

# Performance Characterization in Computer Vision

Robert M. Haralick  
University of Washington  
Seattle WA 98195

## Abstract

Computer vision algorithms are composed of different sub-algorithms often applied in sequence. Determination of the performance of a total computer vision algorithm is possible if the performance of each of the sub-algorithm constituents is given. The problem, however, is that for most published algorithms, there is no performance characterization which has been established in the research literature. This is an awful state of affairs for the engineers whose job it is to design and build image analysis or machine vision systems.

This suggests that there has been a cultural deficiency in the computer vision community: computer vision algorithms have been published more on the merit of an experimental or theoretical demonstration suggesting that some task can be done, rather than on an engineering basis. Such a situation was tolerated because the interesting question was whether it was possible at all to accomplish a computer vision task. Performance was a secondary issue.

Now, however, a major interesting question is how to quickly design machine vision systems which work efficiently and which meet requirements. To do this requires an engineering basis which describes precisely what is the task to be done, how this task can be done, what is the error criterion, and what is the performance of the algorithm under various kinds of random degradations of the input data.

In this paper, we discuss the meaning of performance characterization in general, and then discuss the details of an experimental protocol under which an algorithm performance can be characterized.

## 1 Introduction

A major interesting question is how to quickly design machine vision systems which work efficiently and which meet requirements. To do this requires an engineering basis which describes precisely what is the task to be done, how this task can be done, what is the error criterion, and what is the performance of the algorithm under various kinds of random degradations of the input data. To accomplish this in the general case means propagating random perturbations through each algorithm stage in an open loop systems manner. To accomplish this for adaptive algorithms requires being able to do a closed loop engineering analysis. To perform a closed loop engineering analysis requires first doing an open loop engineering analysis and closing the loop by adding a constraint relation and solving for the output and the output random perturbation parameters.

The purpose of this discussion is to raise our sensitivity to these issues so that our field can more rapidly transfer the research technology to a factory floor technology. To initiate this dialogue, we will first expand on the meaning of performance characterization in general, and then discuss the experimental protocol under which an algorithm performance can be characterized.

## 2 Performance Characterization

What does performance characterization mean for an algorithm which might be used in a machine vision system? The algorithm is designed to accomplish a specific task. If the input data is perfect and has no noise and no random variation, the output produced by the algorithm ought also to be perfect. Otherwise, there is something wrong with the algorithm.

So measuring how well an algorithm does on perfect input data is not interesting. Performance characterization has to do with establishing the correspondence of the random variations and imperfections which the algorithm produces on the output data caused by the random variations and the imperfections on the input data. This means that to do performance characterization, we must first specify a model for the ideal world in which only perfect data exist. Then we must give a random perturbation model which specifies how the imperfect perturbed data arises from the perfect data. Finally, we need a criterion function which quantitatively measures the difference between the ideal output arising from the perfect ideal input and the calculated output arising from the corresponding randomly perturbed input.

Now we are faced with an immediate problem relative to the criterion function. It is typically the case that an algorithm changes the data unit. For example, an edge-linking process changes the data from the unit of pixel to the unit of a group of pixels. An arc segmentation/extraction process applied to the groups of pixels produced by an edge linking process produces fitted curve segments. This data unit change means that the representation used for the random variation of the output data set may have to be entirely different than the representation used for the random variation of the input data set. In our edge-linking/arc extraction example, the input data might be described by the false alarm/mis-detection characteristics produced by the preceding edge operation, as well as the standard deviation in the position and orientation of the correctly detected edge pixels. The random variation in the output data from the extraction process, on the other hand, must be described in terms of fitting errors (random variation in the fitted coefficients) and segmentation errors. Hence, the random perturbation model may change from stage to stage in the analysis process.

