# SPONTANEOUS CLUSTERING IN THEORETICAL AND SOME EMPIRICAL STATIONARY PROCESSES [*]

T. Downarowicz[1], Y. Lacroix[2] and D. Léandri[2]

**Abstract.** In a stationary ergodic process, clustering is defined as the tendency of events to appear in series of increased frequency separated by longer breaks. Such behavior, contradicting the theoretical "unbiased behavior" with exponential distribution of the gaps between appearances, is commonly observed in experimental processes and often difficult to explain. In the last section we relate one such empirical example of clustering, in the area of marine technology. In the theoretical part of the paper we prove, using ergodic theory and the notion of category, that clustering (even very strong) is in fact typical for "rare events" defined as long cylinder sets in processes generated by a finite partition of an arbitrary (infinite aperiodic) ergodic measure preserving transformation.

## 1. INTRODUCTION: RECURRENCE AND CLUSTERING

Let us consider an aperiodic ergodic dynamical system $(X, \mathcal{B}, \mu, T)$, where $(X, \mathcal{B}, \mu)$ is a standard probability space, and $T : X \to X$ is measurable, almost surely $1 - 1$, preserves $\mu$, for which it is an ergodic transformation [10].

Let us next consider a measurable partition $\mathcal{P} = \{P_1, \ldots, P_n\}$ of the space $X$. To this partition we associate the natural symbolic factor of the system, using the coding map

$$c : x \in X \mapsto i \in \{1, \ldots, n\} \Leftrightarrow x \in P_i$$

which generates the factor map [10]

$$\pi : \begin{array}{ccc} X & \to & \{1, \ldots, n\}^{\mathbb{Z}} \\ x & \mapsto & (c(T^n x))_{n \in \mathbb{Z}}. \end{array}$$

Then if $\nu = \pi\mu$ and $\sigma : \{1, \ldots, n\}^{\mathbb{Z}} \to \{1, \ldots, n\}^{\mathbb{Z}}$ denotes the left shift map, we obtain the factor $\pi : (X, \mathcal{B}, \mu, T) \to (\{1, \ldots, n\}^{\mathbb{Z}}, \mathcal{C}, \nu, \sigma)$ which satisfies $\pi \circ T = \sigma \circ \pi$ (the $\sigma$-algebra $\mathcal{C}$ is the obvious one).

This standard procedure, based on the selection of a partition on $X$, transforms the initial abstract dynamical system into a symbolic dynamical system. The coding procedure is somewhat natural in that it produces the accessible dynamics through the selection of a finite partition of the space of observables $X$.

On the resulting symbolic space we inherit of standard subsets of the phase space $\{1, \ldots, n\}^{\mathbb{Z}}$, which we call *cylinders, or blocks*. A typical such, of length $p$, is obtained by selecting a pattern $w \in \{1, \ldots, n\}^p$, and defining

$$[w] = \{y \in \{1, \ldots, n\}^{\mathbb{Z}} : (y_0, \ldots, y_{p-1}) = (w_0, \ldots, w_{p-1})\}.$$

When the length $p$ of the cylinder $[w]$ increases, of course, the measure $\nu([w])$ tends to 0, which results from aperiodicity of the dynamics.

Whence long cylinder sets are prototypical rare events (*i.e.* small measure sets) in the symbolic dynamics we accessed by cross-ruling our initial space $X$ with the partition $\mathcal{P}$. The system is ergodic whence recurrence to any positive measure subset of $X$ must occur, a.s., and happens along a typical trajectory with frequency equal to the measure of the set (this is the ergodic theorem).

Our paper concerns the study of recurrence to rare events generated by partitions as above. If we set for given $B \in \mathcal{B}$ with $\mu(B) > 0$,

$$\tau_B(x) = \min\{k \geq 1 : T^k x \in B\},$$

then $\tau_B$ is $\mu$-a.s. well defined, integer valued. When $x \in B$, it is called the return time of $x$ to $B$, otherwise it is called the entry time of $x$ to $B$. Kac's theorem [6] states that $\sum_{k \geq 1} k\mu(\{\tau_B = k\} \cap B) = 1$, whence the random variable $\mu(B)\tau_B$ defined on the probability space induced on $B$ has expected value equal to 1. We therefore call $\mu(B)\tau_B$ the "normalised return time".

