

AutoQS v1: Automatic parameterization of QuickSampling based on training images analysis

Mathieu Gravey^{1,2,3}, Grégoire Mariethoz¹

¹ University of Lausanne, Faculty of Geosciences and Environment, Institute of Earth Surface Dynamics, Switzerland

² Institute for Interdisciplinary Mountain Research, Austrian Academy of Sciences, Innsbruck, Austria

³ Department of Physical Geography, Faculty of Geosciences, Utrecht University, Utrecht, Netherlands

Correspondence to: Mathieu Gravey (research@mgravey.com)

Highlights

- Adaptative calibration as a function of the simulation progression
- Calibration depends on each training image
- Robust parameterization based on a rapid prior analysis of the training image

Abstract. Multiple-point geostatistics are widely used to simulate complex spatial structures based on a training image. The practical applicability of these methods relies on the possibility of finding optimal training images and parametrization of the simulation algorithms. While methods for automatically selecting training images are available, parametrization can be cumbersome. Here, we propose to find an optimal set of parameters using only the training image as input. The difference between this and previous work that used parametrization optimization is that it does not require the definition of an objective function. Our approach is based on the analysis of the errors that occur when filling artificially constructed patterns that have been borrowed from the training image. Its main advantage is to eliminate the risk of overfitting an objective function, which may result in variance underestimation or in verbatim copy of the training image. Since it is not based on optimization, our approach finds a set of acceptable parameters in a predictable manner by using the knowledge and understanding of how the simulation algorithms work. The technique is explored in the context of the recently developed QuickSampling algorithm, but it can be easily adapted to other pixel-based multiple-point statistics algorithms using pattern matching, such as Direct Sampling or Single Normal Equation Simulation (SNESIM).

1 Introduction

Geostatistics is extensively used in natural sciences to map spatial variables such as surface properties (e.g., soils, geomorphology, meteorology) and subsurface geological features (e.g. porosity, hydraulic conductivity, 3D geological facies). Its main applications involve the estimation and simulation of natural phenomena. In this paper, we focus on simulation approaches.

Traditional two-point geostatistical simulations preserve the histogram and variogram inferred from point data (Matheron, 1973). However, inherent limitations make the reproduction of complex structures difficult (Gómez-Hernández and Wen, 1998; Journel and Zhang, 2006). Multiple-point statistics (MPS), by accounting for more

37 complex relations, enables the reproduction of such complex structures (Guardiano and Srivastava, 1993), but
38 comes with its own limitations (Mariethoz and Caers, 2014). The main requirements for using MPS algorithms
39 are 1) analog images (called training images) and 2) appropriate parametrization. While training images can often
40 be provided by expert knowledge, and several methods have been proposed to automatically select one or a subset
41 of appropriate training images among a set of candidates (Pérez et al., 2014; Abdollahifard et al., 2019). However,
42 the parametrization of an MPS algorithm depends not only on the chosen training image but also on the specifics
43 of the algorithm. This makes the task of finding good parametrization cumbersome, and therefore users often have
44 to resort to trial-and-error approaches (Meerschman et al., 2013). Here we will mainly focus on QuickSampling
45 (QS) (Gravey and Mariethoz, 2020) which has as two main parameters: n that defines the maximum number of
46 conditional data points to consider during the search process, and k which is the number of best candidates from
47 which to sample the simulated value. Additionally, QS supports a kernel that allows weighting each conditioning
48 pixel in the pattern based on its position related to the simulated pixel. Direct Sampling (DS) has for parameters:
49 n which has an identical role as in QS, th that represents the pattern acceptance threshold, or the degree of
50 similarity between local data patterns and the training image, and f the maximum proportion of the image that can
51 be explored for each simulated pixel. In summary, n controls the spatial continuity, and k or th and f control the
52 variability.

53 Over the last few years, several studies have addressed the challenge of automatically finding appropriate
54 parameters for MPS simulation. These can be categorized in two approaches. The first approach is to assume that
55 an optimal parametrization is related to the simulation grid (including possible conditioning data), the training
56 image and the MPS algorithm. In this vein, Dagan et al. (2018) proposed a method that uses the known hard
57 data from the simulation grid as a reference for computing the Jensen-Shannon divergence between histograms.
58 Following this, they employ a simulated annealing optimization to update the MPS parameters until the metrics
59 achieve the lowest divergence. This method is flexible enough to be adapted to any other metric. The second type
60 of approach assumes that the parametrization is only related to the training image and the MPS algorithm. Along
61 these lines, Baninajar et al. (2019) propose the MPS Automatic Parameter Optimizer (MPS-APO) method based
62 on the cross-validation of the training image (TI) to optimize simulation quality and CPU cost. In this approach,
63 artificially generated gaps in the high gradient areas of the training image are created, and a MPS algorithm is used
64 to fill those gaps. The performance of a particular parameterization is quantified by assessing the correspondence
65 between the filled and original training data. By design, this approach is extremely interesting for gap-filling
66 problems. The authors state that it can be used for the parametrization of unconditional simulations; however, the
67 use of limited gaps cannot guarantee the reproduction of long-range dependencies. Furthermore, due to the design
68 of the framework for generating gaps, only MPS algorithms able to handle gap-filling problems can be used.

69 While both approaches yield good results based on their objective functions, they all rely on a stochastic
70 optimization process, therefore the duration of the optimization process cannot be predetermined or controlled by
71 the user. Furthermore, an objective function is needed, which can be difficult because it depends on the training
72 image used: many metrics can be accounted for in the objective function, such as histogram, variogram, pattern
73 histogram, connectivity function, Euler characteristic, etc., (Boisvert et al., 2010; Renard and Allard, 2013; Tan et
74 al., 2013) or a weighted combination of these. Similarly, one has to define meta-parameters linked to the
75 optimization algorithm itself, such as the cooling rate in simulated annealing or maximum number of iterations.
76 As a result, MPS parameter optimization approaches tend to be complex and difficult to use.

77 In this contribution, we propose a simplified optimization procedure for simulating complex systems. Rather than
78 using a complex optimization algorithm, our approach focuses on finding optimal parameters to accurately
79 simulate a single pixel in the system. The underlying principle of our approach is that if each pixel is accurately
80 simulated, the resulting sequence of pixels will converge to an accurate representation of the real-world system
81 being simulated. The goal is therefore to find the optimal parameters to simulate a single pixel using the training
82 image as the only reference. Baninajar et al. (2019) showed that computing the prediction error (i.e., the error
83 between the simulation and the reference) is an appropriate metric to identify optimal parameters. To find the
84 optimal parameters for simulating a single pixel, we propose an exhaustive exploration of the parameter space and
85 a computation of the prediction error between the simulation and the reference image.

86 The remainder of this paper is structured as follows: Section 2 presents the proposed method. Section 3 evaluates
87 the approach in terms of quantitative and qualitative metrics. Finally, section 4 discusses the strengths and
88 weaknesses of the proposed approach and presents the conclusions of this work.

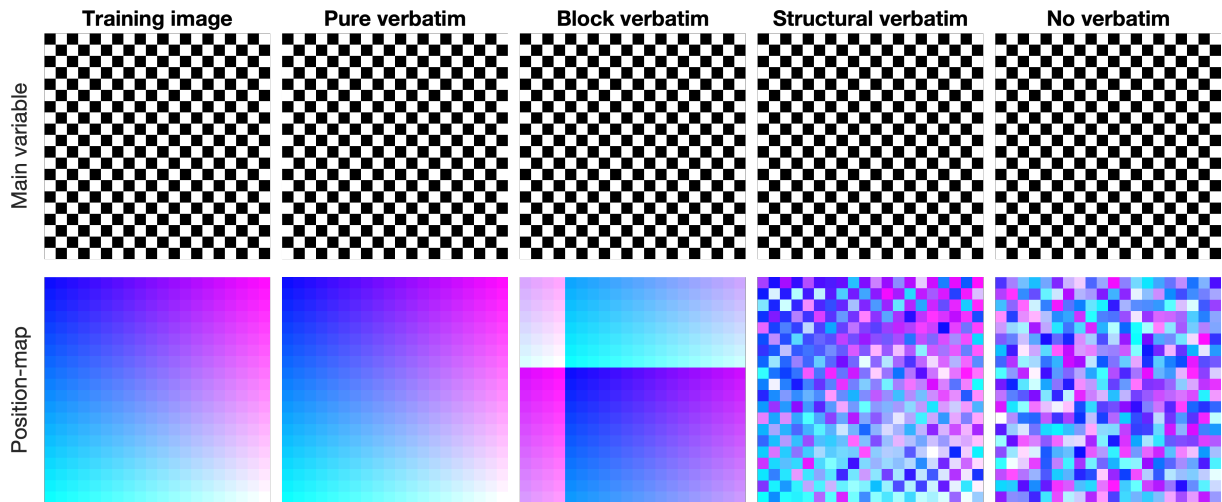
89 **2 Understanding and Addressing Verbatim Copy in Multiple Point Simulation**

90 The principle underlying multiple point simulation is that the neighborhood of a given pixel x (the pattern
91 generated by known or previously simulated pixels) is informative enough to constrain the probability density
92 function of the value $Z(x)$. This requires a training image with several pattern repetitions. The Extended Normal
93 Equation Simulation (ENESIM) algorithm (Guardiano and Srivastava, 1993) computes the full probability
94 distribution for each simulated pixel. To ensure that enough samples are used, the SNESIM (Strebelle, 2002) and
95 the Impala (Straubhaar et al., 2011) algorithms include a parameter to define a minimum number of patterns
96 replicates. Direct Sampling (DS) (Mariethoz et al., 2010) adopts a different strategy by allowing for the interrupted
97 exploration of the training image. It includes a distance threshold parameter that defines what is an acceptable
98 match for a neighborhood, however, too small a threshold typically results in a single acceptable pattern in the
99 training image, leading to exact replication of parts of the training image, a phenomenon known as verbatim copy.
100 To reduce this issue, a parameter f is introduced controlling the fraction of the explored training image.
101 QuickSampling (QS) (Gravey and Mariethoz, 2020) also suffers from verbatim copy when the number of candidate
102 patterns is set to $k = 1$, the authors recommend the use of $k > 1$, and highlight that k is similar to the number of
103 replicates in SNESIM or IMPALA. A value $k = 1.5$ in QS can be seen as SNESIM with a minimum number of
104 replicates of 1 for 50% of the simulated values and 2 for the remaining values.