Consider the case for segmentation errors. The representation of the segmentation errors must be natural and suitable for the input of the next process in high-level vision which might be a model-matching process, for example. What should this representation be to make it possible to characterize the identification accuracy of the model matching as a function of the input segmentation errors and fitting errors? Questions like these, have typically not been addressed in the research literature. Until they are, analyzing the performance of a machine vision algorithm will be in the dark ages of an expensive experimental trial-and-error process. And if the performance of the different

pieces of a total algorithm cannot be used to determine the performance of the total algorithm, then there cannot be an engineering design methodology for machine vision systems.

This problem is complicated by the fact that there are many instances of algorithms which compute the same sort of information but in forms which are actually non-equivalent. For example, there are arc extraction algorithms which operate directly on the original image along with an intermediate vector file obtained in a previous step and which output fitted curve segments. There are other arc extraction algorithms which operate on groups of pixels and which output arc parameters such as center, radius, and endpoints in addition to the width of the original arc.

What we need is the machine vision analog of a system's engineering methodology. This methodology can be encapsulated in a protocol which has a modeling component, an experimental component, and a data analysis component. The next section describes in greater detail these components of an image analysis engineering protocol.

### 3 Protocol

The modeling component of the protocol consists of a description of the world of ideal images, a description of a random perturbation model by which non-ideal images arise, a description of a random perturbation process which characterizes the output random perturbation as a function of the parameters of the input random perturbation and a specification of the criterion function by which the difference between the ideal output and the computed output arising from the imperfect input can be quantified. The experimental component describes the experiments performed under which the data relative to the performance characterization can be gathered. The analysis component describes what analysis must be done on the experimentally observed data to determine the performance characterization.

#### 3.1 Input Image Population

This part of the protocol describes how, in accordance with the specified model, a suitably random, independent, and representative set of images from the population of ideals is to be acquired or generated to constitute the sampled set of images. This acquisition can be done by taking real images under the specified conditions or by generating synthetic images. If the population includes, for example, a range of sizes of the object of interest or if the object of interest can appear in a variety of situations, or if the object shape can have a range of variations, then the sampling mechanism must assure that a reasonable number of images are sampled with the object appearing in sizes, orientations, and shape variations throughout its permissible range. Similarly, if the object to be recognized or measured can appear in a variety of different lighting conditions which create a similar variety in shadowing, then the sampling must assure that images are acquired with the lighting and shadowing varying throughout its permissible range.

Some of the variables used in the image generation process are ones whose values will be estimated by the computer vision algorithm. We denote these

variables by  $z_1, \dots, z_K$ . Other of these variables are nuisance variables. Their values provide for variation. The performance characterization is averaged over their values. We denote these variables by  $w_1, \dots, w_M$ . Other of the variables specify the parameters of the random perturbation and noise process against which the performance is to be characterized. We denote these variables by  $y_1, \dots, y_J$ . The generation of the images in the population can then be described by  $N = J + K + M$  variables. If these  $N$  variables having to do with the kind of lighting, light position, object position, object orientation, permissible object shape variations, undesired object occlusion, environmental clutter, distortion, noise etc., have respective range sets  $R_1, \dots, R_N$ , then the sampling design must assure that images are selected from the domain  $R_1 \times R_2 \times \dots \times R_N$  in a representative way. Since the number of images sampled is likely to be a relatively small fraction of the number of possibilities in  $R_1 \times R_2 \times \dots \times R_N$ , the experimental design may have to make judicious use of a Latin square layout.

### 3.2 Random Perturbation and Noise

Specification of random perturbation and noise is not easy because the more complex the data unit, the more complex the specification of the random perturbation and noise. Each specification of randomness has two potential components. One component is a small perturbation component which affects all data units. It is often reasonable to model this by an additive Gaussian noise process on the ideal values of the data units. This can be considered to be the small variation of the ideal data values combined with observation or measurement noise. The other component is a large perturbation component which affects only a small fraction of the data units. For simple data units it is reasonable to model this by replacing its value by a value having nothing to do with its true value. Large perturbation noise on more complex data units can be modeled by fractionating the unit into pieces and giving values to most of the pieces which would follow from the values the parent data unit had and giving values to the remaining pieces which have nothing to do with the values the original data unit had.

