheta_H$ are the endvalues of $\Theta$ (Eq. (2)), $p$ is the model exponent (Eq. (3)), $d$ the polynomial degree; whereas, $u_L$ and $u_H$ have just been defined. As a consequence one states

**Definition** (*The $\sigma\eta$−derived feature vector*). Let

$$M := u_H - u_L + 1, \tag{12}$$

and

$$u_m := u_L + \frac{u_H - u_L}{M-1} m \quad \text{with } 0 \leq m \leq M-1, \tag{13}$$

then the *feature vector* derived from $g$ by means of $\sigma\eta$ is the $M$-dimensional vector $\vec{w}[g]$, the $m$th entry of which, $w_m[g]$, is the corresponding value of the polynomial:

$$w_m[g] := q[u_m; g, \vec{\psi}], \quad 0 \leq m \leq M-1. \tag{14}$$

The dependence of the feature vector on $\vec{\psi}$ is emphasized by writing $\vec{w}[g, \vec{\psi}]$. Each TAOS pattern is represented by such a $\vec{w}$. The entries of $\vec{\psi}$ range in a (usually convex) admissible polytope, $\Psi_{ad}$ i.e.,

$$\vec{\psi} \in \Psi_{ad}. \tag{15}$$

In fact, as it will be explained in Section 5.1, classification requires more parameters than those listed in Eq. (11), such as $N_T$ (Section 3.4.1, step 1), $K_T$ (Eq. (24)), $P$ (Eqs. (19) and (20)), to drive both the $\sigma\eta$ algorithm and training.

*Historical note*: The $\sigma\eta$ scheme was suggested by an experimental result [23] in Fourier optics and began being converted into an algorithm for image analysis only 10 years ago [24]. It has been the subject of investigation and expansion ever since.

### 3.4. Training, validation, and recognition

To begin, let $\vec{\psi}$ of Eq. (11) be fixed. The polytope $\Psi_{ad}$ (Eq. (15)) and the choice of $\vec{\psi}$ will be described later on. The main difference between standard multivariate classification and the algorithm described herewith is synthesized by the phrase "training for recognition." To attain the goal, two stages are needed: Stage one, {*training*, *validation*}, and stage two, {*tuning*, *recognition*}. Each stage contains one or more steps, which are described and justified below.

### 3.4.1. Stage one: principal components analysis at fixed $\vec{\psi}$

Classification applies to regions of interest (Section 4.3) selected from patterns. The $\sigma\eta$ algorithm operating on a region of interest returns the feature vectors $\{\vec{w}[g,\vec{\psi}]\}$. The linear classifier simultaneously performs two actions: it creates a vector space $\mathcal{Z}$, the coordinates of which depend on the $\{\vec{w}[g,\vec{\psi}]\}$, and assigns each pattern to a point $\vec{z}[g]\in\mathcal{Z}$.

Supervised training includes principal components analysis and class assignment at fixed $\vec{\psi}$ from step (1) to (7) below, and classification rating, based on a figure of merit.

(1) *Formation of the training set*: The training set $\mathfrak{T}$, where $\mathfrak{T}\subset\mathfrak{S}$, is formed by choosing $N_T$ patterns, which belong to at least two pre-defined and known classes.

(2) *Formation of the data matrix* **D**: The $\{\vec{w}[g,\vec{\psi}]\}$, arranged column-wise, form the data matrix **D** of size $M\times N_T$.

(3) *Formation of the correlation matrix* **C**: A correlation matrix can be defined *feature-wise* according to

$$\mathbf{C}^{(fw)}:=\mathbf{D}\cdot\mathbf{D}^{Trs}\in\mathcal{M}[M\times M].\tag{16}$$

As one can show, the rank $\rho_C:=rank[\mathbf{C}^{(fw)}]$ complies with

$$\rho_C\leq\min\{N_T,M\}.\tag{17}$$

(4) *Formation of a linear space*: The eigenvectors $\{\vec{\zeta}_p|1\leq p\leq\rho_C\}$ of $\mathbf{C}^{(fw)}$, known as the *principal components*, span a linear space, the dimension of which is at most $\rho_C$.

(5) *Selection of a P-dimensional subspace*, $\mathcal{Z}$: If $\lambda_1$ stands for the largest eigenvalue, then one chooses the first $P$ terms in the non-increasing sequence $\lambda_1,\ldots,\lambda_P$ and works in a subspace $\mathcal{Z}$, which is at most $P$-dimensional. By the way, the largest eigenvalues correspond to principal components that carry most of the sample variance. A typical choice is

$$2\leq P\leq\max\{2,\rho_C/2\},\tag{18}$$

which makes sense if the problem is non-degenerate.

Principal components analysis is due to K. PEARSON [25] and H. HOTELLING [26]. Some reference book [27] provides the justification of the method.

### 3.4.2. Stage one: class assignment at fixed $\vec{\psi}$

(6) *Transformation of training data at fixed $\vec{\psi}$*: In this section the dependence of all quantities on $\vec{\psi}$ will not show in the notations. By arranging the eigenvectors row-wise, one forms the matrix $\mathbf{L}\in\mathcal{M}[P\times M]$, which characterizes the linear classifier

$$\mathbf{L}:=\begin{bmatrix}\vec{\zeta}_1^{Trs}\\\vec{\zeta}_2^{Trs}\\\ldots\\\vec{\zeta}_P^{Trs}\end{bmatrix}.\tag{19}$$

Let **D** contain the training data. The counterpart of **D** in $\mathcal{Z}$ is a matrix given by

$$\mathbf{Z}:=\mathbf{L}\cdot\mathbf{D}.\tag{20}$$

The matrix **Z** yields the representation of $\mathfrak{T}$ in $\mathcal{Z}$ and allows its graphical display as well.

The classifier is said to be linear, solely because of Eq. (20). However, **Z** in turn, depends non-linearly on the data which form **D**.

(Aside) *Transformation of other data: validation and recognition*

A matrix $\mathbf{D}_{new}$ can be formed by data not used in the training stage, which, therefore, have played no role in the computation of **L**. Classification according to **L** is obtained by the linear transformation

$$\mathbf{Z}_{new}:=\mathbf{L}\cdot\mathbf{D}_{new},\tag{21}$$

which formalizes validation (stage one) and recognition (stage two below).

(7) *Class assignment at fixed $\vec{\psi}$*: Training is supervised, hence the class of belonging, $B[g]\in\mathfrak{B}$, of each pattern $g\in\mathfrak{S}$ is known. Let $\vec{z}_B$ be the empirical center of mass (center, for short) of class $B$, computed by the arithmetic average of all $z[g]$ which belong to $B$. The class empirical standard deviation is computed by the known definition. The class assignment rule implemented herewith is a deterministic one: assign $g$, represented by $\vec{z}[g]$ in $\mathcal{Z}$, to class $A[\vec{z}]$, the center of which is the

closest to $\vec{z}$ among all $\mathfrak{B}-$centers

$$A[\vec{z}]:=\underset{B\in\mathfrak{B}}{\mathrm{argmin}}|\vec{z}[g]-\vec{z}_B|. \tag{22}$$

The classifier, or confusion, matrix $\mathbf{M}_T[\vec{\psi}]$ is thus computed. Its row indices $B$ are classes of belonging, and its column indices $A$ are assigned classes. The entry $M_{BA}$ counts how many patterns of class $B$ have been assigned to class $A$.