105 The definition of verbatim copy is the unintended pasting of a large section from the training image to the
106 simulation (patch-based approaches do so intentionally, e.g. (Rezaee et al., 2013)). This means that the relative
107 position of the simulated values is the same as that in the training image. This occurs when the neighborhood
108 constraints on the simulated pixels are too strong and only the exact same patterns as those in the training image
109 are acceptable. To detect this issue, a common strategy is to create a position map (similar to the index map),
110 which represents the provenance of simulated values by mapping their original coordinates in the training image,
111 as shown in Figure 1.

112 Figure 1 illustrates the most common forms of verbatim copy. The pure verbatim (the most common type of
113 verbatim copy) is a simple copy of a large part of the image, with all pixels in the same order inside of the patches.
114 Block verbatim typically appears when there are many replicates of a very specific type of pattern in the training
115 image and few replicates of all other patterns. Consequently, the MPS algorithm uses common patterns for

116 transitioning between copied blocks resulting from rare patterns. Structural verbatim occurs when the copied
 117 portion spreads throughout the simulation without giving a direct impression of copying (e.g., pure verbatim over
 118 a subset of pixels). Structural verbatim tends to appear when large-scale structures are unique in the training image,
 119 which often allows a visually satisfying image to be quickly obtained, but with large non-stationary features
 120 identical to the training image. Often, users are willing to allow verbatim on large-scale structures, but this can
 121 easily introduce bias between simulations. This is one of the hardest types of verbatim to detect. Typically, this
 122 can occur when the maximum neighborhood radius is too large, leading to the duplication of large structures in
 123 the initial phase of the simulation. Finally, no verbatim, which is the expected result of simulations, occurs when
 124 the position of pixels does not have any particular structure (i.e. their position is unpredictable).



125
 126 **Figure 1 Visualization of verbatim copies using a position map. This is an extreme case that highlights that verbatim is**
 127 **not defined by the values simulated but by their position in the training image.**

128 **3 Method**

129 The objective of the approach presented here is to find an optimal set of parameters using only the training image
 130 and knowledge of the simulation algorithm's mechanics. The simulation algorithm is not used in this context; in
 131 fact, simulations are not required to obtain a proper calibration with the proposed method. The main target
 132 application of the presented approach is the pattern matching simulation algorithm QuickSampling (QS), where
 133 the values, at a pixel scale, are directly sampled from the training image. The method is suitable for the simulation
 134 of continuous and/or categorical variables.

135 Simulation algorithms such as QS, can be summarized by Algorithm 1. The key operation occurs at Line 3, which
 136 is when the algorithm searches for an optimal match based on the neighboring conditioning data.

137 **Algorithm 1 The sequential simulation algorithm. In gray the parametrization for QS.**

138
139 **Inputs:**
140 T : training images
141 S : simulation grid, including the conditioning data
142 P : simulation path
143 θ : parametrization
144 n : number of neighbors
145 k : the number of best candidates
146 ω : the kernel, by default uniform
147 1. **For** each unsimulated pixel x following the path P :
148 2. Find the neighborhood $N(x)$ in S composed of the $n(\theta)$ closest neighbors
149 3. Find a candidate in T those matches $N(x)$ using the parametrization θ
150 4. Assign the value of the selected candidate to x in S
151 **5. End**

152
153 Here, we propose a divide and conquer approach that splits any pixel-based sequential simulation into its atomic
154 operation: the simulation of a single pixel. We assume that if all pixels are perfectly simulated, then the resulting
155 simulation should also be good. By a perfectly simulated pixel, we mean a pixel that respects the conditional
156 probability distribution. When simulating a pixel, there may be numerous potential valid values, but at the very
157 least, there should be one valid value, i.e., the conditional probability distribution should be represented in the
158 data. This can be formalized by the following condition:

$$|\{A|P(A|N(x)) > 0\}| \geq 1 \quad (1)$$

159 where $|\cdot|$ represents the cardinality of a set. $P(A|N(x))$ denotes the probability of A (a given value) knowing
160 $N(x)$, the neighborhood.

161 The proposed approach consists of finding a set of parameters that results in accurate samples for each pattern. At
162 the same time, we want to avoid systematically sampling perfect matches (the exact same neighborhood is
163 available in the training image), which results in verbatim copy.

164 The search for the optimal parametrization is carried out by exhaustive exploration (**Error! Reference source not
165 found.**), and the choice of optimal parameters is based on a prediction error defined as the difference between the
166 original value of the pattern and the value of the selected pattern in the training image.

168 **Algorithm 2**

169 **Inputs:**
170 D : list of stages of the simulation (i.e. pattern decimation levels, equivalent to fractions of the simulation path)
171 θ : list of discretized parameters
172 T the training images
173 \mathcal{V} a set of random positions (in practice we generated the random position on the fly)
174 1. **For** each possible combination of D and θ **do** for all $\nu \in \mathcal{V}$:
175 2. Sample a neighborhood $N(\nu)$ from T and decimate it according to stage D
176 3. Using θ , find a candidate in T that matches $N(\nu)$, excluding for ν itself
177 4. Compute the error ε between the selected candidate and $Z(\nu)$
178 **5. End**
179 **6.** Analyze the errors ε to determine the best θ for each D .

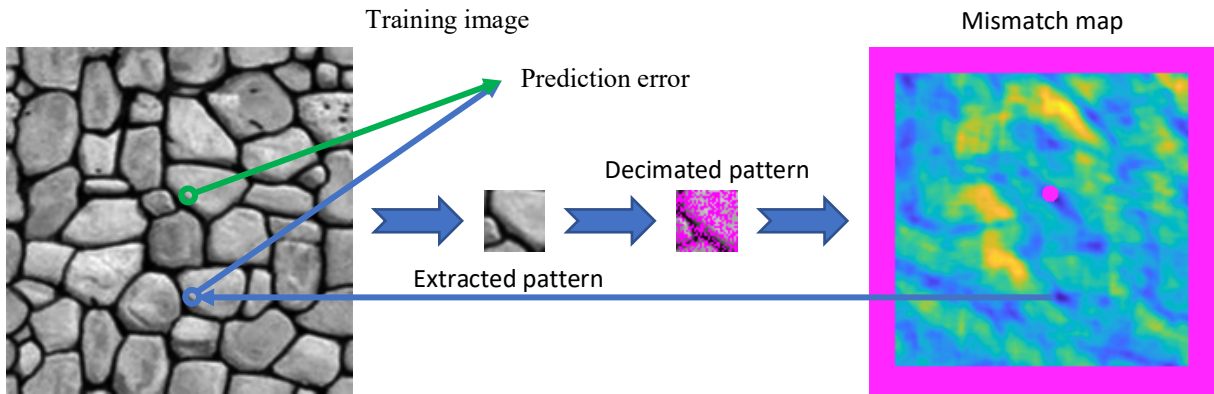
180
181 The proposed algorithm explores a discretized parameter space θ (**Error! Reference source not found.**, Line 1)
182 (e.g., for QS: n, k, ω). While this discretization is natural for some parameters, such as n that is an integer, it can
183 require an explicit discretization for other parameters, such as the kernel in QS (or th in DS). Furthermore, a key
184 component of our method is the exploration of the parameter space for several representative stages D of the

185 simulation (**Error! Reference source not found.**, Line 1). In the case of a random path, the progress of the
 186 simulation is directly related to the density of the neighborhoods, i.e., when $x\%$ of the pixels are simulated, in
 187 average $x\%$ of neighbors are informed. To reproduce this behavior, at each stage D , we randomly decimate patterns
 188 extracted from the TI by keeping only $x\%$ pixels informed. For each combination D and θ , multiple measures over
 189 a set of random locations \mathcal{V} ($500 < |\mathcal{V}| < 10000$) are computed in Lines 1-5 in **Error! Reference source not**
 190 **found.**, with their mathematical expression shown in Equation 2:

$$191 \quad \varepsilon(\theta, D, T) = \sqrt{\frac{1}{|\mathcal{V}|} \sum_{\nu \in \mathcal{V}} \left(Z(\nu) - Z\left(\text{Cand}(\theta, N(\nu, D))\right) \right)^2} \quad (2)$$

192 where $\text{Cand}(\theta, N)$ returns a single candidate position for a given neighborhood N and follows the parametrization
 193 θ . $N(\nu, D)$ denotes a neighborhood around ν that is decimated according to stage D . \mathcal{V} represents a random set
 194 of positions in the training image, and $Z(\nu)$ refers to the actual value at position $\nu \in \mathcal{V}$ in the training image. To
 195 avoid parameters that generate verbatim copy of the training image, the position ν and its direct neighbors (in a
 196 small radius (here 5 pixels) are excluded from the set of potential candidates. The set of candidates considering
 197 this exclusion is denoted by $T \setminus \{\nu\}$ in Equation 2. Furthermore, in the case of equality between several optimal
 198 options, we set as a rule to take the cheapest parameter set in terms of computational cost (e.g., the smallest n).
 199 **Error! Reference source not found.** graphically represents the entire algorithm. Finally, for each stage
 200 considered, the set of parameters with the minimum associated error ε is considered optimal (**Error! Reference**
 201 **source not found.**, Line 6):

$$202 \quad \varepsilon(\theta_{optimal}, D, T) = \min_{\theta} \varepsilon(\theta, D, T) \quad (3)$$



203
 204 **Figure 2 All steps for a single pattern, summarizing Error! Reference source not found., Lines 2-4.**
 205

206 4 An efficient implementation

207 In practice, the implementation of Algorithm 2 separates θ into two parameter subsets: θ_h and θ_s . The θ_h subset
 208 consists of all parameters that influence the calculation of a single pattern match, which varies depending on the
 209 algorithm used. For instance, in QS, it includes the number of neighbors n and the kernel ω , while in DS, it
 210 comprises the threshold th and n . The other hand, θ_s encompasses parameters related to the sampling process of

211 the training image. For QS, this includes the number of candidates to keep k , while for DS, it involves the fraction
 212 f of the training image being scanned.