This kind of large random perturbation affecting a small fraction of units is replacement noise. It can be considered to be due to random occlusion, linking, grouping, or segmenting errors. Algorithms which work near perfectly on small amounts of random perturbation on all data units, often fall apart with large random perturbation on a small fraction of the data units. Much of the performance characterization of a complete algorithm will be specified in terms of how much of this replacement kind of random perturbation the algorithm can tolerate and still give reasonable results. Algorithms which have good performance even with large random perturbation on a small fraction of data units can be said to be robust.

### 3.3 Performance Characterization

Some of the variables used in the image generation are those whose values are to be estimated by the machine vision algorithm. Object kind, location, and orientation are prime examples. The values of such variables do not make the recognition and estimation much easier or harder, although they may have some

minor effect. For example, an estimate of the surface normal of a planar object viewed at a high slant angle will tend to have higher variance than an estimate produced by the planar object viewed at a near normal angle. The performance characterization of an image analysis algorithm is not with respect to this set of variables. From the point of view of what is to be calculated, this set of variables is crucial. From the point of view of performance characterization, the values for the variables in this set as well as the values in the nuisance set are the ones over which the performance is averaged.

Another set of variables characterize the extent of random perturbations which distort the ideal input data to produce the imperfect input data. These variables represent variations which degrade the information in the image, thereby increasing the uncertainty of the estimates produced by the algorithm. Such variables may characterize object contrast, noise, extent of occlusion, complexity of background clutter, and a multitude of other factors which instead of being modeled explicitly are modeled implicitly by the inclusion of random shape perturbations applied to the set of ideal model shapes.

Finally, there may be other variables governing parameter constants that must be set in the image analysis algorithm. The values of these variables may to a large or small extent change the performance of the algorithm.

The variables characterizing the input random perturbation process and the variables which are the algorithm tuning constants constitute the set of variables in terms of which the performance characterization must be measured. Suppose there are  $I$  algorithm parameters  $x_1, \dots, x_I$ , which can be set,  $J$  different variables  $y_1, \dots, y_J$  characterizing the random perturbation process, and  $K$  different measurements  $\hat{z}_1, \dots, \hat{z}_K$  to be made on each image. There will be a difference between the true ideal values  $z_1, \dots, z_K$  of the measured quantities and the measured values  $\hat{z}_1, \dots, \hat{z}_K$  themselves. The nature of this difference can be characterized by the parameters  $q_1, \dots, q_L$  of the output random perturbation process:  $(q_1, \dots, q_L) = f(x_1, \dots, x_I, y_1, \dots, y_J, z_1, \dots, z_K)$ .

The last step of a total algorithm not only has a characterization of the output random perturbation parameters, but also an error criterion  $e$  which is application and domain specific. The error criterion,  $e(z_1, \dots, z_K, \hat{z}_1, \dots, \hat{z}_K)$ , must state how the comparison between the ideal values and the measured values will be evaluated. Its value will be a function of the  $I$  algorithm parameters and the  $J$  random perturbation parameters.

An algorithm can have two different dimensions to the error criterion. To explain these dimensions, consider algorithms which estimate some parameter such as position and orientation of an object. One dimension the error criterion can have is reliability. An estimate can be said to be reliable if the algorithm is operating on data that meets certain requirements and if the difference between the estimated quantity and the true but known value is below a user specified tolerance. An algorithm can estimate whether the results it produces are reliable by making a decision on estimated quantities which relate to input data noise variance, output data covariance, and structural stability of calculation. Output quantity covariance can be estimated by estimating the input data noise variance and propagating the error introduced by the noise variance into the calculation of the estimated quantity. Hence the algorithm itself can provide an indication of whether the estimates it produces have an uncertainty below a given value. High uncertainties would occur if the algorithm can determine that the assumptions about the environment producing the data or the

assumptions required by the method are not being met by the data on which it is operating or if the random perturbation in the quantities estimated is too high to make the estimates useful.