### 3.4.3. Stage one: classification rating as a function of $\vec{\psi}$ on a fixed training set

Let the $\mathfrak{T}$ be fixed, but $\vec{\psi}$ be now allowed to vary in $\Psi_{\mathrm{ad}}$. After each sequence of steps (1) to (7), the following figure of merit $F[\mathbf{M}_T[\vec{\psi}]]$ is computed to rate the classification result:

$$F[\mathbf{M}[\vec{\psi}]]:=\left(\underline{\delta}:\mathbf{M}[\vec{\psi}]\right)/\sum_{AB}M_{AB}[\vec{\psi}], \quad \vec{\psi}\in\Psi_{\mathrm{ad}} \tag{23}$$

i.e., trace of $\mathbf{M}[\vec{\psi}]$ over sum of all entries of $\mathbf{M}[\vec{\psi}]$.

### 3.4.4. Stage one: training sets and design of the training experiments

From $\mathfrak{S}$, given card$[\mathfrak{B}_T]$ training classes, one forms $K_T$ training sets

$$\mathfrak{T}_k\subset\mathfrak{S}, \quad 1\le k\le K_T, \tag{24}$$

containing $N_T/\mathrm{card}[\mathfrak{B}_T]$ patterns per class. Every $\mathfrak{T}_k$ can differ from another training set by $\Delta N_T/\mathrm{card}[\mathfrak{B}_T]$ patterns per class, where

$$1\le\Delta N_T\le N_T. \tag{25}$$

Training based on $\mathfrak{T}_k$ is the $k$th *training experiment*. The sequence of such $K_T$ experiments is the *design-wide training*.

Let $\tau(>0)$ be a threshold value. If $F[\mathbf{M}_{Tk}[\vec{\psi}]]$ is the figure of merit of the $k$th training experiment, and $F[\mathbf{M}_{Vk}[\vec{\psi}]]$ is the figure of merit of the simultaneous validation experiment (Section 3.4.5 below), then the design-wide figure of merit $F_D$ is defined by

$$F_D[\vec{\psi}]:=\sum_{\kappa}F[\mathbf{M}_{T\kappa}[\vec{\psi}]], \quad \vec{\psi}\in\Psi_{\mathrm{ad}}. \tag{26}$$

For an index $\kappa$ to be included in the sum, *both* figures of merit shall comply with

$$F[\mathbf{M}_{T\kappa}[\vec{\psi}]]>\tau, \quad F[\mathbf{M}_{V\kappa}[\vec{\psi}]]>\tau, \quad 1\le\kappa\le K_T, \quad \vec{\psi}\in\Psi_{\mathrm{ad}}, \tag{27}$$

although the sum only contains training figures of merit. The values of $\tau$ are provided in Sections 5.1.1 and 5.2.1. The best, design-wide 12-tuple of Eq. (11), which needs not be unique, is defined by

$$\vec{\psi}^*:=\underset{\vec{\psi}}{\mathrm{argmax}}\in\Psi_{\mathrm{ad}} \ F_D[\vec{\psi}]. \tag{28}$$

### 3.4.5. Stage one: design evaluation

At the same time as the $\mathfrak{T}_k$, the validation- or $\mathfrak{B}_k-$sets are formed. $\mathfrak{B}_k$ shall contain patterns that have not been used for training i.e., $\forall k$, $\mathfrak{B}_k\cap\mathfrak{T}_k=\varnothing$. The corresponding $\{\vec{w}[g,\vec{\psi}]\}$ are not submitted to principal component analysis. Instead, they are projected onto $\mathcal{Z}_k$ according to Eq. (21). Each pattern of $\mathfrak{B}_k$ is assigned to a class according to Eq. (22). The corresponding matrix $\mathbf{M}_{Vk}[\vec{\psi}]$ is formed. Validation is rated by the counterpart of Eq. (23). If Ineq. (27) is not met, then the $k$th experiment is rejected. A design is evaluated by means of the $F_D[\vec{\psi}^*]$ of Eq. (26). The $\vec{\psi}^*$ of Eq. (28) is passed on to stage two, together with the list of $\tau-$selected training sets $\{\mathfrak{T}_\kappa, 1\le\kappa\le K_T\}$.

### 3.4.6. Stage two: tuning

The output from stage one includes the $\vec{\psi}^*$ of Eq. (28). Given the same training sets $\mathfrak{T}_\kappa$, $1\le\kappa\le K_T$ (Eqs. (24) and Ineq. (27)), tuning consists of re-computing the matrix $\mathbf{L}_\kappa$ of Eq. (20), at $\vec{\psi}=\vec{\psi}^*$ only. The subspaces $\mathcal{Z}_\kappa$ are thus recreated.

### 3.4.7. Stage two: recognition

The patterns to be recognized undergo spectrum enhancement controlled by $\vec{\psi}^*$ and their feature vectors are projected onto $\mathcal{Z}_\kappa$ as in the validation stage (Section 3.4.5).

Since experiment-wise class assignment of the unknown pattern may vary, then results have to be assembled design-wide.

### 3.4.8. Stage two: information fusion

Each experiment $\kappa$ assigns a pattern to be recognized, $g_R$, represented by $\vec{z}_\kappa[g_R]$, to one of the training classes. By means of two fusion rules, not described herewith in detail, one, and only one design-wide class $A_D\in\mathfrak{B}$ is assigned to each $g_R$.

### 3.4.9. Stage two: conclusion

The main result of stage two is the design-wide assignment matrix, $\mathbf{M}_{DA}$. The latter is formed by scanning the list of recognized patterns, record by record, and comparing the class to which a pattern belongs to the class to which it has been assigned according to Section 3.4.8.

$$
\mathbf{M}_{DA} := 
\begin{array}{|c||c|c|c|}
\hline
\text{design} - \text{wide assigned class} \rightarrow & & & \\
- - - - - - - - - - - - & A_{1D} & A_{2D} & \ldots \\
\text{class of belonging} & & & \\
\text{or experiment data set} \downarrow & & & \\
\hline\hline
B_1 & \ldots & \ldots & \ldots \\
\hline
B_2 & \ldots & \ldots & \ldots \\
\hline
\ldots & \ldots & \ldots & \ldots \\
\hline
\end{array}
\tag{29}
$$

Of independent interest is the design-wide statistics of class assignments i.e., the scoring of recognition patterns which are always assigned to the same class throughout the design. Details are not provided herewith.

## 4. TAOS pattern classification

### 4.1. Obstacle inversion vs. statistical classification

Strictly speaking, the interpretation of a *TAOS* pattern means the determination of size, shape and complex refractive index of the particle that gave rise to the pattern. This is an inverse problem of electromagnetics.

The scatterer is a material particle contained in a bounded domain $D{\subset}\mathbb{R}^3$ and may be described either by a potential $q[.]$ or by a complex refractive-index function $n[.]$. In principle one could try to reconstruct $q[.]$ or $n[.]$ from knowledge of the incident wave and the scattering pattern. However, there are discrepancies between the TAOS experimental conditions and the ideal conditions required by scattering theory for unique reconstruction:

(a) real-world scatterers have arbitrary shapes and may have a complicated internal structure, which translates into a position-dependent refractive index, $n[\vec{x}\,]$, $\vec{x} \in D$;
(b) in TAOS experiments the particle is illuminated by only one plane wave (experiments with two incident waves have been reported, though [28]);
(c) even if the scatterer has a symmetry axis, orientation of the latter with respect to the laboratory frame when scattering occurs is generally unknown (exceptions are asbestos fibers [29] and distorted droplets [3,4], because their major axis can be oriented along with the airstream);
(d) in spite of the advancements in the optics, data are always collected within a limited aperture $\tilde{s}^2$, not on the whole surface of the unit sphere, $S^2$, i.e., do not cover the full range of scattering angles;
(e) TAOS patterns consist of intensity values without information about either phase or polarization.