213 Our implementation precomputes and stores all matches for a specific θ_h parameterization (e.g., a value of n and
 214 all matches for k). Consequently, the saved matches of θ_h can be employed to swiftly evaluate all options for the
 215 parameters in $\theta = \theta_h \times \theta_s$ (e.g., we can process for $k = 1, 2, 3, \dots k_{max}$). This two-phase approach considerably
 216 decreases redundant calculations.

217 The algorithm can be further accelerated by terminating the estimation of ε if the error remains at a high level after
 218 assessing only a small amount of samples from \mathcal{V} (here set to 500). To this end, we increase \mathcal{V} for the parameter
 219 combinations of interest, i.e., parametrization with potentially the lowest ε . This entails iterating and verifying at
 220 each step whether additional computations are required. Only places respecting following inequality are refined
 221 with extra measures:

$$222 \quad \varepsilon(\theta, D, T) - \varepsilon(\theta_{min}, D, T) < \frac{1}{2}\sigma(\theta, D, T) + \frac{1}{2}\sigma(\theta_{min}, D, T) \quad (4)$$

223 With

$$224 \quad \varepsilon(\theta_{min}, D, T) = \min_{\theta} \varepsilon(\theta, D, T)$$

$$225 \quad \sigma(\theta, D, T) = \sqrt{\frac{1}{|\mathcal{V}|} \sum_{\vartheta \in \mathcal{V}} \left(\left(Z(\vartheta) - Z\left(\text{Cand}(\theta, N(\vartheta, D)) \right) \right) - \varepsilon(\theta, D, T) \right)^2}$$

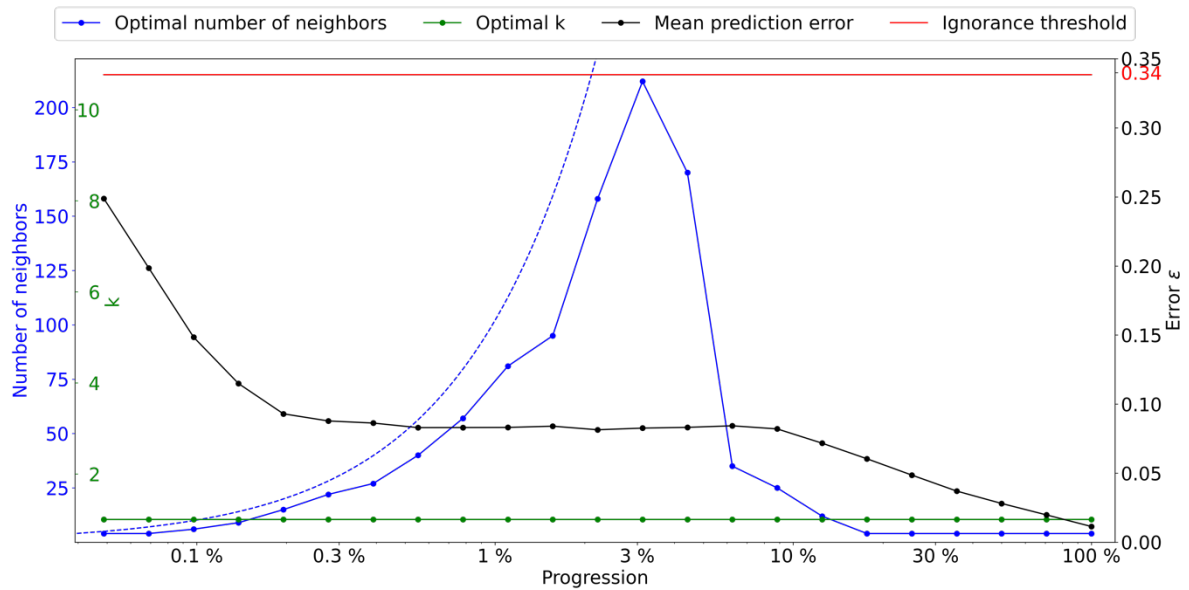
226 With $\varepsilon(\cdot)$ the error, and $\sigma(\cdot)$ represent the standard deviation of all differences, between estimated and true values.

227 5 Results

228 5.1 Optimization of 2 parameters

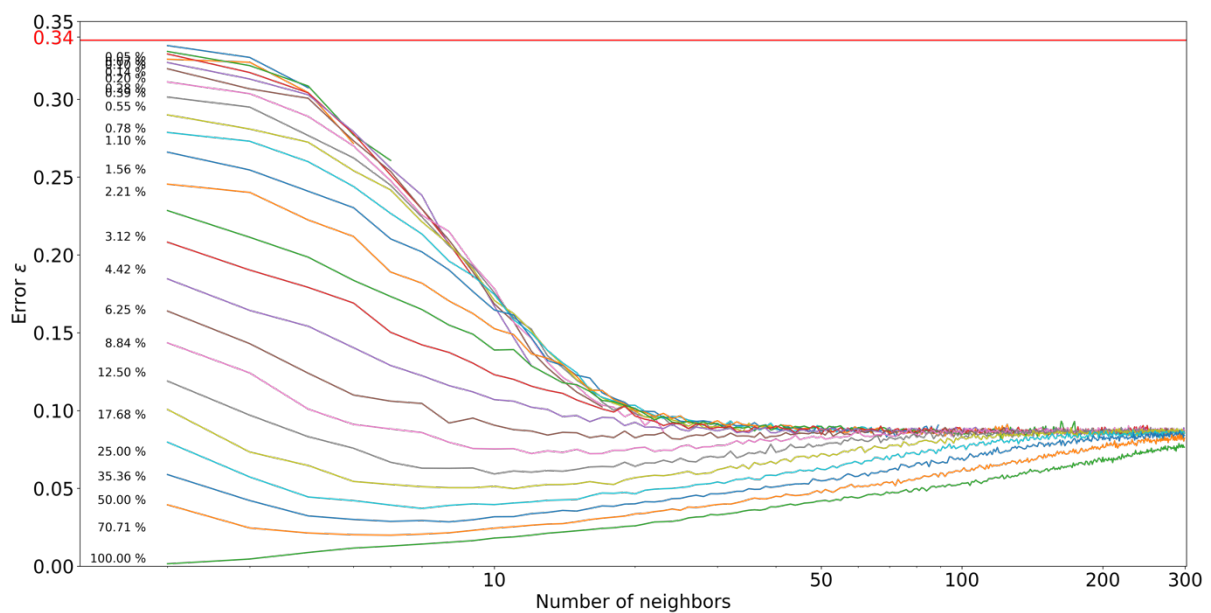
229 All experimental tests in this section are performed using the training image shown in **Error! Reference source**
 230 **not found.**, and the stages D are distributed following a logarithmic scale.

231 As a first test, we use the configuration $\theta_h = \{n\}$, and $\theta_s = \{k\}$. The kernel ω is defined as uniform, meaning that
 232 it has a constant value and is not part of the optimization. The outcome is represented in Figure 3, with the optimal
 233 number of candidates k and number of neighbors n as a function of the density D , which is assimilated to the
 234 progression during the simulation. The ignorance threshold is defined as the average error between elements of
 235 the marginal distribution. It represents the error value at which no further information can be derived from the
 236 neighborhood, meaning that the simulated values can equivalently be drawn from the marginal distribution.



237
 238 **Figure 3 Optimal parameters for QS (k in green and number of neighbors in blue) as a function of the progression,**
 239 **with the associated prediction error (in black). The red line represents the ignorance threshold. The dashed blue line**
 240 **indicates the average maximal number of neighbors..**

241
 242 The optimal k remains small (in fact 1) throughout the simulation, which is probably due to the limited size of the
 243 training image in this case. It seems important to use many neighbors in the early stages of the simulation. The
 244 number of neighbors increases until approximately 3% of the simulation. This is followed by a subsequent drastic
 245 reduction, indicating that once the large structures are informed, only the few direct neighbors are important. It
 246 seems logical that MPS algorithms simulate large structure first and then smaller patterns in a hierarchical manner
 247 where each smaller structure is part of the larger one. We however note that it remains generally difficult to predict
 248 the optimal settings as a function of the simulation stage. This indicates that the use of a single parametrization for
 249 the entire MPS simulation is generally suboptimal, and the parameters should be adapted as the simulation
 250 progresses.