Characterizing this dimension can be done by two means. The first is by the probability that the algorithm claims reliability as a function of algorithm parameters and parameters describing input data random perturbations. The second is by misdetection false alarm operating curves. A misdetection occurs when the algorithm indicates it has produced a reliable enough result when in fact it has not produced a reliable enough result. A false alarm occurs when the algorithm indicates that it has not produced a reliable enough result when in fact it has produced a reliable enough result. A misdetection false alarm rate operating curve results for each different noise and random perturbation specification. The curve itself can be obtained by varying the algorithm tuning constants, one of which is the threshold by which the algorithm determines whether it claims the estimate it produces is reliable or not.

The second dimension of the error criterion would be related to the difference between the true value of the quantity of interest and the estimated value. This criterion would be evaluated only for those cases where the algorithm indicates that it produces a reliable enough result. A scalar error criterion would weight both of these dimensions in an appropriate manner.

Each estimated quantity  $\hat{z}_k$  is a random variable which is a function of the ideal input data, the values of the algorithm tuning parameters  $x_1, \dots, x_I$  and the random perturbation parameters  $y_1, \dots, y_J$  characterizing the random perturbation process distorting the ideal input.

Each ideal quantity  $z_k$  is a function only of the algorithm constants  $x_1, \dots, x_I$ . The expected value  $E$  of  $e(z_1, \dots, z_K, \hat{z}_1, \dots, \hat{z}_K)$  is taken over the input data set subpopulation consistent with  $z_1, \dots, z_K$  and the random perturbation process. It is, therefore, a function of  $x_1, \dots, x_I$  and  $y_1, \dots, y_J$ . Performance characterization of the estimated quantity with respect to the error criterion function then amounts to expressing in graph, table or analytic form  $E[e(z_1, \dots, z_K, \hat{z}_1, \dots, \hat{z}_K)]$  for each  $z_1, \dots, z_K$  as a function of  $x_1, \dots, x_I$  and  $y_1, \dots, y_J$ .

### 3.4 Experiments

In a complete design, the values for the algorithm constants  $x_1, \dots, x_I$  and the values for the random perturbation parameters  $y_1, \dots, y_J$  will be selected in a systematic and regular way. The values for  $z_1, \dots, z_K$  and the values for the nuisance variables  $w_1, \dots, w_M$  will be sampled from a uniform distribution over the range of their permissible values.

The values for  $z_1, \dots, z_K$  specify the equivalence class of ideal images. The values for  $y_1, \dots, y_J$  characterize the random perturbations and noise which are randomly introduced into the ideal image and/or object(s) in the ideal image. In this manner, each noisy trial image is generated. The values for  $x_1, \dots, x_I$  specify how to set the tuning constants required by the algorithm. The algorithm is then run over the trial image producing estimated values  $\hat{z}_1, \dots, \hat{z}_K$  for  $z_1, \dots, z_K$ .

The data analysis plan for the characterization of the output random per-

turbation process generates records of the form

$$\mathbf{x}_1, \dots, \mathbf{x}_I, \mathbf{y}_1, \dots, \mathbf{y}_J, z_1, \dots, z_K, \hat{z}_1, \dots, \hat{z}_K$$

From an assumed model of the output random perturbation process, the data analysis plan will have a way of estimating the parameters  $q_1, \dots, q_L$  of the output random perturbation process from the  $z_1, \dots, z_K, \hat{z}_1, \dots, \hat{z}_K$  part of the records. Thus  $q_1, \dots, q_L$  will be a function of  $\mathbf{x}_1, \dots, \mathbf{x}_I, \mathbf{y}_1, \dots, \mathbf{y}_J, z_1, \dots, z_K$ . The data analysis plan must also specify how this dependence will be determined by an estimating or fitting procedure.