In view of all these deviations from the ideal context, and in spite of remarkable progress occurring in the field of shape reconstruction over the last 15 years [30,31], the inversion of TAOS data remains an open and difficult problem.

On the brighter side, TAOS patterns can be obtained in large quantities. The visual inspection of patterns from known materials suggests significant morphological differences. Even in the absence of perceivable symmetry, coarser and finer features are immediately seen. In the image-processing parlance, these features are respectively called "structure" and "texture" and give rise to the OSHER–RUDIN model of an image [32].

As a consequence there are sufficient arguments to leave aside, for the moment at least, the solution of an inverse obstacle problem and search for an automated procedure, based on machine learning, that analyzes the TAOS patterns quantitatively on a morphological basis and applies linear classification.

Historically, one of the first attempts at morphological analysis [2] did rely on scoring intensity peaks and valleys in the pattern: the score was related to the scatterer size, and TAOS patterns were grouped by morphology. No machine learning scheme was however implemented. The early implementations of a classifier, with separate training and recognition stages, were described years ago [12,13] and shown to work well, although only applied to a relatively small number (tens) of patterns. To the best of the authors' knowledge, no other trainable classifier applied to TAOS patterns has been described before.

The machine learning approach marked a change in the strategy aimed at analyzing TAOS patterns, which can be formalized by

*Ansatz* 1: particles of a given material form classes.

It follows:

**Corollary.** *The corresponding TAOS patterns form classes as well.*

Practically this translates into three rules.

(1) *Instead of relating a TAOS pattern to the original particle, assign the pattern to a class.*
(2) *If the class B to which the particle belongs is known to a supervisor, then the learning machine shall return A=B.*
(3) *When dealing with large numbers of patterns, the success condition A=B shall be met in most cases.*

These are specifications for a statistical pattern classifier, to be designed according to Section 3. The TAOS pattern inversion problem has been replaced by a pattern classification problem, to be solved by machine learning: a transition has occurred from inverse scattering to statistical learning.

What one shall expect from automated classification depends on the available data sets, on training design, on the ability of the feature extraction module and of the related control parameters to capture the distinctive pattern morphology, and, eventually, on the linear classification module, which shall first learn from features, then recognize patterns.

### 4.2. Data set and training sets

Lack of knowledge about the probability density $f_G[g]$ of the generator, and about the conditional probability density $f_S[B|g]$ of the supervisor, mentioned in Section 3.1 require, in addition to Ansatz 1 of Section 4.1, assumptions about the raw pattern sets.

*Ansatz* 2: all available *Bq* patterns do come from *B. subtilis* spores.
*Ansatz* 3: neither the outdoor sampling sets (*K0*, …, *K5*), nor the *sq* set, contain any pattern from *B. subtilis* spores.
*Ansatz* 4: training-wise, the same amount of information is carried by

- any *Bq* pattern, regardless of orientation of the spore in the laboratory frame (Fig. 1b),
- any *Dq* or *Pq* pattern, as if the respective microsphere had exactly the preset radius (Table 1),
- all training patterns, regardless of class of belonging, provided their signal strength exceeds a given threshold (Ineq. (34)).

Supervised training shall be based on these assumptions, none of which, in turn, can either be disproved or challenged. The data set is presented by Table 1. Material preparation and pattern collection were described by Aptowicz et al. [14]. In particular,

- *Gq* patterns were recorded under the same experimental conditions as a TAOS pattern, but without any particle being hit by the illuminating laser; the photo–electron counts in each pixel of the ICCD in this data set are mainly produced from background scattering and stray light in the system, as well as the thermal noise and read-out detector noise;
- *Bq* patterns come from single *B. subtilis* spores [9];
- *Dq* patterns come from dioctyl-phthalate droplets of 2.8 μm diameter;
- *Pq* and *Qq* come from two separate experiments with 1.034 μm polystyrene latex spheres;
- *K* sets come from outdoor sampling carried out in Maryland on 2004 October 6;
- *sq* patterns were obtained by scraping Diesel engine soot from a tailpipe, resuspending the material and feeding it into the *Royco* aerosol generator.

In view of Ansätze 2 and 4, patterns {Bq, *Dq*, *Pq* and Qq} are regarded as *homogeneous* by origin and fit for machine learning. Instead, physics suggests to regard *K* patterns as *heterogeneous* by origin; therefore Ansatz 4 does not apply and the inclusion of *K* into a training set would not be advisable. The attitude towards *sq* is less clear-cut. For this reason, two Designs have been worked on (Sections 4.3 and 5), one where *sq* is submitted to stage two only (Design 06, Section 5.1), another where *sq* has been included into the training set as well (Design 07, Section 5.2).

Ansatz 4 practically amounts to ignoring the particle size distribution, which can, in principle, be available from experiments. For example, polystyrene latex sphere aggregates and *B. subtilis* spore aggregates were characterized statistically by Holler et al. [12]. The aerosol generators currently in use [14] also control the particle size distribution of the given material.

### 4.3. Classification designs and preprocessing stages

The aim is to discriminate airborne biological threat particles, represented by pathogenic bacterial spores (*Bq*), from all other patterns. Those from natural background, which are by definition heterogeneous and are represented herewith by the *K* set, are not further analyzed herewith for possible morphological classification or source apportionment. Discrimination of *Bq* from *sq* is required too, because diesel soot aggregates are weak scatterers: some of their TAOS patterns may be feature-poor and as such more easily mistaken for patterns of other classes, *Bq* in particular.

Pattern preprocessing consists of selecting a region of interest, applying a median filter, accepting the region by signal strength, transforming coordinates and normalizing the gray scale.

*Regions of interest*: four regions of interest are defined i.e., ring sectors with the same center as the pattern, $\mathbf{O} \equiv \{511, 511\}$ pixels, and the following endpoints:

$$\tilde{S}^2_{NE} = \{r_B \leq r \leq r_E, 0.392 \leq \varphi \leq 1.342\}, \tag{30}$$

$$\tilde{S}^2_{NW} = \{r_B \leq r \leq r_E, 1.800 \leq \varphi \leq 2.750\}, \tag{31}$$

$$\tilde{S}^2_{SW} = \{r_B \leq r \leq r_E, -2.750 \leq \varphi \leq -1.800\}, \tag{32}$$

$$\tilde{S}^2_{SE} = \{r_B \leq r \leq r_E, -1.342 \leq \varphi \leq -0.392\}, \tag{33}$$

where the radius values are in pixels and the $\varphi$ values in radians. Two *R*-intervals will be chosen, $155 \leq r \leq 470$ pixel and $200 \leq r \leq 495$ pixel, as specified in the designs below (Table 3 and Section 5.2.1). Ring sectors are processed independently by the feature extraction algorithm of Section 3.3 and by the linear classifier.

*Filtering and thresholding*: noise is reduced by applying to each sector a median filter, the mask of which is a square of $\gamma^2$ pixels. Then the sector-averaged intensity $\langle I \rangle_{\text{sector}}$ is compared to a threshold, $\zeta$. If

$$\langle I \rangle_{\text{sector}} > \zeta, \tag{34}$$

the sector is accepted for classification.