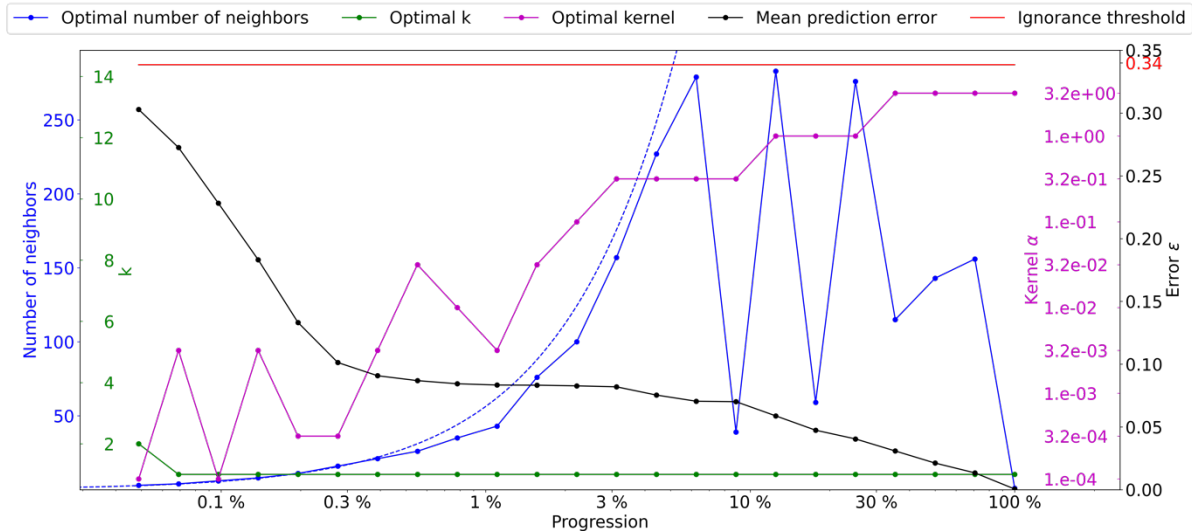
251
 252 **Figure 4 Pattern error as a function of the number of neighbors n , with $k = 1$, where each curve represents a**
 253 **neighborhood density D .**

254
 255
 256
 257
 258
 259
 260
 261
 262

Figure 4 shows the evolution of ε as a function of the number of neighbors n and the simulation progression D . Two regimes are visible: in the first percentages of the simulation, each extra neighbor is informative and improves simulation quality. However, as the neighborhoods become denser, the importance of spatial continuity takes over, and only the few neighbors are really informative. This two-step process is expected, as random large-scale features are generated first, and then the image is filled with consistent fine-scale structures. Furthermore, it shows that using a large number of neighbors at the end of the simulation generates suboptimal results, which could explain the small-scale noise that is sometimes visible in some MPS simulations.

263 **5.2 Optimization of 3 parameters**

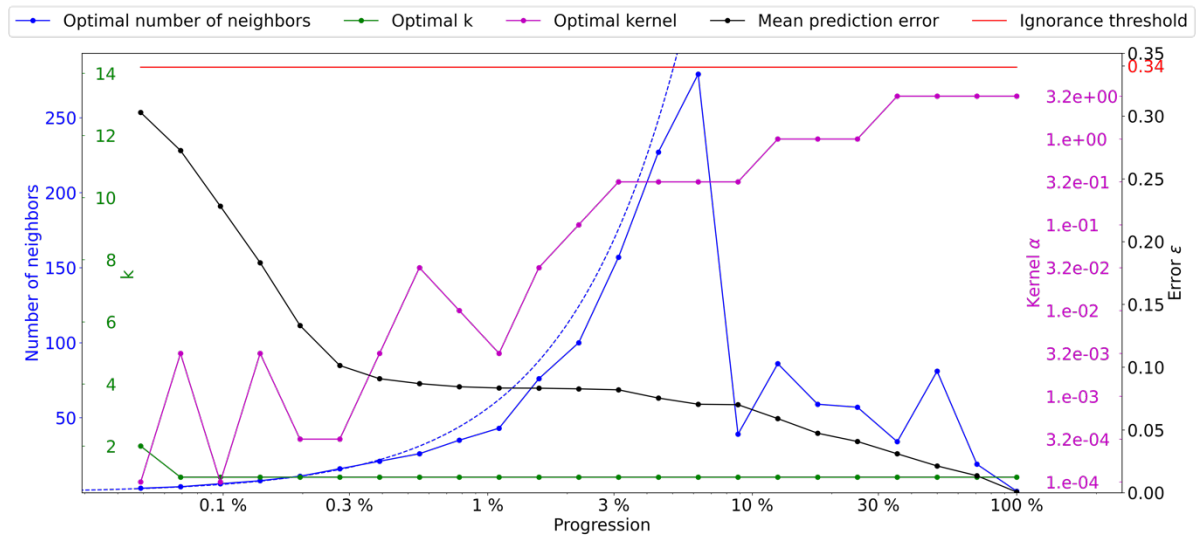
264 Here, we use the following configuration $\theta_h = \{n, \alpha\}$ and $\theta_s = \{k\}$, and we consider kernels as having a radial exponential shape, i.e. $\omega_i = e^{-\alpha \cdot d_i}$. The weight of a given position i in the kernel ω is defined as ω_i , and its distance to the kernel center as d_i .



267
 268 **Figure 5 Optimal parameters for QS (k in green, number of neighbors in blue, and best kernel in magenta), as a function**
 269 **of the simulation progress, with the associated prediction error (in black). The dashed blue line indicates the average**
 270 **density for the neighborhood considered. The ignorance threshold in red.**

271
 272 The results presented in Figure 5 demonstrate the impact of the number of neighbors and narrow kernels
 273 (characterized by high α values) on the evolution of the QS parameters. Specifically, it can be observed that
 274 interactions arise between these two factors, resulting in slightly erratic calibrated parameters. As the number of
 275 neighbors increases, the weights assigned to the furthest neighbors become negligible with larger α values. This
 276 means that these far away neighbors, despite being considered, have very little influence. This insensitivity only
 277 occurs for large n values, leading to minimal differences between possible configurations and noise in the metric.
 278 As expressed in the methodology section, in cases of a similar error, the cheapest solution is considered. In the
 279 case of QS, having a large number of neighbors can marginally increase the computational time, therefore, we
 280 introduce a small tolerance that results in favoring small n values. It is formulated as a small cost for each extra
 281 neighbor, i.e., by adding $5e-5 \times (\max(T) - \min(T))$ for each extra neighbor. However, the speed-up during
 282 simulation was limited to up to 10 %. Figure 6 shows similar quality (ε curves) as in Figure 5, but with the added

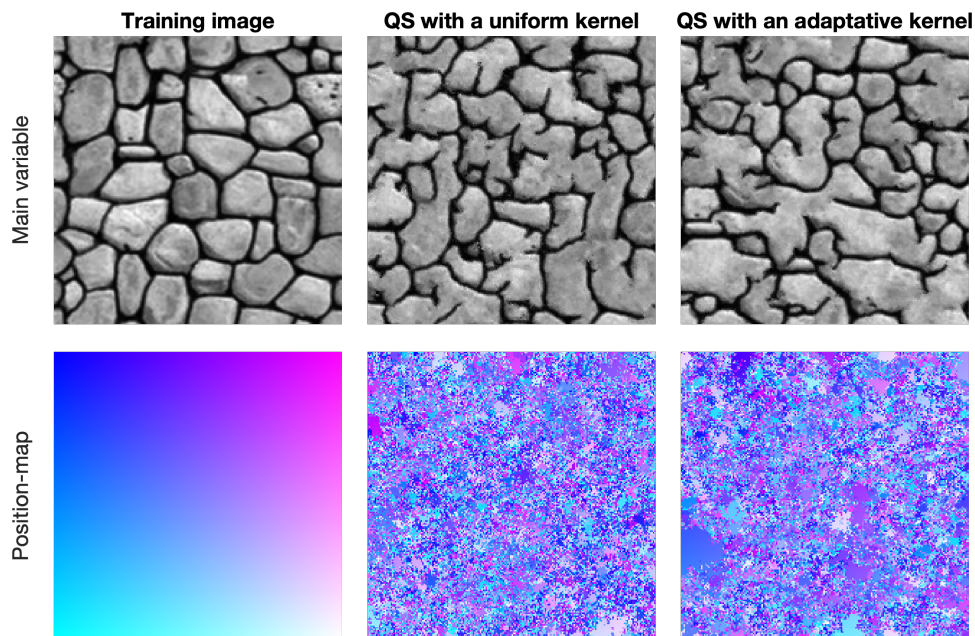
283 tolerance. As expected, the number of neighbors required during the simulation drastically decreases as advanced
 284 simulation stages, and the fluctuations in n are avoided.



285
 286 **Figure 6** Optimal parameters for QS (k in green and number of neighbors in blue, and the best kernel in magenta)
 287 as a function of the progression, with the associated prediction error (in black). The dashed blue line is the average
 288 density for the neighborhood considered. The ignorance threshold is in red.

289 5.3 Sequential simulation using automatic calibration

290 Figure 7 shows qualitative results using the evolutive parametrization resulting of the proposed autocalibration,
 291 using a case study that was published in Gravey and Mariethoz (2020). QS with an adaptive kernel refers to the
 292 use of different values of α for the kernel as a function of the simulation progression. In this case, the results are
 293 similar to state-of-the-art simulations using a manual calibration. Tests using QS with a uniform kernel fail to
 294 reproduce some structures, in particular the size of the objects is incorrect. Each position-map shows few
 295 homogenous areas; therefore, realizations are produced with a low rate of verbatim copy.