For rare events such as cylinder sets the distribution of return times or entry times to such events, though by the ergodic theorem has on the overall frequential distribution, may change quite a lot depending on the dynamics of the system. What is oftenly looked after is weak convergence of such distributions as the measure of the sets $B$ (resp. lengths of the cyinders) shrink to 0 (resp. go to $\infty$). This is because if such convergence holds then the limiting distribution provides information about recurrence to rare events in the dynamical system, as it is an approximation of distributions converging to it weakly.

Such a limit distribution is called an asymptotic for return times. We denote by $\tilde{F}_B$ the distribution function of the normalised return time to $B$, and by $F_B$ the distribution function of the entry time to $B$, *i.e.* of the random variable $\mu(B)\tau_B$. Following the same lines, asymptotics for entry times are analogously defined.

Possible asymptotics have been characterized only recently ([8] for return times [7], for entry times). An integral formula connecting asymptotics for entry times and the one for return times has been provided in [5], where it is proved that weak convergence of distributions, whenever it holds, must hold simultaneously for entry and return times. If $\tilde{F}$ is the weak limit distribution function for return times and $F$ is the one for entry times, then additionally

$$F(t) = \int_0^t (1 - \tilde{F}(s))\mathrm{d}s, \ t \geq 0. \tag{1.1}$$

Many research papers have been devoted to the study of return times of specific events, so-called cylinder sets, in stationary ergodic processes. We refer also the reader to expository papers [3] and [1] for further information. Most asymptotics along cylinder sets were found in mixing enough dynamical systems [10], but in all cases were proved to be exponential with parameter one, the distribution of which is the only fixed point of (1.1). Furthermore, in the treated cases, entropy of the system was positive [10].

An essential progress in understanding some phenomena towards asymptotics in processes with positive entropy has been obtained recently in [4], where it is proved that whatever the system, as soon as its entropy is positive, then any asymptotic $F$ for entry times along cylinder sets must satisfy

$$F(t) \leq 1 - \mathrm{e}^{-t}, \ t \geq 0. \tag{1.2}$$

This was interpreted as a first explanation to a complicated and misunderstood common-sense phenomenon known as "*the law of series*": indeed in any ergodic system, if $t > 0$ and $\mu(B) > 0$, if we define the variable

$$I(x) = \#\{0 \leq n \leq \frac{t}{\mu(B)} : T^n x \in B\},$$

then invariance of the measure implies that

$$\mathbb{E}(I) \approx t, \tag{1.3}$$

with uniform accuracy [D]. Considering an independent symbolic process so as to be one for which most randomness occurs, and using the well-known fact that in such system the asymptotics along cylinder sets exist and are exponential, we are led to argue that a small cylinder set $B$ has a tendency to occur presenting clusters more frequently than it would in an independent process, if given that $B$ occurs, the conditional expectation of $I$ increases, compared to what it is in the independent case.

That is to say, $B$ clusters (or appears "in series") in distribution, if

$$\mathbb{E}(I | I > 0) \geq \mathbb{E}_{Ind}(I_{Ind} | I_{Ind} > 0),$$

where the subscript $Ind$ refers to the independent process. Now using (1.3) and the fact that $\mathbb{E}(I | I > 0) = \mathbb{E}(I) \,/\, \mu(I > 0)$, observing moreover that $\mu(I > 0) = F_B(t)$, and finally approximating $F_{B,Ind}(t) \approx 1 - e^{-t}$, we deduce

$$B \text{ clusters} \Leftrightarrow F_B(t) \leq 1 - e^{-t}.$$

In [4] we introduced the notions of attracting, strong attracting, repelling, and neutral recurrences for a set $B$, in reference to the above explained comparison to the case of the independent process. Below is a purely naive illustration of the phenomena we are describing, interpreted as "along a typical orbit".