*Coordinate transformation, gray scale normalization*: polar coordinates are transformed to Cartesian, in order to obtain a tile according to Section 3.3

$$\tilde{S}^2 \rightarrow \Omega$$
$$\{r, \varphi\} \mapsto \{x_1, x_2\}. \tag{35}$$

Accordingly, the median-filtered, ring sector-supported intensity is bi-linearly interpolated and becomes the function $g[x_1, x_2]$. The gray scale of $g[.]$ is finally normalized.

Two training design groups are accounted for herewith, Design 06 and Design 07, which in turn contain different designs made of experiments.

| Design 06 : |
| --- |
| reject faint patterns,  train by $\{Bq, Dq, Pq\}$,  recognize $\{Bq, K5, sq\}$. |

As will be explained in Section 5.1, *Bq* patterns enter both stage one and stage two.

| Design 07 : |
| --- |
| reject faint patterns,  train by $\{Bq, Dq, Pq, sq\}$,  recognize $\{Bq, K0, K1, K5, sq\}$. |

The new addition here are training sets which include DIESEL soot patterns (the *sq* set). In both Designs, the number of patterns admitted to stage two is arbitrary and can be large ($\approx$ thousands).

## 5. Results

### 5.1. Design 06

| Design 06 : |
| --- |
| reject faint patterns,  train by $\{Bq, Dq, Pq\}$,  recognize $\{Bq, K5, sq\}$ |

#### 5.1.1. Training

The four ring sectors of Eqs. (30)–(33), as explained in Section 4.3, are treated separately. The threshold $\zeta$ is determined by processing the *Gq* patterns (Table 1). After a few exploratory attempts, the value $\zeta = 470$ a.u. has been assigned, which has been kept fixed hereinafter. The training sets are formed by the scheme of Table 2, which pertains to $N_T = 60$, $\Delta N_T = 15$ (hence $\Delta N_T / \text{card}[\mathfrak{B}_T] = 5$ patterns per class) and $K_T = 7$. For different $N_T$ and $K_T$ the scheme is similar.

In Table 2 bullets between square brackets represent sequences of contiguous patterns. E.g., in training set $k = 1$, patterns $\mathcal{N}.6$ to $\mathcal{N}.25$ are selected. However, pattern numbers in the actual training lists differ from those in the raw input list, because some patterns do not pass the test of Ineq. (34) and are dropped. As can be deduced from the construction of Table 2, since $N_T = 60$ and $\text{card}[\mathfrak{B}] = 3$, then, in stage one, 50 *Bq* patterns are used at least once, 10 patterns ($\mathcal{N}.6$ to $\mathcal{N}.10$ and $\mathcal{N}.51$ to $\mathcal{N}.55$) are used exactly once, 10 more patterns ($\mathcal{N}.11$ to $\mathcal{N}.15$ and $\mathcal{N}.46$ to $\mathcal{N}.50$) are used exactly twice, and 20 patterns ($\mathcal{N}.21$ to $\mathcal{N}.40$) are used 4 times. The same rule applies to *Dq* patterns. Instead, all the 39 *Pq* patterns are used more than once in stage one. All *Bq*s are submitted for recognition to stage two, from which *Dq*s and *Pq*s are excluded. Stage two instead analyzes *K5*s and *sq*s.

Other than depending on $K_T$ (kept fixed herewith) and $\{\mathfrak{T}_k | 1 \leq k \leq K_T\}$, this design includes four values of $\gamma$ : 3, 5, 7, and 9 [pixels], radial intervals $R := [r_B, r_E]$ with two different endpoints, and a few different integration arcs ($\equiv \Theta$) in the reciprocal domain.

**Table 2**
Design 06: formation of the training sets when $N_T=60$. $\mathcal{N}$: pattern number; $k$: training set number.

| $\mathcal{N}.\rightarrow$ | 6 | . | 11 | . | 16 | . | 21 | . | 25 | 26 | . | 30 | 31 | . | 35 | 36 | . | 40 | 41 | . | 45 | . | 50 | . | 55 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| $k\downarrow$ | | | | | | | | | | | | | | | | | | | | | | | | | |
| 1 | [• | • | • | • | • | • | • | • | •] | | | | | | | | | | | | | | | | |
| 2 | | | [• | • | • | • | • | • | • | • | • | •] | | | | | | | | | | | | | |
| 3 | | | | | [• | • | • | • | • | • | • | • | • | • | •] | | | | | | | | | | |
| 4 | | | | | | | [• | • | • | • | • | • | • | • | • | • | • | •] | | | | | | | |
| 5 | | | | | | | | | [• | • | • | • | • | • | • | • | • | • | • | • | •] | | | | |
| 6 | | | | | | | | | | | | [• | • | • | • | • | • | • | • | • | • | • | •] | | |
| 7 | | | | | | | | | | | | | | | | [• | • | • | • | • | • | • | • | • | •] |

**Table 3**
Design 06: $\gamma = 7$.

| $N_T$ | $R\rightarrow$ $\zeta\rightarrow$ $\Theta\downarrow$ | [155,470] 470 | [200,495] 470 |
|---|---|---|---|
| 21 | | M7H16, M7H56, M7H96, M7HD6 | |
| 30 | [80, 100] | OH718, OH758, OH798, OH7D8 | OH721, OH761, OH7A1, OH7E1 |
| 60 | | M7H15, M7H55, M7H95, M7HD5 | MzH08, MzH48, MzH88, MzHC8 |
| 30 | [75, 105] | ON719, ON759, ON799, ON7D9 | ON720, ON760, ON7A0, ON7E0 |
| 60 | | | MzN09, MzN49, MzN89, MzNC9 |
| 30 | [70, 110] | Ok718, Ok758, Ok798, Ok7D8 | |
| 60 | | | MzI11, MzI51, MzI91, MzID1 |
| 30 | [65, 115] | | |
| 60 | | | MzO14, MzO54, MzO94, MzOD4 |
| 30 | [60, 120] | Ok717, Ok757, Ok797, Ok7D7 | |
| 60 | | | |
| Design 06: $\gamma = 9$ | | | |
| 21 | | MzH16, MzH56, HzH96, MzHD6 | |
| 30 | [80, 100] | OH919, OH959, OH999, OH9D9 | OH921, OH961, OH9A1, OH9E1 |
| 60 | | MzH15, MzH55, MzH95, MzHD5 | MzH11, MzH51, MzH91, MzHD1 |
| 30 | [75, 105] | | |
| 60 | | | MzN14, MzN54, MzN94, MzND4 |
| 30 | [70, 110] | Ok918, Ok958, Ok998, Ok9D8 | |
| 60 | | | MzI13, MzI53, MzI93, MzID3 |

Typical design matrices for $\gamma = 7$ and 9 are displayed by Table 3. Design names, such as $\{OH718, OH758, OH798, OH7D8\}$, reflect the sector sequence, $\{NE, NW, SW, SE\}$. In summary, the "parameters" (numbers and sets) which control training form the $n$-tuple

$$\wp = \{K_T, \{\mathfrak{T}_k\}, r_B, r_E, \theta_B, \theta_E, \gamma, \zeta, \vartheta_L, \vartheta_H, p, d, u_L, u_H, \tau\}. \tag{36}$$

Throughout this Design, Ineq. (27) has been applied with $\tau = 0.4$. However, only the entries of the 7-tuple, still denoted by $\overrightarrow{\psi}$ as in Eq. (11)

$$\overrightarrow{\psi} = \{\gamma, \vartheta_L, \vartheta_H, p, d, u_L, u_H\}, \tag{37}$$

are iteratively scanned to maximize the design-wide figure of merit (Eq. (26)). The other entries of $\wp$ in Eq. (36) give rise to different classification designs, which have to be programmed and launched separately.