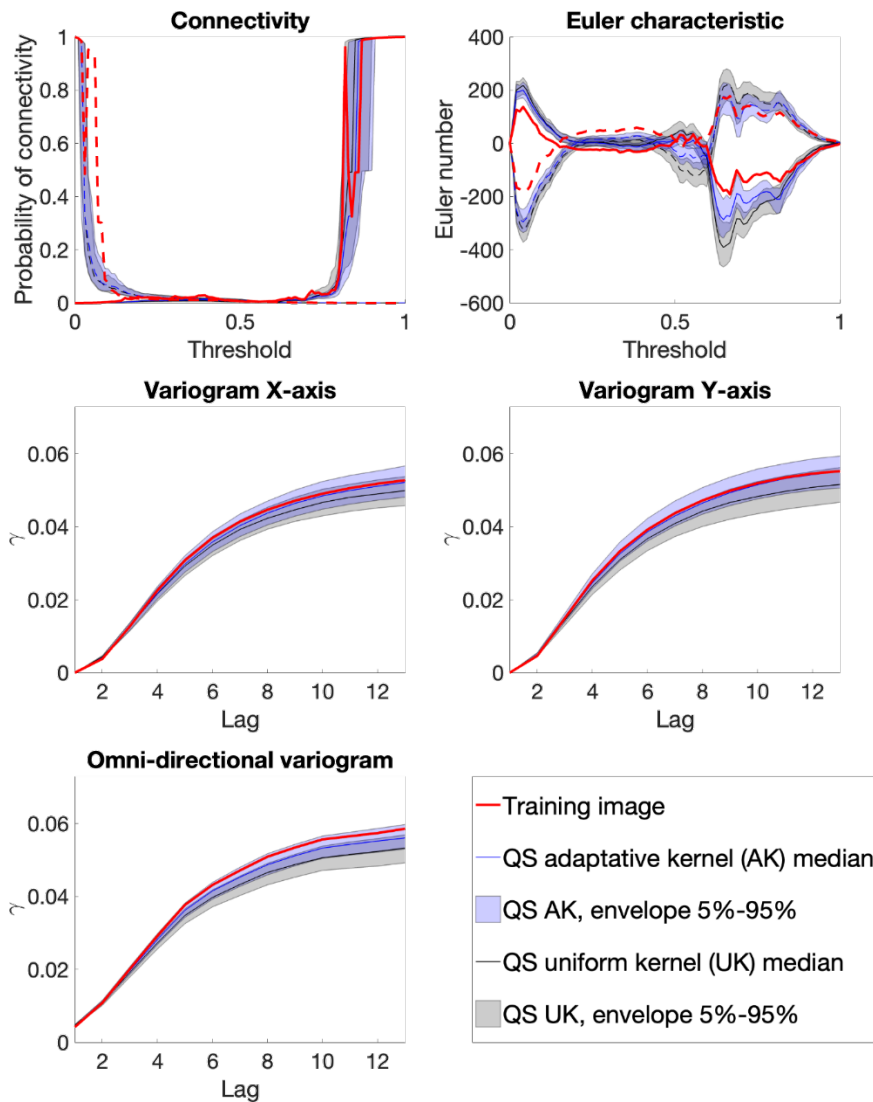
296
 297 **Figure 7** Simulation using QS with parameters generated by the automatic calibration.

298

299 From a quantitative evaluation, Figure 8 illustrates different metrics (variograms, connectivity as a structural
 300 indicator, and the Euler characteristic as noise indicator) (Renard and Allard, 2013) across a set of 100 realizations.
 301 The automatic calibration method proposed here allows obtaining better quality simulations than in Gravey and
 302 Mariethoz (2020).

303 Figure 9 shows that variogram and connectivity metrics are well reproduced, although they have not been directly
 304 constrained in the calibration process. Indeed, the parameter optimization only considers the simulation of single
 305 pixels and never computes global metrics over an entire grid.

306



307

308 **Figure 8 Benchmark between QS with an adaptive kernel (Figure 6) and a uniform (without) kernel (Figure 3) over 100**
 309 **simulations for 5 different metrics.**

310 6 Discussion and conclusion

311 The proposed method allows for the automatic calibration of QS and potentially similar pixel-based MPS
 312 approaches, reaching a similar or better quality as that of manual parameterization from both quantitative and
 313 qualitative points of view. Furthermore, it demonstrates that the optimal parametrization should not remain
 314 constant and instead needs to evolve with the simulation progression. The metrics confirm the good reproduction

315 of training patterns and the method finds a calibration that avoids verbatim copy. One major advantage of our
316 approach is the absence of a complex objective function, which often itself requires calibration.
317 A limitation of our approach is that it cannot be used to determine an optimal simulation path because it focuses
318 on the simulation of a single pixel. It also does not optimize the computational cost required for a simulation.
319 The computation time necessary to identify the appropriate parameters is contingent upon the expected quality.
320 However, the maximum time required for completion is predictable and depends on the number of patterns tested.
321 If required, the calibration can be further refined based on prior outcomes without restarting the entire process;
322 this can be achieved by adjusting D , incorporating additional kernels, or increasing $|\mathcal{V}|$. In certain instances,
323 adjusting the kernel parameter offers only minor improvements while necessitating a substantial number of
324 computations. Employing a more streamlined parameter space can yield comparable, and significantly reduce the
325 computational cost. This streamlined parameter space can be established, for instance, by subsampling the number
326 of neighbors according to a squared function (2,4,9,16,25,...) or by leveraging external or expert knowledge.
327 The proposed methodology was evaluated in multivariate scenarios, resulting in a more expansive parameter space
328 compared to single-variable cases. Although the approach yields satisfactory parameters, the inclusion of extra
329 parameters significantly extends the computation time, rendering the process impractical, particularly when
330 dealing with four or more variables.
331 In the context of testing the generality of our approach, calibration was computed on multiple training images
332 (found in the Supplementary material). The calibration pattern with two regimes (n large, then n small) seems to
333 be universal, at least for univariate simulations. While the position of the abrupt transition between regimes seems
334 to vary greatly (between 0.5% and 20% of the path), the overall shape remains the same. Therefore, the approach
335 proposed by Baninajar et al. (2019), in which long ranges are not considered, could be extended by using large n
336 values in the early stages of the simulation.
337 While show that it is possible to calibrate a parametric kernel, in future work one can envision the optimization of
338 a nonparametric kernel where the weight of each individual neighbor w_i is considered a variable to optimize using
339 ε as an objective function (e.g., using a machine learning regression framework).
340 The study of the evolution of parameters shows a smooth behavior of the average error. Therefore, the use of
341 multivariate fitting approaches to estimate the error surface with fewer evaluations could be an interesting solution
342 to speed up the parametrization. The use of machine learning to take advantage of transfer learning between
343 training images also has a high potential.
344

345 **Code availability**

346 The source code of the AutoQS algorithm is available as part of the G2S package at: [https://github.com/GAIA-](https://github.com/GAIA-UNIL/G2S)
347 [UNIL/G2S](https://github.com/GAIA-UNIL/G2S) (last access: 1st May 2023) under the GPLv3 license. And permanently available at
348 <https://doi.org/10.5281/zenodo.7792833>. Platform: Linux/macOS/Windows 10+. Language: C/C++. Interfacing
349 functions in MATLAB, Python3, and R.

350 **Author contributions.**

351 MG proposed the idea, implemented and optimized the autoQS approach and wrote the article. GM provided
352 supervision, methodological insights and contributed to the redaction.

353 **Competing interests**

354 The authors declare that they have no conflict of interest.

355 **Acknowledgements**

356 This research was funded by the Swiss National Science Foundation.

357 **6.1 Financial support**

358 This research has been supported by the Swiss National Science Foundation (grant no. 200021_162882).

359 **7 References**

- 360 Abdollahifard, M. J., Baharvand, M., and Mariéthoz, G.: Efficient training image selection for multiple-point
361 geostatistics via analysis of contours, *Computers & Geosciences*, 128, 41–50,
362 <https://doi.org/10.1016/j.cageo.2019.04.004>, 2019.
- 363 Baninajar, E., Sharghi, Y., and Mariéthoz, G.: MPS-APO: a rapid and automatic parameter optimizer for multiple-
364 point geostatistics, *Stoch Environ Res Risk Assess*, 33, 1969–1989, <https://doi.org/10.1007/s00477-019-01742-7>,
365 2019.
- 366 Boisvert, J. B., Pyrcz, M. J., and Deutsch, C. V.: Multiple Point Metrics to Assess Categorical Variable Models,
367 *Nat Resour Res*, 19, 165–175, <https://doi.org/10.1007/s11053-010-9120-2>, 2010.
- 368 Dagan, Y., Renard, P., Straubhaar, J., Erten, O., and Topal, E.: Automatic Parameter Tuning of Multiple-Point
369 Statistical Simulations for Lateritic Bauxite Deposits, *Minerals*, 8, 220, <https://doi.org/10.3390/min8050220>, 2018.
- 370 Gómez-Hernández, J. J. and Wen, X.-H.: To be or not to be multi-Gaussian? A reflection on stochastic
371 hydrogeology, *Advances in Water Resources*, 21, 47–61, [https://doi.org/10.1016/s0309-1708\(96\)00031-0](https://doi.org/10.1016/s0309-1708(96)00031-0), 1998.
- 372 Gravey, M. and Mariéthoz, G.: QuickSampling v1.0: a robust and simplified pixel-based multiple-point simulation
373 approach, *Geosci. Model Dev.*, 13, 2611–2630, <https://doi.org/10.5194/gmd-13-2611-2020>, 2020.
- 374 Guardiano, F. B. and Srivastava, R. M.: Multivariate Geostatistics: Beyond Bivariate Moments, in: *Quantitative*
375 *Geology and Geostatistics*, Springer Netherlands, 133–144, https://doi.org/10.1007/978-94-011-1739-5_12, 1993.
- 376 Journé, A. and Zhang, T.: The Necessity of a Multiple-Point Prior Model, *Math Geol*, 38, 591–610,
377 <https://doi.org/10.1007/s11004-006-9031-2>, 2006.
- 378 Mahmud, K., Mariéthoz, G., Caers, J., Tahmasebi, P., and Baker, A.: Simulation of Earth textures by conditional
379 image quilting, *Water Resour. Res.*, 50, 3088–3107, <https://doi.org/10.1002/2013wr015069>, 2014.
- 380 Mariéthoz, G., Caers, J., 2014. *Multiple-point geostatistics: stochastic modeling with training images*. Wiley.
- 381 Mariéthoz, G., Renard, P., and Straubhaar, J.: The Direct Sampling method to perform multiple-point geostatistical
382 simulations, *Water Resour. Res.*, 46, <https://doi.org/10.1029/2008wr007621>, 2010.
- 383 Matheron, G.: The intrinsic random functions and their applications, *Advances in Applied Probability*, 5, 439–
384 468, <https://doi.org/10.2307/1425829>, 1973.
- 385 Meerschman, E., Pirot, G., Mariéthoz, G., Straubhaar, J., Van Meirvenne, M., and Renard, P.: A practical guide
386 to performing multiple-point statistical simulations with the Direct Sampling algorithm, *Computers &*
387 *Geosciences*, 52, 307–324, <https://doi.org/10.1016/j.cageo.2012.09.019>, 2013.