This is the interpretation along which it was stated in [4] that (1.2) can be understood differently, saying that in positive entropy systems, clustering for rare events must be at least what it reveals to be in the neutral independent case. In other words, "laws of series" (more frequent clustering for rare events) correspond to the natural behaviour for positive entropy systems.

This was completed in [4] by the study of asymptotics along cylinder sets for varying partitions $\mathcal{P}$ of the space $X$, so as to understand what generically (in a Baire category setting) could happen to be the case (the set of partitions can be turned to a Polish structured space).

It was proved that generically (on $\mathcal{P}$) there exists an upper density one set $\mathcal{N} \subset \mathbb{N}$ of lengths of cylinder sets along which asymptotics exist, and for entry times, converge to the degenerated distribution function $F \equiv 0$, the one corresponding to strong attracting above (enormous clusters).

In the present paper we go forward in this last direction and prove that this genericity holds without the positive entropy assumption. We also addressed the question to specialists, since the "law of series", frequently referred to under more pessimistic interpretations like "Murphy laws", is used by engineers and manufacturers and formalised under the label "Murphy proof devices", in aeronautics for instance. In fact other occurrences of Murphy proof procedures appear, and we reveal such in the last part of this paper, were we relate to some rare but clustering disturbances that occur on low speed flights of underwater gliders, without further understanding at this point of knowledge.

The only reasonable understanding for this generic clustering phenomenon we can talk out, at this point, more like a guess, is the observation that rare events, to occur, need favourable conditions, probably even rarer, and that therefore when they are collected, conditionally, the event has a better measure and comes through easier.

The paper is organised as follows: the next section introduces rigourous formulations about genericity. The next section presents the formulation and proof of our main result (Thm. 1). In the formulation we use the notion of strong clustering explained above in the Introduction. The last section is a presentation of an empirical example of a strongly clustered process occurring in marine technology, preceded by a short description of the performed experiment and its technical background.

## 2. Typicality of strong attraction without entropy assumptions

Let us go back to the original dynamical system $(X, \mathcal{B}, \mu, T)$ on which the observed symbolic dynamics is defined by the partition $\mathcal{P}$.

Denote by $\mathfrak{P}_l$ the collection of all measurable partitions $\mathcal{P}$ of $X$ into $l$ cells. In general there is no canonical measure on the space $\mathfrak{P}_l$. The meaning of "almost surely" with respect to a random partition does not have a definite meaning. Instead, we will adopt the topological approach, and "typicality" defined in terms of so-called *category*.

Recall that in a complete metric space a subset is called *residual* if it contains a dense $G_\delta$ set, equivalently, if its complement if of first category (*i.e.*, is a countable union of nowhere dense sets). The Baire category theorem asserts that the intersection of any countable collection of residual sets is still residual. Because similar property is enjoyed by sets of measure 1 in probability spaces, residual sets are considered a topological analog of sets of measure 1. Given a fixed residual set, its elements are referred to as "typical".

If $(X, \mathcal{B}, \mu)$ is a probability space, the *Rokhlin metric* endows $\mathfrak{P}_l$ with a structure of a complete metric space. The distance in this metric between two $l$-element partitions $\mathcal{P}, \mathcal{Q}$ is defined as

$$d(\mathcal{P}, \mathcal{Q}) = \inf_\pi \sum_{A \in \mathcal{P}} \mu(A \triangle \pi(A)),$$

where $\pi$ ranges over all bijections from $\mathcal{P}$ to $\mathcal{Q}$ and $\triangle$ denotes the symmetric difference of sets. A partition $\mathcal{Q}$ which is very close to $\mathcal{P}$ in this distance may be considered a slight perturbation of $\mathcal{P}$. Our theorem below applies to processes typical in the following sense: given an ergodic dynamical system $(X, \mathcal{B}, \mu, T)$ and some $l \in \mathbb{N}$, the set of all $l$-element partitions $\mathcal{P}$ of $X$ such that the generated process $(\{1, \ldots, \ell\}^\mathbb{Z}, \mathcal{C}, \nu_\mathcal{P}, \sigma)$ satisfies the assertion of the theorem is residual in the Rokhlin metric in $\mathfrak{P}_l$.