### 5.1.2. Discrimination results

Results from a classifier trained by different sets can be of two types: design-wide classification matrices and experiment-wise graphical displays. The former are yielded by fusion rules, such as those of Section 3.4.8. The latter pertain to the $\mathcal{Z}_\kappa$ space created from each $\mathfrak{T}_\kappa$ as explained in Sections 3.4.6 and 3.4.8. Usually, one value of $\kappa$ is selected in representation of the whole design. According to Section 3.4.1, Eq. (18) in particular, all $\mathcal{Z}_\kappa$ have the same dimension, $P = 10$ here. Usually the first two coordinates (principal components) of $\mathcal{Z}_\kappa$ carry the largest amount of information in the sample. Display in the $\{z_1, z_2\}_\kappa$ plane is therefore sufficient to visualize where the original TAOS patterns have been positioned by the classifier.

Since principal components analysis does not tell anything about the possible continuous dependence of coordinates in $\mathcal{Z}$ on the training set, then there is no general rule to determine an "overall" $\mathcal{Z}$ space by "assembling" the different $\mathcal{Z}_\kappa$. Moreover, the class assigned to a given pattern may change along with $\kappa$, as pointed out by the captions of Figs. 4 and 6.
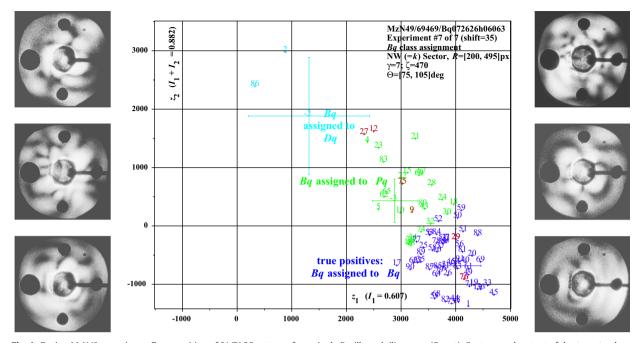
**Fig. 4.** Design *MzN49*, experiment 7: recognition of 91 TAOS patterns from single *Bacillus subtilis* spores (*Bq* set). Centre panel: cutout of the $\{z_1, z_2\}_7$ plane. $I_1$: sample variance carried by the first principal component; $I_2$: sample variance carried by the first and second together. Blue labels (lower right corner): *Bq* patterns recognized as such (true positives, $\simeq 80\%$); green labels (centre): *Bq*'s mistaken for *Pq*; cyan labels (top left): *Bq*'s mistaken for *Dq*. The color-coded centers of mass (step 7 of Section 3.4.2) of each training class are shown, together with their error bars. This training experiment discriminates all three classes *Bq*, *Dq* and *Pq* from one another, because error bars do not overlap with any other along both directions. Patterns shown in the side stacks are identified by red labels in the centre panel. Counterclockwise from top left: *Bq*130 (label=27) classified as *Dq*, *Bq*181 (75) classified as *Pq*, *Bq*133 (29) classified as *Bq*; similar rules for patterns *Bq*182 (76), *Bq*112 (9), and *Bq*115 (12). Class assignment is based on the *NW* pattern sector: a bowtie (as in Fig. 1, leftmost panel) is detected in *Bq*133 (29) and *Bq*133 (76), in spite of asymmetry, whereas it is missed in *Bq*112 (9), which exhibits it clearly. The almost-random patchwork of *Bq*130 (27) and *Bq*115 (12) assigns them to *Dq*, although near the divide from the *Pq* class. Patterns *Bq*130, *Bq*181, *Bq*133, *Bq*115 and *Bq*182 are stably assigned to the same class by the 7 experiments in the design. Pattern *Bq*112 is assigned to *Pq* by 6 experiments and to *Bq* by 1.

In general terms, a result is judged by the misclassification errors (a.k.a., loss) it brings about. Ideally, design-wide false negatives (misrecognized *Bq* patterns) and false positives (patterns from other materials mistaken for *Bq*) should be the fewest at the same time. Herewith, the best result and some other (the "next best") are summarized for comparison.

The format chosen to represent design-wide class assignment from stage two has been simplified with respect to the actual output from the computer program. Namely, it is limited to the classifier matrix augmented by two columns. The latter carry the counts of actually analyzed patterns and the input totals.

Design *MzN49* has yielded the best performance within the explored set of control parameters. The relevant ones are: $N_T = 60$ i.e., 20 patterns per class, $\tilde{S}^2 = NW$ ring sector with $R = [200, 495]$ pixels, $\gamma = 7$ pixels, $\zeta = 470$ a.u., and $\Theta = [75, 105]$ degrees. Classification, according to the *Bq* column (assigned class) of Table 4 features $18/91 \simeq 20\%$ false *Bq* negatives (*Bq* patterns not recognized as such), $98/957 \simeq 10\%$ false *K*5 positives (*K*5 patterns mistaken for *Bq*'s) and $12/114 \simeq 11\%$ false *sq* positives.

Experiment-wise, the displayed class assignment pertains to $\kappa = 7$ (Experiment # in the figure heading) and is split into three figures, Figs. 4–6, as many as the classes of belonging, *Bq*, *K*5 and *sq*, submitted to stage two. Displays are zoomed-in projections of $\mathcal{Z}_7$ onto a rectangle in the plane $\{z_1, z_2\}_7$. The rectangle contains the majority of patterns, but not all. Details and comments are in the respective captions.

The sensitivity of classification with respect to some control parameters can be inferred from the following examples.

Dependence on ring sector at fixed $N_T = 60$, $\gamma = 7$, $\zeta = 470$, $\Theta = [75, 105]$: the other three sectors have more false negatives, up to 57%, and false positives, up to 23% from *sq*.

Dependence on $\Theta$ at fixed $N_T = 60$, $\gamma = 7$, $\zeta = 470$: regardless of ring sector, wider arcs e.g., [70, 100], [65, 115] degree, deteriorate classification remarkably. The narrower arc, $\Theta = [80, 100]$ yields, in one instance (*SW* ring sector) 18% false *Bq* negatives and 23% false *K*5 positives.

Dependence on $N_T$ at fixed $\gamma = 7$ and $\zeta = 470$: the best result from a smaller training set, $N_T = 30$, over all four ring sectors, where $\Theta$ has been assigned the values [80,100] and [75,105], yields in the case of design *OH7A1* (*SW* sector, $\Theta = [80, 100]$), fewer false *Bq* negatives, $15/90 \simeq 17\%$ but $213/969 \simeq 22\%$ false *K*5 positives. The classifier matrix of *OH7A1* is shown in Table 5.