388 Pérez, C., Mariethoz, G., and Ortiz, J. M.: Verifying the high-order consistency of training images with data for
389 multiple-point geostatistics, *Computers & Geosciences*, 70, 190–205,
390 <https://doi.org/10.1016/j.cageo.2014.06.001>, 2014.

391 Renard, P. and Allard, D.: Connectivity metrics for subsurface flow and transport, *Advances in Water Resources*,
392 51, 168–196, <https://doi.org/10.1016/j.advwatres.2011.12.001>, 2013.

393 Rezaee, H., Mariethoz, G., Koneshloo, M., and Asghari, O.: Multiple-point geostatistical simulation using the
394 bunch-pasting direct sampling method, *Computers & Geosciences*, 54, 293–308,
395 <https://doi.org/10.1016/j.cageo.2013.01.020>, 2013.

396 Straubhaar, J., Renard, P., Mariethoz, G., Froidevaux, R., and Besson, O.: An Improved Parallel Multiple-point
397 Algorithm Using a List Approach, *Math Geosci*, 43, 305–328, <https://doi.org/10.1007/s11004-011-9328-7>, 2011.

398 Strebel, S.: *Mathematical Geology*, 34, 1–21, <https://doi.org/10.1023/a:1014009426274>, 2002.

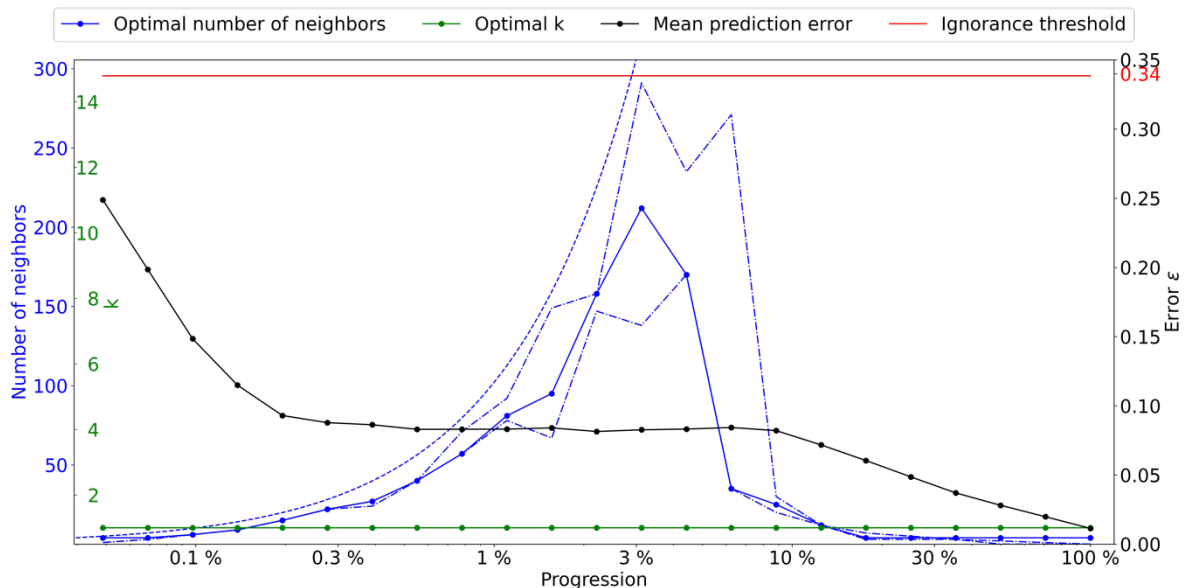
399 Tan, X., Tahmasebi, P., and Caers, J.: Comparing Training-Image Based Algorithms Using an Analysis of
400 Distance, *Math Geosci*, 46, 149–169, <https://doi.org/10.1007/s11004-013-9482-1>, 2013.

401
402
403

404 Appendix

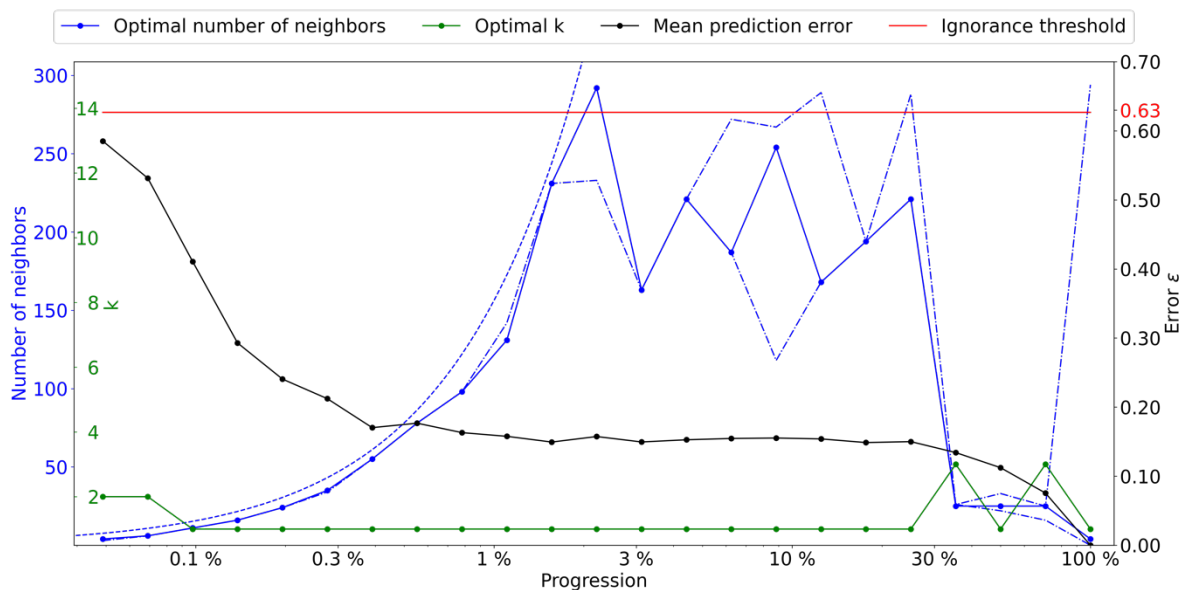
405 This supplementary material contains a similar calibration for other training images.
406

407 A. Stone

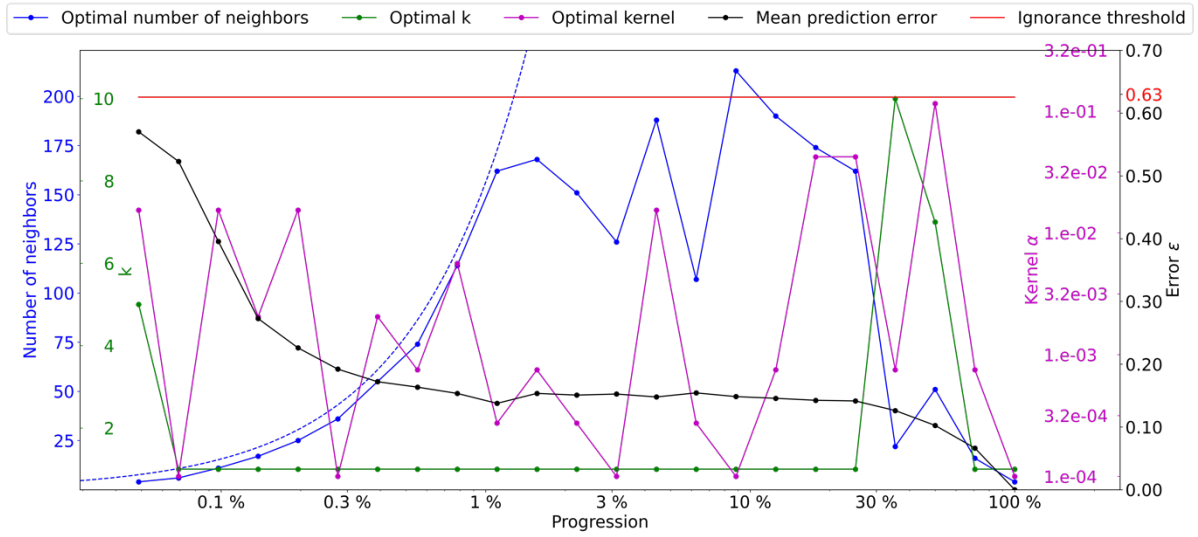


408
409 **Figure A.1** Optimal parameters for QS (k in green and number of neighbors in blue) as a function of the progression,
410 with the associated prediction error (in red). The red line represents the ignorance threshold. The dashed blue line is
411 the average density for the neighborhood considered. The dot-dashed line represents the variability in 1% of the error.