If we apply the error criterion to each record, we then produce the values  $e(z_1, \dots, z_K, \hat{z}_1, \dots, \hat{z}_K)$ . The data produced by each trial then consists of a record

$$\mathbf{x}_1, \dots, \mathbf{x}_I, \mathbf{y}_1, \dots, \mathbf{y}_J, e(z_1, \dots, z_K, \hat{z}_1, \dots, \hat{z}_K)$$

The data analysis plan for the error criterion describes how the set of records produced by the experimental trials will be processed or analyzed to compactly express the performance characterization. For example, an equivalence relation on the range space for  $y_1, \dots, y_J$  may be defined and an hypothesis may be specified stating that all combinations of values of  $y_1, \dots, y_J$  in the same equivalence class have the same expected error. The data analysis plan would specify the equivalence relation and give the statistical procedure by which the hypothesis could be tested. Performing such tests are important because they can reduce the number of variable combinations which have to be used to express the performance characterization. For example, the hypothesis that all other variables being equal, whenever  $y_{J-1}/y_J$  has a ratio of  $k$ , then the expected performance is identical. In this case, the performance characterization can be compactly given in terms of  $k$  and  $y_1, \dots, y_{J-2}$ .

Once all equivalence tests are complete, the data analysis plan would specify the kinds of graphs or tables employed to present the experimental data. It might specify the form of a simple regression equation by which the expected error, the probability of claimed reliability, the probability of misdetection, the probability of false alarm, and the computational complexity or execution time can be expressed in terms of the independent variables  $\mathbf{x}_1, \dots, \mathbf{x}_I, \mathbf{y}_1, \dots, \mathbf{y}_J$ . As well it would specify how the coefficients of the regression equation could be calculated from the observed data. Finally, when error propagation can be done analytically using the parameters associated with input data noise variance and the ideal noiseless input data, the data analysis plan can discuss how to make the comparison between the expected error computed analytically and the observed experimental error.

Finally, if the computer vision algorithm must meet certain performance requirements, the data analysis plan must state how the hypothesis that the algorithm meets the specified requirement will be tested. The plan must be supported by a theoretically developed statistical analysis which shows that an experiment carried out according to the experimental design and analyzed according to the data analysis plan will produce a statistical test itself having a given accuracy. That is, since the entire population of images is only sampled, the sampling variation will introduce a random fluctuation in the test results. For some fraction of experiments carried out according to the protocol, the hypothesis to be tested will be accepted but the algorithm, in fact, if it were tried on the complete population of image variations, would not meet the specified

requirements; and for some fraction of experiments carried out according to the protocol, the hypothesis to be tested will be rejected but if the algorithm were tried on the complete population of image variation, it would meet the specified requirements. The specified size of these errors of false acceptance and missed acceptance will dictate the number of images to be in the sample for the test. This relation between sample size and false acceptance rate and missed acceptance rate of the test for the hypothesis must be determined on the basis of statistical theory. One would certainly expect that the sample size would be large enough so that the uncertainty caused by the sampling would be below 20%.

For example, suppose the error rate of a quantity estimated by a machine vision algorithm is defined to be the fraction of time that the estimate is further than  $\epsilon_0$  from the true value. If this error rate is to be less than  $\frac{1}{1,000}$ , then in order to be about 85% sure that the performance meets specification, 10,000 tests will have to be run. If the image analysis algorithm performs incorrectly 9 or fewer times, then we can assert that with 85% probability, the machine vision algorithm meets specification [1].

## 4 Conclusion

We have discussed the problem of the lack of performance evaluation in the published literature on computer vision algorithms. This situation is causing great difficulties to researchers who are trying to build up on existing algorithms and to engineers who are designing operational systems. To remedy the situation, we suggested the establishment of a well-defined protocol for determining the performance characterization of an algorithm. Use of this kind of protocol will make using engineering system methodology possible as well as making possible well-founded comparisons between machine vision algorithms that perform the same tasks. We hope that our discussion will encourage a thorough and overdue dialogue in the field so that a complete engineering methodology for performance evaluation of machine vision algorithms can finally result.

## References

- [1] Haralick, R.M., "Performance Assessment of Near Perfect Machines," *Machine Vision and Applications*, Vol. 2, No. 1, 1989, pp. 1-16.