Results at $\gamma = 9$, $\zeta = 470$ and $N_T = 60$, with variable sector and $\Theta$: too narrow ($\Theta = [80, 100]$) or too wide ($\Theta = [70, 100]$) an arc yield at least $24/92 \simeq 26\%$ *Bq* false negatives, although, for a different sector and arc, the *K*5 false positives can be as low as
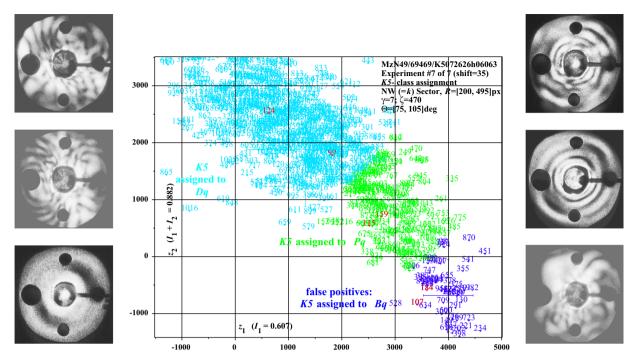
**Fig. 5.** Design *MzN*49, experiment 7: recognition of 957 TAOS patterns from outdoor sampling (*K*5 set). Centre panel: cutout of the $\{z_1, z_2\}_7$ plane. Blue labels (lower right corner): *K*5 patterns mistaken for *Bq*'s (false positives, ≃10%); green labels (centre): *K*5's assigned to *Pq*; cyan labels (top left): *K*5's assigned to *Dq*. Ccw from top left: *K*5034 (label=124) classified as *Dq*, *K*5024 (115) classified as *Pq*, *K*5016 (107) mistaken for *Bq*, *K*5094 (184, mistaken for *Bq*), *K*5069 (159), and *K*5008 (99). Coarse features in the *NW* sectors of patterns *K*5016 (107) and *K*5094 (184) may have caused the false alarm. The other four patterns have almost surely been assigned according to their ring and fringe structures. All patterns are stably assigned to the same class by the 7 experiments in the design.

74/957≃8%. The best result, within the explored parameter subset, is design *MzN*54, sector *NW* and $\Theta = [75, 105]$, the classifier matrix of which is shown in Table 6.

### 5.2. Design 07

| Design 07 : |
| --- |
| reject faint patterns,  train by $\{Bq, Dq, Pq, sq\}$,  recognize $\{Bq, K0, K1, K5, sq\}$ |

#### 5.2.1. Training

The input set is $\mathfrak{S} = \{Bq, Dq, Pq, sq\}$, which contains 98 *Bq*, 200 *Dq*, 39 *Pq* and 288 *sq* patterns. Training sets have been formed, with due changes, by the same scheme as in Table 2, by letting $N_T = 28$ (7 patterns per class) or $N_T = 60$ (15 patterns per class). In both cases $\Delta N_T / \text{card}[\mathfrak{B}_T] = 5$ patterns per class and $K_T = 7$ experiments have been kept fixed. The following control parameter values have been tested: $\tau = 0.5$ (Ineq. (27)), fixed, $[r_b, r_E] = [155, 470]$ and $[200, 495]$, $\gamma = 3$, 5 and 7, $\zeta = 470$, fixed, $\Theta = [50, 130]$, $[60, 120]$ and $[70, 110]$. The overview of Design 07 would be provided by the counterpart of Table 3, which is not displayed for the sake of conciseness. A total of 11 designs of 4 sectors each have been run. Training by four classes, one of which, *sq*, known to be morphologically inhomogeneous, may cause poor performance.

#### 5.2.2. Discrimination results

The design which returns the fewest (19/93≃20%) false *Bq* negatives is *MzI*71, which is now described. The control parameters are $N_T = 60$ i.e., 15 patterns per class, $\tilde{S}^2 = NW$ ring sector with $R = [200, 495]$, $\gamma = 5$, $\zeta = 470$ a.u., and $\Theta = [70, 110]$. Design wide, 6 out of the $K_T = 7$ experiments pass the Ineq. (27) test with $\tau = 0.5$. As it can be deduced from Table 7, the false positives from *K*5, 83/957≃9%, and *sq*, 12/114≃11%, stay low. Discrimination of *K*0 from *Bq*, which returns 139/1960≃9% is stable, in spite of the larger number of processed patterns, as compared to Design 06.

Experiment-wise, the expected poor performance translates in the confusion of classes *Dq*, *Pq*, *sq*. Fig. 7 is limited to the classification of *Bq* patterns. Further details are provided in the caption. As any other, this experiment implements all provisions stated as ANSÄTZE 2–4 in Section 4.2.

In analogy with the previous Design, sensitivity of discrimination with respect to a few control parameters has been carried out within Design 07. The corresponding classification matrices are not shown. Only the highlights are provided. If training occurs with smaller sets ($N_T = 28$), when the *SW* sector is chosen, $\gamma = 3$ and $\Theta = [50, 130]$, then, at best, 25/93≃27% false *Bq* negatives are found, whereas the most false positives (15/117≃13% from *sq*) stay under control. This may be the result of undertraining. Training with $N_T = 60$, based on the *NW* sector, with $\gamma = 7$ and $\Theta = [70, 110]$, does not seem to bring
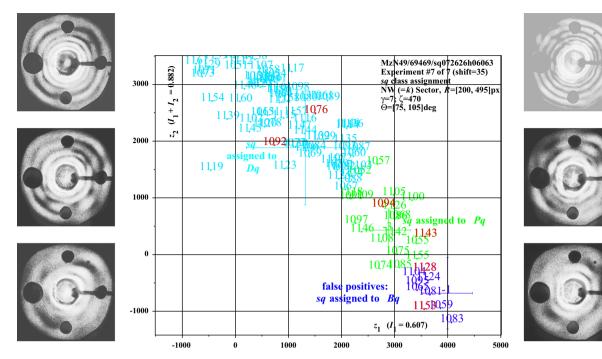
**Fig. 6.** Design *MzN*49, experiment 7: recognition of 114 TAOS patterns from resuspended diesel soot particles (*sq* set). Centre panel: cutout of the $\{z_1, z_2\}_7$ plane. Blue labels (lower right corner): *sq* patterns mistaken for *Bq*'s (false positives, $\simeq 11\%$); green labels (centre): *sq*'s assigned to *Pq*; cyan labels (top left): *sq*'s assigned to *Dq*. Ccw from top left: *sq*109 (label=1092) classified as *Dq*, *sq*116 (1094) classified as *Pq*, *sq*204 (1128) mistaken for *Bq*, *sq*261 (1153, mistaken for *Bq*), *sq*234 (1143), and *sq*068 (1076). Coarse features in the *NW* sectors of patterns *sq*204 (1128) and *sq*261 (1153) may have been mistaken for bowties. The other four patterns have been assigned either to *Pq* or to *Dq* according to the spatial frequency of their concentric rings, regardless of signal strength. Patterns *sq*109, *sq*116, *sq*204, *sq*068 and *sq*261 are stably assigned to the same class by the 7 experiments in the design. Pattern *sq*234 is assigned to *Bq* by 4 experiments and to *Pq* by 3, including the displayed one.

**Table 4**
Design *MzN*49 (*NW* sector, $N_T = 60, R = [200, 495], \gamma = 7, \zeta = 470, \Theta = [75, 105]$): augmented classifier matrix.