412 B. Strebelle (Strebelle, 2002)

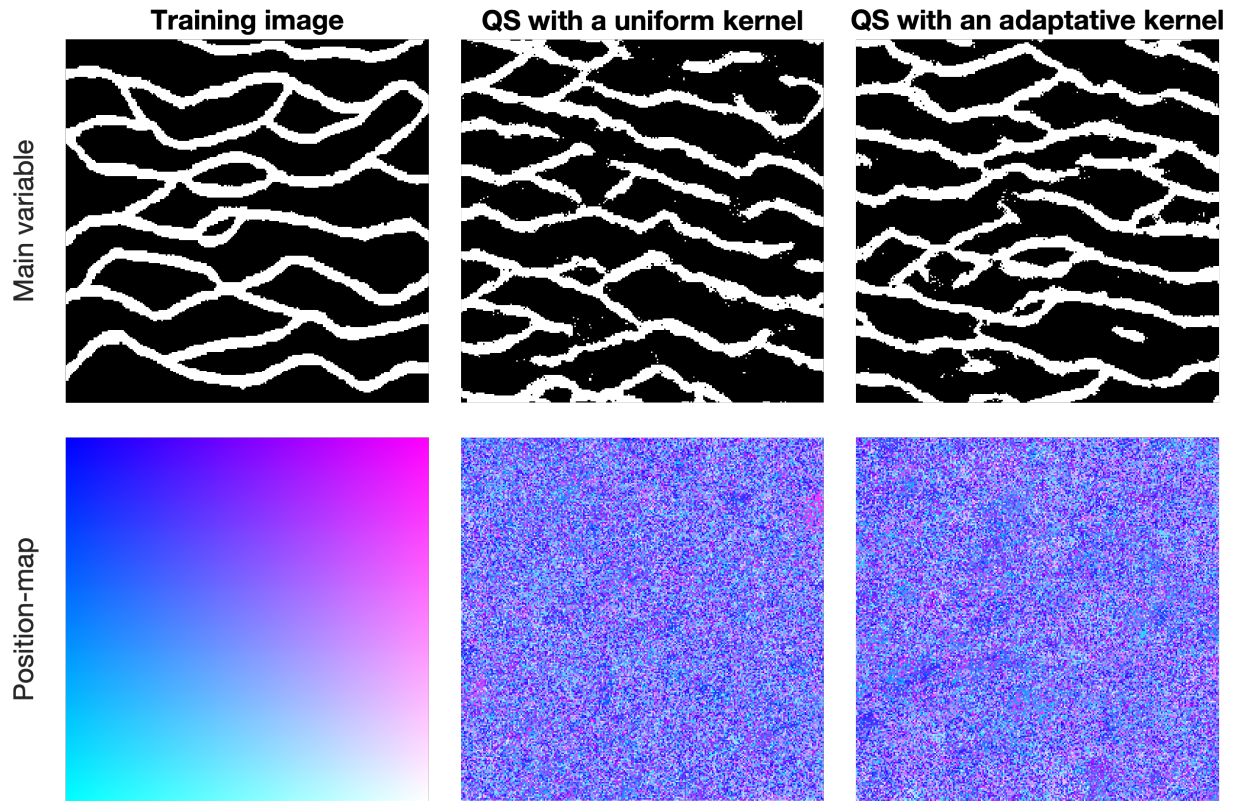


413
414 **Figure B.1** Optimal parameters for QS (k in green and number of neighbors in blue) as a function of the progression,
415 with the associated prediction error (in red). The red line represents the ignorance threshold. The dashed blue line is
416 the average density for the neighborhood considered. The dot-dashed line represents the variability in 1% of the error.



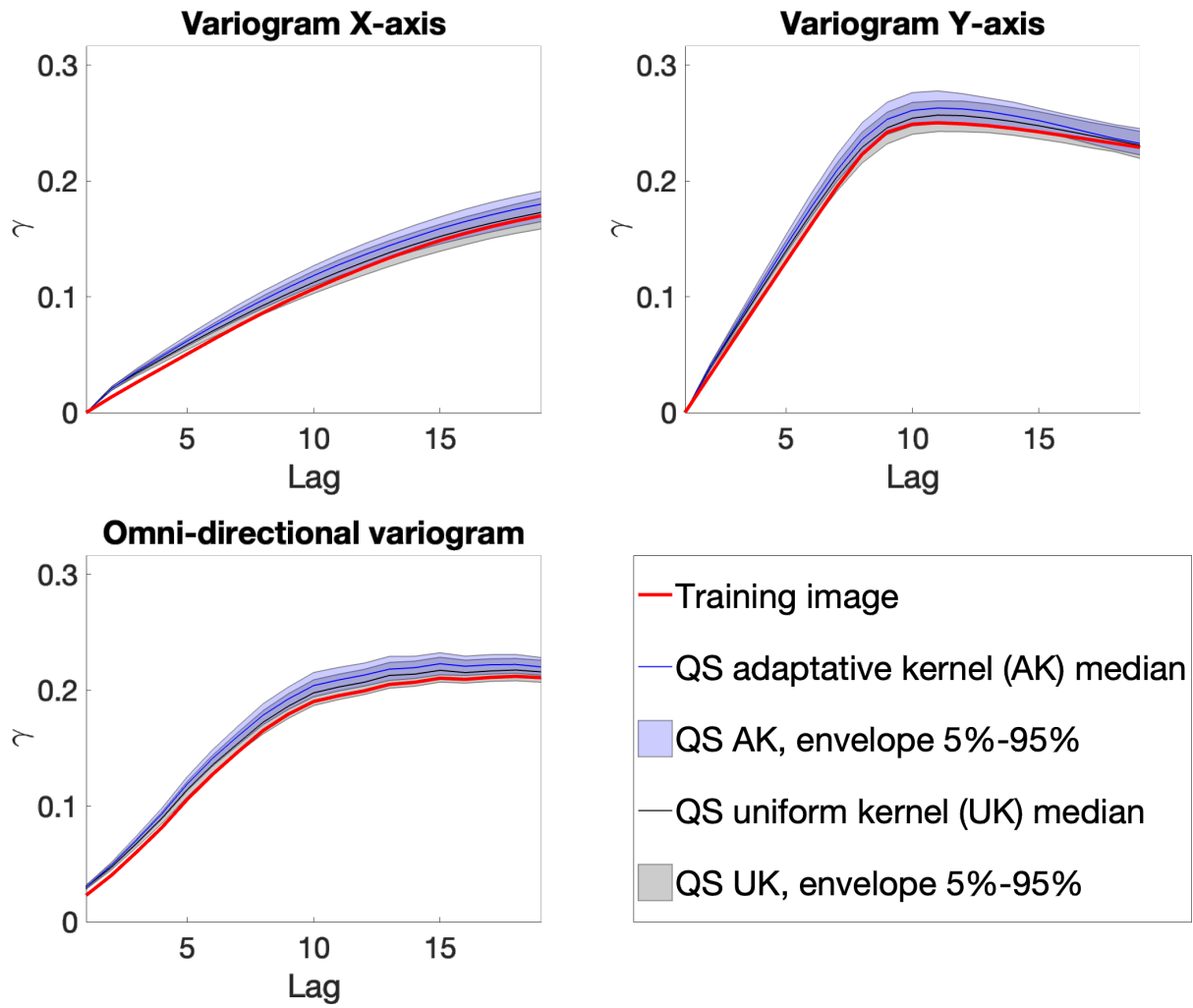
417
418
419
420

Figure B.2 Optimal parameters for QS (k in green and number of neighbors in blue, and the best kernel in magenta) as a function of the progression, with the associated prediction error (in red). The dashed blue line is the average density for the neighborhood considered.



421
422

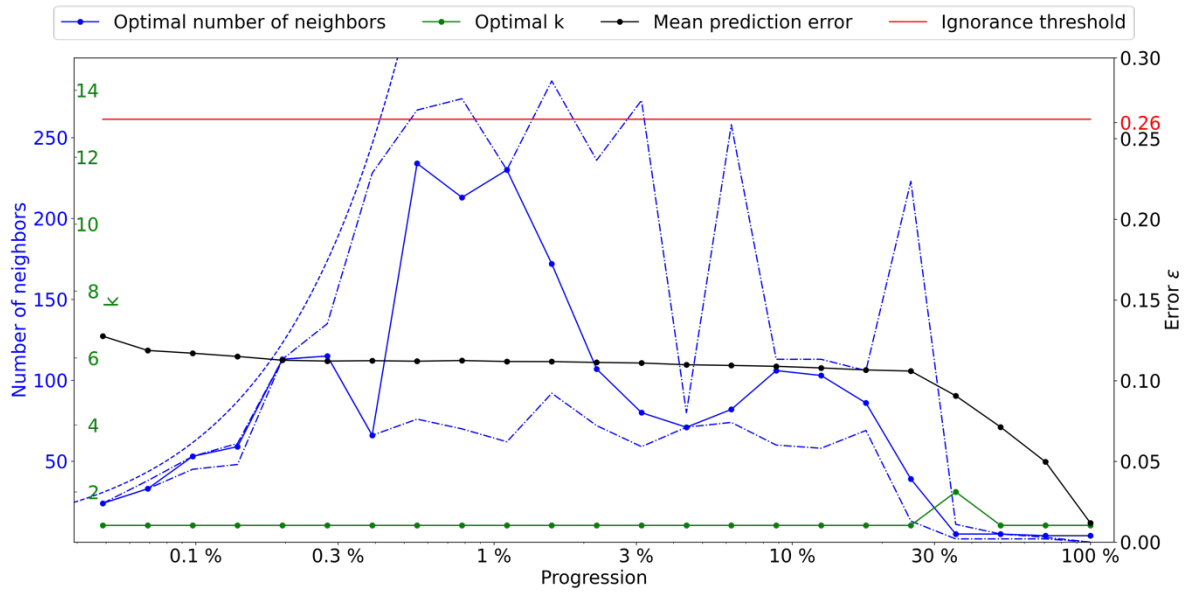
Figure B.3 Simulation using QS using parameters generated by the automatic calibration.