We will also use the phrase that a property $\Phi$ holds for all blocks of "majority" of lengths. By this we mean that there exists a set $\mathbb{N}_\Phi$ of *upper density* 1 such that if $n \in \mathbb{N}_\Phi$ then all blocks of of length $n$ satisfy $\Phi$. Upper density is defined as

$$\overline{D}(\mathbb{N}_\Phi) = \limsup_{n \to \infty} \frac{\#(\mathbb{N}_\Phi \cap [1, n])}{n}.$$

## 3. Formulation of the theorem and the proof

Below we give two formulations of the main result. The first one is short thanks to the terminology introduced above and appeals to the intuitive understanding of the subject. The latter formulation is the rigorous version without using the shortcut terminology.

**Theorem 3.1.** *In a typical ergodic process, for majority of lengths $n$, all rare elementary events of length $n$ reveal strong clustering.*

**Theorem 3.2.** *Let $(X, \mathcal{B}, \mu, T)$ be an ergodic not periodic dynamical system. Fix some natural $l \geq 2$. Then in the space $\mathfrak{P}_l$ of all $l$-element measurable partitions of $X$ endowed with the Rokhlin metric there exists a residual subset $\mathfrak{C}$ such that for every $\mathcal{P} \in \mathfrak{C}$ the generated process $(\{1, \ldots, \ell\}^\mathbb{Z}, \mathcal{C}, \nu_\mathcal{P}, \sigma)$ has the following property $\Phi$: there exists a set $\mathbb{N}_0 \subset \mathbb{N}$ of upper density one, such that for every $\epsilon > 0$ there is $n_\epsilon \in \mathbb{N}$ such that for every $n \in \mathbb{N}_0$, $n > n_\epsilon$ and every block $B$ of length $n$, the $\tilde{F}_B(\varepsilon) < \epsilon^2$.*

*Proof.* Fix $\epsilon > 0$ and $N \in \mathbb{N}$. Suppose a partition $\mathcal{P} \in \mathfrak{P}_l$ satisfies the following property $\Phi_{\epsilon,N}$: "For every $n \in [N, N^2]$ and every block $B$ of length $n$ holds $\tilde{F}_B(\varepsilon) < \epsilon^2$."

Clearly, if we perturb $\mathcal{P}$ very little, the property will still be satisfied. Thus $\Phi_{\epsilon,N}$ holds on an open set $\mathfrak{C}_{\epsilon,N}$ of partitions. Of course, the set

$$\mathfrak{C}_\epsilon = \bigcup_{N \geq 1} \mathfrak{C}_{\epsilon,N},$$

of partitions such that the same property holds for some $N$, is also open. The main effort in the proof will be to show that this set is also dense. Once this is done, the proof is complete, because then the dense $G_\delta$ set $\mathfrak{C}$ of partitions which fulfill the hypothesis of the theorem is obtained by intersecting the sets $\mathfrak{C}_\epsilon$ over countably many parameters $\epsilon$ converging to zero. Every element of this intersection satisfies $\Phi_{\epsilon,N}$ for arbitrarily small $\epsilon$ and some $N$ depending on $\epsilon$. Clearly, $N$ must grow to infinity as $\epsilon$ decreases to zero. Notice that for any infinite sequence of natural numbers $N$ the set $\bigcup[N, N^2]$ has upper density 1 in $\mathbb{N}$.

It remains to prove the density of the open set $\mathfrak{C}_\epsilon$. The proof (as most proofs of typicality) may look a bit artificial, as it is done by perturbing an arbitrarily chosen partition $\mathcal{P}$ in a very specific way, so that a highly particular partition is created. This perturbation is then shown to belong to the set $\mathfrak{C}_\epsilon$. It is important to realize that once the density is proved, the "largeness" properties of open dense sets imply that $\mathfrak{C}_\epsilon$ is represented in the vicinity of $\mathcal{P}$ by many more partitions, not only the artificially constructed perturbation.