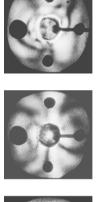| Assigned class→ | Bq | Dq | Pq | $\{\gamma, \zeta\}$−selected total | out of |
|---|---|---|---|---|---|
| Class of belonging↓ | | | | | |
| Bq | 73 | 5 | 13 | 91 | 98 |
| K5 | 98 | 747 | 112 | 957 | 993 |
| sq | 12 | 87 | 15 | 114 | 288 |

**Table 5**
Design *OH7A1* (*SW* sector, $N_T = 30, R = [200, 495], \gamma = 7, \zeta = 470, \Theta = [80, 100]$): augmented classifier matrix.

| Assigned class→ | Bq | Fq | Pq | $\{\gamma, \zeta\}$−selected total | out of |
|---|---|---|---|---|---|
| Class of belonging↓ | | | | | |
| Bq | 75 | 8 | 7 | 90 | 98 |
| K5 | 213 | 537 | 219 | 969 | 993 |
| sq | 16 | 95 | 5 | 116 | 288 |

**Table 6**
Design *MzN*54 (*NW* sector, $N_T = 60, R = [200, 495], \gamma = 9, \zeta = 470, \Theta = [75, 105]$): augmented classifier matrix.

| Assigned class→ | Bq | Fq | Pq | $\{\gamma, \zeta\}$−selected total | Out of |
|---|---|---|---|---|---|
| Class of belonging↓ | | | | | |
| Bq | 70 | 7 | 14 | 91 | 98 |
| K5 | 84 | 768 | 105 | 957 | 993 |
| sq | 13 | 89 | 12 | 114 | 288 |

**Table 7**

Design $MzI71$ $N_T=60$, pattern sector $\tilde{s}^2=NW(\equiv k)$, $\gamma = 5$ pixels, $\Theta=[70,110]$ degrees.

| Assigned class→<br>Class of belonging↓ | $Bq$ | $Fq$ | $Pq$ | $sq$ | $\{\gamma,\zeta\}$−selected total | out of |
|---|---|---|---|---|---|---|
| $Bq$ | 74 | 0 | 10 | 9 | 93 | 98 |
| $K5$ | 83 | 240 | 54 | 580 | 957 | 993 |
| $K0$ | 139 | 483 | 107 | 1231 | 1960 | 1999 |
| $sq$ | 12 | 16 | 13 | 73 | 114 | 288 |



**Fig. 7.** Design $MzI71$, experiment 7: recognition of 93 TAOS patterns from single *Bacillus subtilis* spores ($Bq$ set). Centre panel: cutout of the $\{z_1,z_2\}_7$ plane. $I_1$: sample variance carried by the first principal component; $I_2$: sample variance carried by the first and second together. Blue labels (lower right corner): $Bq$ patterns recognized as such (true positives, $\simeq 80\%$); cyan green labels (centre): $Bq$'s mistaken for $Pq$; green labels (middle left): $Bq$'s mistaken for $sq$; blue labels (top left): $Bq$'s mistaken for $Dq$. The color-coded centers of mass (step 7 of Section 3.4.2) of each training class are shown, together with their error bars. This training experiment only discriminates $Bq$ from the remaining ones. Namely, $Dq$, $sq$ and $Pq$ are confused, because error bars overlap. Patterns shown in the side stacks are identified by red labels in the centre panel. Ccw from top left: $Bq105$ (label=2) classified as $Dq$, $Bq107$ (4) classified as $sq$, $Bq113$ (10) classified as $Pq$, $Bq168$ (63) recognized as $Bq$; similar rules for patterns $Bq122$ (19), $Bq167$ (62), $Bq126$ (23) and $Bq192$ (86). Class assignment is based on the $NW$ pattern sector: the only clearly visible bowtie (as in Fig. 1, leftmost panel) in $Bq113$ (10) is missed. Patterns $Bq105$ (2) and $Bq192$ (86) may result from laser misfiring, not from actual spores. However, they pass Ineq. (34): their removal would bring prejudice into the learning algorithm, in contradiction to ANSÄTZE 2 and 4.

about improvement with respect to the above described $MzI71$ design: namely, the false $Bq$ negatives increase to $20/93\simeq 21\%$, false positives from $K5$ and $\{K0,K1\}$ decrease to $73/957\leq 8\%$ and to $118/1959\simeq 6\%$ respectively; however, false positives from $sq$ increase to $16/114\leq 14\%$. Finally, when $\gamma$ is increased to 9 then, the $NW$ sector and $\Theta=[60,120]$ at best yield the same false $Bq$ negatives, whereas all false positives drop by a small fraction: $68/957\simeq 7\%$ from $K5$, $109/959\leq 6\%$ from $\{K0,K1\}$ and $15/114\simeq 13\%$ from $sq$.

## 6. Discussion

As one may have deduced from Section 5, the design of classification experiments is affected by decisions at various levels and depends on many numerical parameters. The top-level decision has affected the learning machine architecture: decomposition into two interacting modules, feature extraction and classification, has been chosen. Feature extraction has

relied on the spectrum enhancement algorithm of Section 3.3, which is a general-purpose algorithm for the separation of image structure from texture; as such, it is not specific to the physical process from which images arrive. Namely, a process-specific algorithm for TAOS patterns would be an inverse problem solver based on scattered intensity data, without knowledge of either phase or polarization. To the best of these authors' knowledge said algorithms are not available.

About machine learning, a wide variety of methods is available to choose from. Linear classification by standard multivariate statistical analysis (Section 3.4.1) and a supervised training scheme (Sections 3.4.2 and 3.4.3) have been chosen for ease of integration with the feature extraction module.

The three most relevant task specific functions added to the linear classifier are: the automated formation of training sets (Section 3.4.4 and Table 2), the computation of design-wide figures of merit (Eq. (26)), and information fusion rules (Section 3.4.8). The latter assemble the results of all experiments in a design and form the corresponding classification matrix (Eq. (29)) of which a few instances have been provided (Tables 4–7).

At the implementation level, many parameters have been chosen to control the algorithm. Namely, the designs of Sections 4 and 5 depend on the $n$-tuple defined by Eq. (36). As pointed out in Section 5.1.1, and as can be deduced from inspecting the results of Section 5, even if the entries of $\wp$ are assigned a few different values, one obtains a very large set. The latter has been explored in part. Physical intuition and the results obtained step by step have played a role in selecting $n$-tuples of values. A key role in reducing false $sq$ positives has been played by intensity thresholding ($\zeta$, Eq. (34)).

The role of the $\Theta$ arc (Eq. (2)) deserves comment. Some training patterns, $Dq$ and $Pq$ in particular, exhibit concentric rings i.e., have structure. Because of the coordinate transformation applied to a sector, $r$ becomes the $x_2$ axis and rings become bright and dark stripes parallel to the $x_1$ axis. Power spectral density is thus higher in the vicinity of the $u_2$ axis. Integration along an arc $\Theta$, which is symmetric with respect to that axis, collects most of the power and captures structure-related information, which makes classification effective. Conversely, enhanced spectra derived from integration over $\Theta$'s symmetric about $u_1$ do not carry such information. As expected, no classification has been possible.

The basic assumptions that this investigation cannot disprove have been summarized by Ansätze 2–4 in Section 4.3.

Improvements in classification by the current algorithm may occur after a wider exploration of the control parameter set. Among classification methods, data and information fusion algorithms have to be developed in order to simultaneously deal with two or more sectors in the pattern. On the applications side, the recognition of airborne environmental materials is still awaiting a solution. Namely, the heterogeneity exhibited by $K$ patterns is apparent from Figs. 3 and 5. Different training designs, with possible input from a wider set of reference materials, may eventually lead to morphological classification or source apportionment of the $K$ material.

Finally, the extent to which the spectrum enhancement algorithm captures information relevant to obstacle inversion remains to be assessed.

## 7. Conclusion

The interpretation of TAOS patterns requires, in principle, the solution of an inverse problem in electromagnetics. A general inverse problem solver is beyond current abilities. As a consequence, the problem has been restated in statistical learning terms: rather than identifying the scatterer from the TAOS pattern, a pattern has been assigned to a class by means of a learning machine, where feature extraction interacts with linear classification.