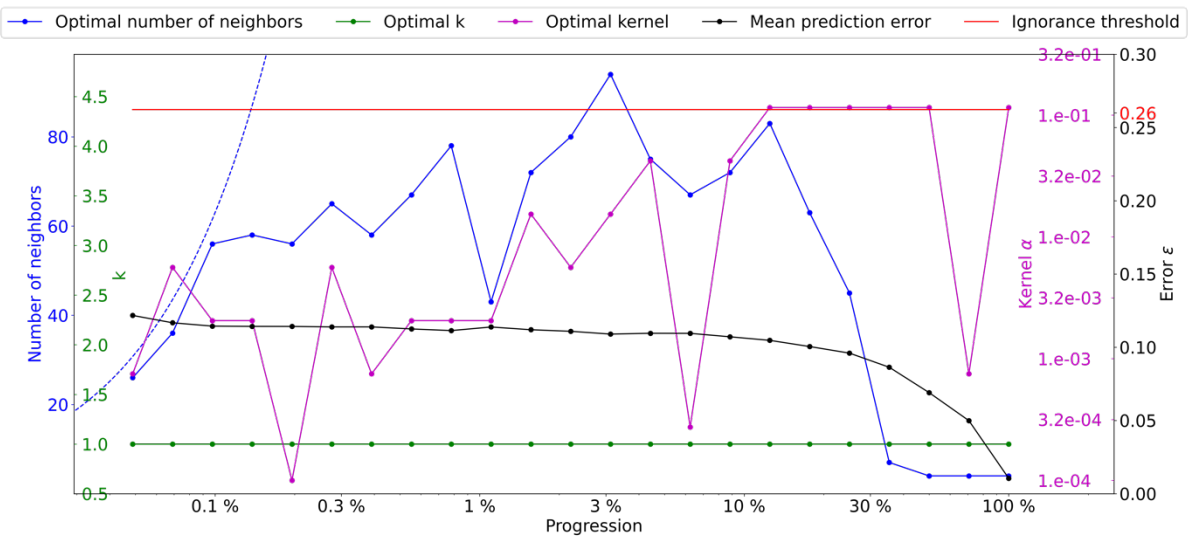
423

424 Figure B.4 Benchmark between QS with adaptative kernel (Figure B.2) and uniform (without) kernel (Figure B.1)
 425 over 100 simulations for 5 different metrics.

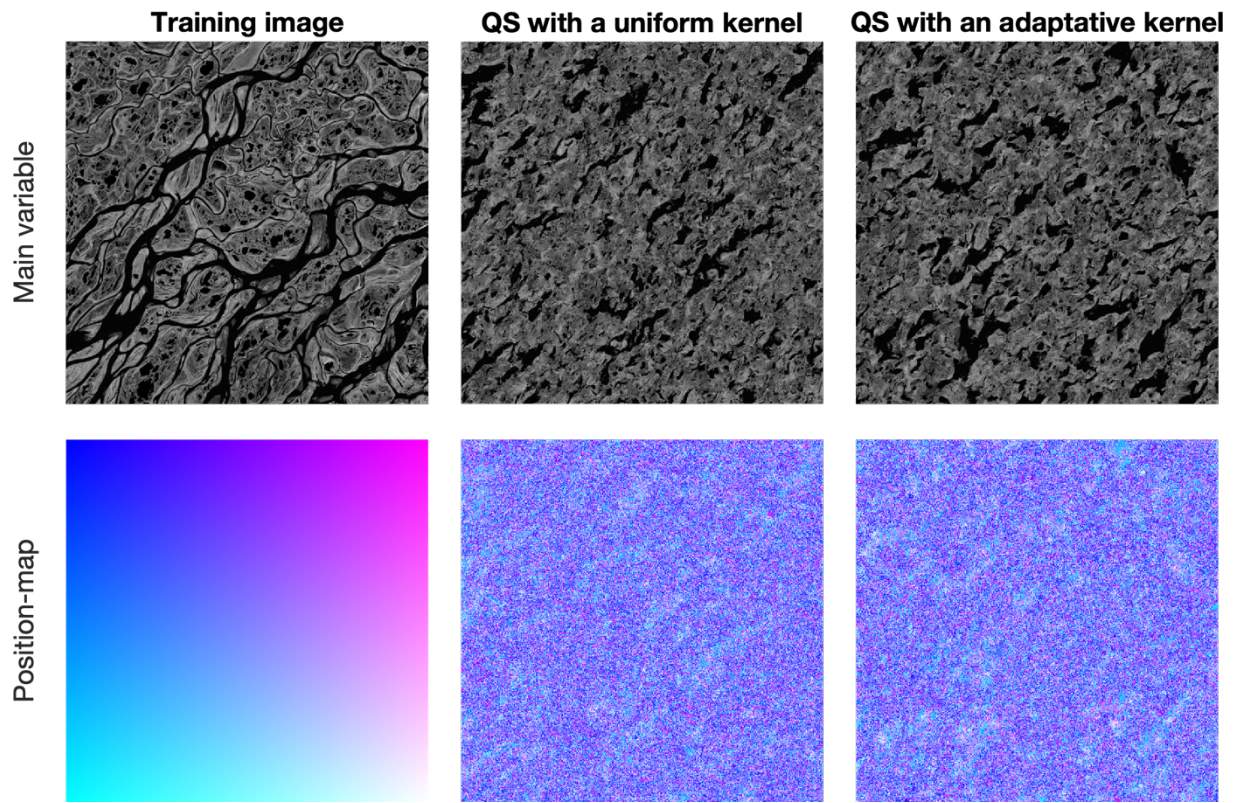
426



428
 429 **Figure C.1** Optimal parameters for QS (k in green and number of neighbors in blue) as a function of the progression,
 430 with the associated prediction error (in red). The red line represents the ignorance threshold. The dashed blue line is
 431 the average density for the neighborhood considered. The dot-dashed line represents the variability in 1% of the error.



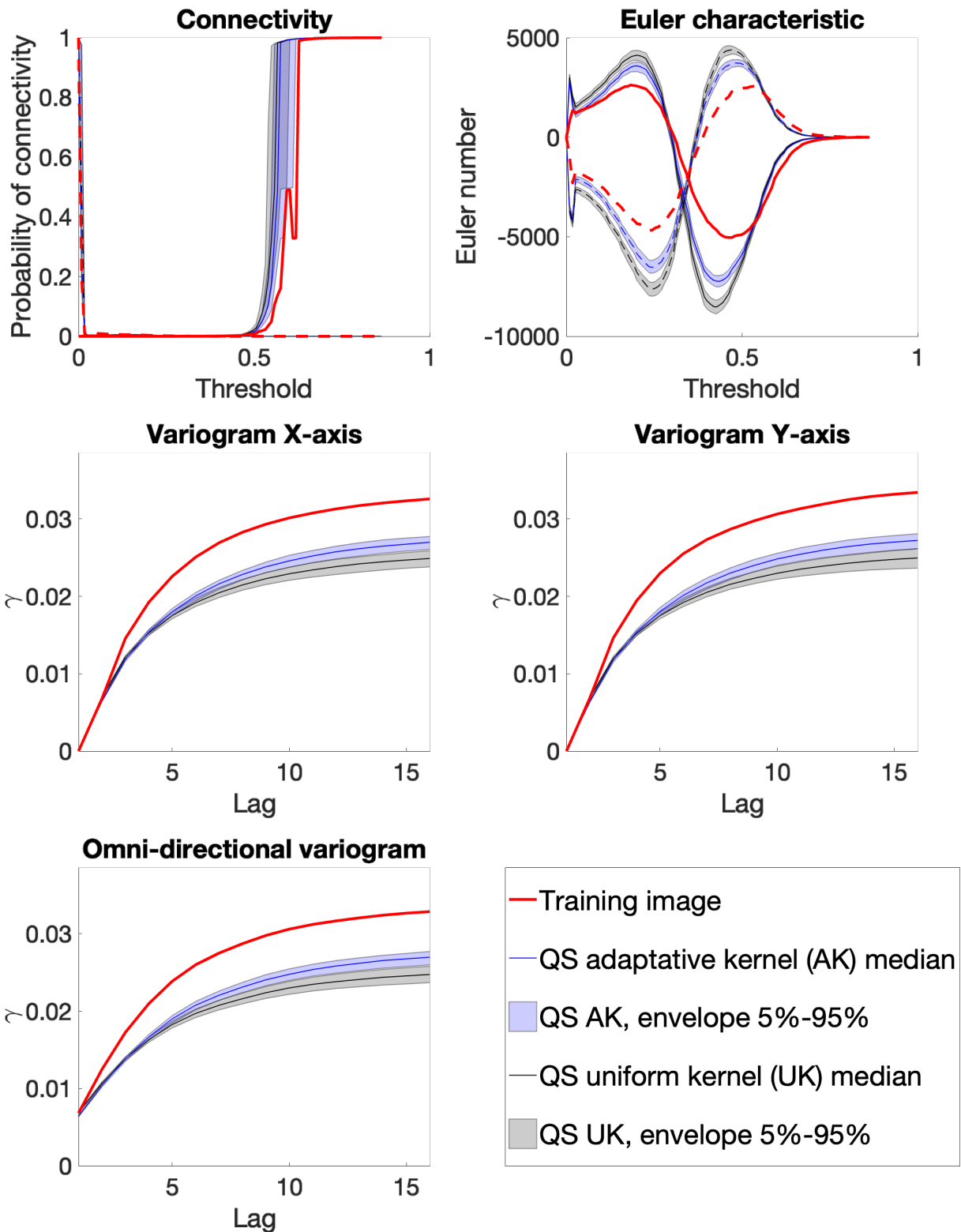
432
 433 **Figure C.2** Optimal parameters for QS (k in green and number of neighbors in blue, and the best kernel in magenta)
 434 as a function of the progression, with the associated prediction error (in red). The dashed blue line is the average density
 435 for the neighborhood considered.



436

437

Figure C.3 Simulation using QS using parameters generated by the automatic calibration.



438

439 **Figure C.4** Benchmark between QS with adaptative kernel (Figure C.2) and uniform (without) kernel (Figure C.1)
 440 over 100 simulations for 5 different metrics.

441