We begin with a technical lemma.

**Lemma 3.3.** *In every ergodic not periodic dynamical system* $(X, \mathcal{B}, \mu, T)$, *for each sufficiently large* $r \in \mathbb{N}$ *there exists a "semiperiodic $r$-marker", i.e., a measurable set* $F_r$ *such that the return time to* $F_r$ *assumes almost surely only two values:* $r$ *and* $r + 1$.

*Proof.* We need a measurable and shift-invariant procedure dividing the trajectory of almost every point $x \in X$ into intervals of the two lengths $r$ and $r + 1$. (Shift-invariant means that the division of the trajectory of $Tx$ coincides with the division of trajectory of $x$ with shifted enumeration.) Once this is done, the set $F_r$ is defined as the collection of all points which have a division marker at the coordinate zero.

Using ergodicity and nonperiodicity (or the Rokhlin theorem, see [9]) it is very easy to construct a set of positive measure such that the return time to this set assumes almost surely only values larger than or equal to $r^2$. The times of visits to this set divide each trajectory into intervals of lengths at least $r^2$ (see also [8]). Every such interval can be further divided into intervals of two lengths $r$ and $r + 1$, and we can fix one such way for every length $m \geq r^2$ (for example, we can choose the division which minimizes the number of longer intervals and all longer intervals appear to the right of the shorter ones). This (clearly measurable) technique divides almost every trajectory into desired pieces in a shift-invariant way, as required. $\square$

We continue with the main proof; we are proving that $\mathfrak{C}_\epsilon$ is dense. Fix any $l$-element partition $\mathcal{P}$ and some $\delta > 0$. We will construct a perturbation $\mathcal{P}'$ within the distance $\delta$ from $\mathcal{P}$, which belongs to $\mathfrak{C}_{\epsilon,N}$ for some $N$. Since the partition has at least two elements, we select two of them and label them 0 and 1. Pick $L$ so large such that there exist $K = \frac{2}{\epsilon^2}$ different blocks $W'_k$ ($k = 1, 2, \ldots, K$) of length $\frac{L-1}{2}$, none of them equal to $000\ldots0$, each of measure at most $\frac{\delta}{2LK}$ (here we admit also blocks that do not appear in the system and hence have measure zero). Denote by $W_k$ the block $1W'_k1000\ldots0$ (with $\frac{L-1}{2}$ zeros) of length $L$. The blocks $W_1, \ldots, W_K$ also have measures at most $\frac{\delta}{2LK}$. Notice that any two (different or equal) blocks from this family never occur with an overlap. Let $r$ be an integer larger than $\frac{2L}{\delta}$. Let $N$ be larger than $2r + 2$. Every cylinder set of positive measure over a block $B$ of a length $n \geq N$ decomposes into a finite number of sets depending on the positioning of the $r$-markers (there is at least one such marker in every occurrence of $B$). We denote these sets by $B_i$ ($i = 1, 2, \ldots, j(B)$). Let $M$ be so large that all points $x$, except in a set $Z$ whose measure is smaller than $\frac{1}{3}$, satisfy the following: for every block $B$ of positive measure and length between $N$ and $N^2$, and for every $i = 1, 2, \ldots, j(B)$ the orbit of $x$ visits $B_i$ between times 0 and $M$ at least $\frac{3}{\epsilon}$ times. Now let $r_1 = KM$.