Feature extraction has relied on spectrum enhancement, which includes the discrete cosine Fourier transform and some nonlinear operations. Linear classification has relied on principal components analysis and supervised training, based on the maximization of a suitable figure of merit.

All algorithms that analyse images, organize feature vectors, design classification experiments, carry out supervised training, assign unknown patterns to classes, fuse information from different training and recognition experiments, have been tested on large TAOS pattern data sets ( $> 3000$ patterns). The role of the many parameters that control the algorithms at different stages has been explored.

Classification experiments have been aimed at discriminating *B. subtilis* TAOS patterns ($Bq$ class) from patterns of other airborne materials and from interfering materials (diesel engine soot). The most satisfactory result corresponds to $\approx 20\%$ false $Bq$ negatives and $< 11\%$ false positives from all other materials.

The most effective operations have been thresholding TAOS patterns by intensity aimed at rejecting faint ones, and the formation of training sets from three or four pattern classes. The capabilities and limitations of the present automated classifier have been assessed to some extent.

In order to reduce false negatives, directions of improvement have been suggested, which include data and information fusion at the TAOS pattern level. Supervised training by more reference materials may lead to classification of the $K$ material.

Since the TAOS instrument has been designed to run unassisted for weeks, the automated classifier described herewith can be adapted for integration and to real-time operation.

## Acknowledgments

## References

[1] Nicoli DF, Hasapidis K, O'Hagan P, McKenzie DC, Wu JS, Chang YJ, et al. High resolution particle size analysis of mostly submicrometer dispersions and emulsions by simultaneous combination of dynamic light scattering and single particle optical sensing. In: Provder T, editor. Particle size distribution IIIassessment and characterization Washington, DC: American Chemical Society; 1998. p. 52–76.

[2] Holler S, Pan Y-L, Chang RK, Bottiger JR, Hill SC, Hillis DB. Two-dimensional angular optical scattering for the characterization of airborne microparticles. Opt Lett 1998;23(18):1489–91.

[3] Secker DR, Kaye PH, Greenaway RS, Hirst E, Bartley DL, Videen GL. Light scattering from deformed droplets and droplets with inclusions. I. Experimental results. Appl Opt 2000;39(27):5023–30.

[4] Videen GL, Sun W, Fu Q, Secker DR, Greenaway RS, Kaye PH, et al. Light scattering from deformed droplets and droplets with inclusions. II. Theoretical treatment. Appl Opt 2000;39(27):5031–8.

[5] Ramm AG. Multidimensional inverse scattering problems. Harlow: Longman; 1992.

[6] Colton D, Kress R. Inverse acoustic and electromagnetic scattering theory. Berlin: Springer; 1998.

[7] Mishchenko MI, Travis LD, Mackowski DW, T-matrix codes for computing electromagnetic scattering by nonspherical and aggregated particles ⟨http://www.giss.nasa.gov/crmim/t_matrix.html⟩.

[8] Mishchenko MI, Travis LD, Lacis AA. Scattering, absorption, and emission of light by small particles. Cambridge, UK: Cambridge University Press; 2002.

[9] Auger J-C, Aptowicz KB, Pinnick RG, Pan Y-L, Chang RK. Angularly resolved light scattering from aerosolized spores: observations and calculations. Opt Lett 2007;32(22):3358–60.

[10] Crosta GF, Camatini MC, Zomer S, Pan Y-L, Holler S, Bhaskara P, et al. Optical scattering (TAOS) by tire debris particles: preliminary results. Opt Exp 2001;8(6):302–7 ⟨http://www.opticsexpress.org⟩.

[11] Crosta GF, Zomer S, Pan Y-L, Holler S. Classification of single-particle two-dimensional angular optical scattering patterns and heuristic scatterer reconstruction. Opt Eng 2003;42:2689–701.

[12] Holler S, Zomer S, Crosta GF, Pan Y-L, Chang RK, Bottiger JR. Multivariate analysis and classification of two-dimensional angular optical scattering patterns from aggregates. Appl Opt 2004;43(33):6198–206.

[13] Crosta GF. Classification of single particle optical scattering patterns by the spectrum enhancement algorithm. In: Sedlacek A III, Christesen S, Combs R, Vo-Dinh T. Chemical and biological sensors for industrial and environmental security proceedings of SPIE 5994, SPIE, Bellingham, WA; 2005. p. 599402-1–12.

[14] Aptowicz KB, Pinnick RG, Hill SH, Y-L Pan, Chang RK. Optical scattering patterns from single aerosol particles at Adelphi, Maryland, USA: a classification relating to particle morphologies. J Geophys Res 2006;111:D12212.

[15] Marr D. Vision. San Francisco, CA: Freeman; 1982.

[16] Vapnik VN. The nature of statistical learning theory. Berlin: Springer; 1995.

[17] Cucker F, Smale S. On the mathematical foundations of learning. Bull AMS (NS) 2001;39(1):1–49.

[18] Crosta GF, Urani C, Fumarola L. Classifying structural alterations of the cytoskeleton by spectrum enhancement and descriptor fusion. J Biomed Opt 2006;11:024020-1–18.

[19] Crosta GF. Image analysis and classification by spectrum enhancement: new developments. In: Astola JT, Egiazarian KO. editors, Image processing: algorithms and systems VIII—proceedings of SPIE, vol. 7532, SPIE, Bellingham, WA; 2010. p. 75320L-01–12.

[20] Russ JC. Fractal surfaces. New York: Plenum; 1994.

[21] Meyer Y. Oscillating patterns in image processing and nonlinear evolution The fifteenth Dean Jacqueline B. Lewis Memorial Lectures, AMS, Providence, RI; 2001.

[22] Press WH, Teukolsky SA, Vetterling WT, Flannery BP. Numerical recipes in C—the art of scientific computing. second editionCambridge, UK: Cambridge University Press; 1992.

[23] Crosta GF. Texture analysis of phosphor screens. J Phys E: Sci Instr 1977;10:187–90.

[24] Crosta GF, Sung C, Kang B, Ospina C, Stenhouse P. Quantitative morphology of aluminum silicate nanoaggregates. Materials Research Society symposia proceedings, vol. 738, Materials Research Society, Warrendale, PA; 2003. p. G.15.1–G.15.6.

[25] Pearson K. On lines and planes of closest fit to systems of points in space. Lond Edinburgh Dublin Philos Mag J Science, Sixth Ser 1901;2:559–72.

[26] Hotelling H. Analysis of a complex of statistical variables into principal components. J Educ Psychol 1933;24:417–41.

[27] Bishop CM. Pattern recognition and machine learning.New York, NY: Springer Science + Business Media; 2006.

[28] Aptowicz KB, Pan Y-L, Chang RK, Pinnick RG, Hill SC, Tober RL, et al. Two-dimensional angular optical scattering patterns of microdroplets in the mid infrared with strong and weak absorption. Opt Lett 2004;29(17):1965–7.

[29] Kaye PH. Spatial light scattering analysis as a means of characterizing and classifying non-spherical particles. Meas Sci Technol 1998;9:141–9.

[30] Kress R, Rundell W. Inverse obstacle scattering using reduced data. SIAM J Appl Math 1998;59(2):442–54.

[31] Kress R. Uniqueness and numerical methods in inverse obstacle scattering, inverse problems in applied sciences. J Phys Conf Ser. London: IOP Publishing, 2007;73:012003_1–012003_16.

[32] Osher S, Rudin L, Fatemi E. Nonlinear total variation-based noise removal algorithms. Physica D 1992;60:259–68.