At this point we modify the partition $\mathcal{P}$ as follows: In the $\mathcal{P}$-name of (almost) every $x$ the interval between two consecutive $r_1$-markers splits into $K$ intervals of length $M$ (the last one may be of length $M + 1$). We call

them "sectors". In the sector number $k$ $(k = 1, 2, \ldots, K)$ we put the block $W_k$ immediately to the right of every $r$-marker (replacing whatever was there), and next we replace all occurrences of all blocks $W_j$ $(j = 1, 2, \ldots, K)$ also by $W_k$. Because the blocks $W_j$ do not overlap, such exchange happens on disjoint intervals. Notice that we have changed the partition $\mathcal{P}$ only on the $r$-marker set $F_r$ and its $L$ consecutive preimages by $T$ (jointly a set of measure $\frac{\delta}{2}$), and on the cylinders corresponding to the blocks $W_1, \ldots, W_K$, and also their $L$ preimages by $T$ (jointly another set of measure $\frac{\delta}{2}$). Thus the Rokhlin distance between $\mathcal{P}$ and the modified partition $\mathcal{P}'$ is at most $\delta$.

Let $B'$ be an arbitrary block of some length $n$ between $N$ and $N^2$ appearing with positive probability in the process generated by the modified partition. Consider an $x \in X$. By ergodicity, $B'$ occurs almost surely in the $\mathcal{P}'$-name of $x$ at some position $n$. The block $B'$ is long enough so that at least one $r$-marker occurs within its length. Next to this marker $B'$ contains one of the blocks $W_k$. Because $W_k$ does not appear in any sector other then the sector number $k$ (of an interval between two $r_1$-markers) $B'$ also can occur only in sectors that carry the number $k$. In the $\mathcal{P}$-name of $x$ (*i.e.*, before modification) at the position $n$ there is some block $B$ and also there is a fixed positioning of $r$-markers along this block, *i.e.*, $T^n(x)$ belongs to some $B_i$. We know that the orbit of $x$ visits $B_i$ at least $\frac{3}{\epsilon}$ times in each sector, except the cases when at the beginning of the sector the orbit falls into $Z$ (which happen with probability $< \frac{1}{3}$). If this sector happens to be the sector number $k$, in the modified partition every such visit generates an occurrence of $B'$. Thus we conclude the following: The block $B'$ occurs in the $\mathcal{P}'$-name of every $x \in Z$ at least $\frac{3}{\epsilon}$ times in more than $\frac{2}{3}$ of all sectors number $k$ and zero times in sectors carrying other numbers (see Fig. 1; stars indicate the $r_1$-markers).

$$
\begin{array}{c}
\overset{*}{\phantom{x}} \qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \overset{*}{\phantom{x}} \\
\cdots|\cdots\cdots|\cdots\cdots|\cdots\cdots|\cdots\cdots|_{B'B'B'}|\cdots\cdots|\cdots\cdots|\cdots\cdots|\cdots\cdots|\cdots\cdots|\cdots\cdots|\cdots\cdots|\cdots\cdots|_{B'B'B'}|\cdots\cdots|\cdots \\
\quad s.\,1 \quad\ s.\,2 \qquad\quad s.\,k \qquad\qquad\qquad\qquad\quad s.\,K \quad\ s.\,1 \quad\ s.\,2 \qquad\quad s.\,k
\end{array}
$$

FIGURE 1. The distribution of the occurrences of $B'$.

It remains to compute the intensity, normalize the waiting time, and prove that $\tilde{F}_{B'}(\epsilon) \leq \epsilon^2$. The intensity of the occurrences of $B'$ is at least $\frac{2}{3}\frac{3}{\epsilon KM} = \frac{2}{\epsilon KM}$. (the multiplier $\frac{2}{3}$ takes into account the visits in $Z$). The waiting time for the signal is smaller than $M$ only when the zero coordinate falls in a sector number $k$ or in the preceding sector. This happens with probability at most $\frac{2}{K} = \epsilon^2$. Otherwise the normalized waiting time is larger than $M\frac{2}{\epsilon KM} = \epsilon$. This ends the proof. $\qquad\square$

## 4. Clustering in an empirical process

One of the multiple illustrations of the law of series appears in the precise study of the behavior of an "underwater glider", one of the newest inventions in marine technology. The underwater glider is an autonomous robot able to cross thousands of miles without any propeller [2]. By changing the buoyancy of the ballasts, the vehicle tends to "fly" alternatively downwards or upwards to the surface. As the vehicle is equipped with wings – hence has some gliding ability, the variations of buoyancy are changed into horizontal motion. Very little energy is needed to modify the buoyancy and, as a consequence, the linear flight between the top and bottom points of the trajectory of the glider is costless in terms of energy (see Fig. 2).

This is the positive aspect of this brand new technology. However, during our numerous and lengthy tests at the sea we discovered some strange side effects, mainly linked to the fact that the glider flies at a very low speed and therefore is very sensitive to hydrodynamic phenomena.

It was discovered [2] that even though one makes no doubt that the flow around the glider is linear, and the physical system (glider plus flow) is not chaotic (zero entropy), some heavy turbulences appear during short periods of time. They are encountered under a conjunction of circumstances with on ocean scale are pretty rare. The point is that it has been observed that these rare circumstances produce heavy turbulences, and those turbulences have a tendency to repeat overexpectedly. After a while, the stable flow returns. The authors
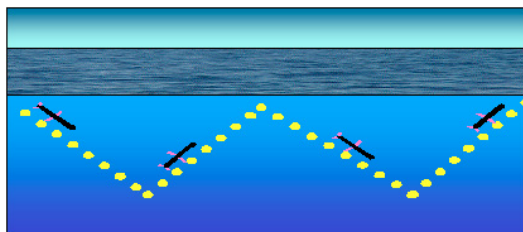
FIGURE 2. (Color online) The trajectory of a low energy underwater glider.

decided to call theses turbulences "germs", due to some similarities with the fact that the cavitations originating the perturbations could be activated by some "germs" in the water, and "pulsed out" somewhat once at a time.

The authors were unable to predict, and still aren't, when and why these turbulences appear in a given homogeneous layer of water in the ocean. What they know by experience, is that once an unexpected turbulence appears, then others follow repeatedly with short gaps in time until the series stops for a new long period of smooth linear flow. This corresponds exactly to the pattern described as "strong clustering".

This series type behavior of appearances of turbulences has been integrated so as to produce, when the first turbulence of a series appears, some warning, and consequently for autonomous vehicles, the control law of the underwater vehicle has been reinforced in these flight specific phases.

The next step would be to produce some statistical inference, and derive from such risk estimation for stability of trajectory under reinforced control laws of different levels of security. This has of course immediate industrial consequences. But it is a difficult mathematical challenge.

## REFERENCES

[1] M. Abadi and A. Galves, Inequalities for the occurrence times of rare events in mixing processes. The state of the art. Inhomogeneous random systems. *Markov Process. Relat. Fields* **7** (2001) 97–112.

[2] D. Brutzman, C. Deltheil, E. Hospital and D. Leandri, An optical guidance system for the recovery of an unmanned underwater vehicle. *IEEE – J. Oceanic Engineering* (2000).

[3] Z. Coelho, Asymptotic laws for symbolic dynamical systems, Topics in symbolic dynamics and applications. *London Math. Soc. Lect. Note Ser.* **279**, 123–165. Cambridge University Press (2000).

[4] T. Downarowicz and Y. Lacroix, The Law of Series, `http://arXiv.org/abs/math/0601166`.

[5] N. Haydn, Y. Lacroix and S. Vaienti, Entry and return times in ergodic aperiodic dynamical systems. *Ann. Probab.* **33** (2005) 2043–2050.

[6] M. Kac, On the notion of recurrence in discrete stochastic processes. *Bull. Amer. Math. Soc.* **53** (1947) 1002–1010.

[7] M. Kupsa and Y. Lacroix, Asymptotics for hitting times. *Ann. Probab.* **33** (2005) 610–619.

[8] Y. Lacroix, Possible limit laws for entrance times of an ergodic aperiodic dynamical system. *Israel J. Math.* **132** (2002) 253–263.

[9] V.A. Rokhlin, Selected topics from the metric theory of dynamical systems. *Amer. Math. Soc. Transl.* **2** (1966) 171–240.

[10] P. Walters, *Ergodic theory – Introductory lectures*, in *Lect. Notes Math.* **458**. Springer-Verlag, Berlin (1975).