Color Constancy Algorithms: Psychophysical Evaluation on a New Dataset

Javier Vazquez, C. Alejandro Párraga, Maria Vanrell[^] and Ramon Baldrich

Department of Computer Science, Centre de Visió per Computador, Universitat Autònoma de Barcelona, Edifíci O, Campus UAB (Bellaterra), C.P. 08193 Barcelona, Spain E-mail: javier.vazquez@cvc.uab.es

I man. Javien valgaez eve. aub.es

AQ: #1 3

4

5

6

7 8 Abstract. The estimation of the illuminant of a scene from a digital 9 image has been the goal of a large amount of research in computer 10 vision. Color constancy algorithms have dealt with this problem by 11 defining different heuristics to select a unique solution from within 12 the feasible set. The performance of these algorithms has shown 13 that there is still a long way to go to globally solve this problem as a 14 preliminary step in computer vision. In general, performance evalu-15 ation has been done by comparing the angular error between the 16 estimated chromaticity and the chromaticity of a canonical il-17 luminant, which is highly dependent on the image dataset. Recently, 18 some workers have used high-level constraints to estimate il-19 luminants; in this case selection is based on increasing the perfor-20 mance on the subsequent steps of the systems. In this paper the 21 authors propose a new performance measure, the perceptual angu-22 lar error. It evaluates the performance of a color constancy algorithm 23 according to the perceptual preferences of humans, or naturalness 24 (instead of the actual optimal solution) and is independent of the 25 visual task. We show the results of a new psychophysical experi-26 ment comparing solutions from three different color constancy algo-27 rithms. Our results show that in more than half of the judgments the 28 preferred solution is not the one closest to the optimal solution. Our 29 experiments were performed on a new dataset of images acquired 30 with a calibrated camera with an attached neutral gray sphere, 31 which better copes with the illuminant variations of the scene. 32 © 2009 Society for Imaging Science and Technology. 33 [DOI: XXXX] 34

35 INTRODUCTION

 Color constancy is the ability of the human visual system to perceive a stable representation of color despite illumination changes. Like other perceptual constancy capabilities of the visual system, color constancy is crucial for succeeding in many ecologically relevant visual tasks such as food collec- tion, detection of predators, etc. The importance of color constancy in biological vision is mirrored in computer vision applications, where success in a wide range of visual tasks relies on achieving a high degree of illuminant invariance. In the last 20 years, research in computational color constancy has tried to recover the illuminant of a scene from an ac-quired image.

48 This has been shown to be a mathematically ill-posed 49 problem, which therefore does not have a unique solution. A

1062-3701/2009/53(3)/1/0/\$20.00.

common computational approach to illuminant recovery ⁵⁰ (and color constancy in general) is to produce a list of pos- 51 sible illuminants (feasible solutions) and then use some as- 52 sumptions, based on the interactions of scene surfaces and 53 illuminants to select the most appropriate solution among 54 all possible illuminants. A recent extended review of compu- 55 tational color constancy methods was provided by Hordley.¹ 56 In this review, computational algorithms were classified in 57 five different groups according to how they approach the 58 problem. These were (a) simple statistical methods, 2 (b) 59 neural networks,³ (c) gamut mapping,^{4,5} (d) probabilistic 60methods,⁶ and (e) physics-based methods.⁷ Comparison 61 studies^{8,9} have ranked the performance of these algorithms, 62 which usually depend on the properties of the image dataset 63 and the statistical measures used for the evaluation. It is 64 generally agreed that, although some algorithms may per- 65 form well on average, they may also perform poorly for spe- 66 cific images. This is the reason why some authors¹⁰ have 67 proposed a one-to-one evaluation of the algorithms on in- 68 dividual images. In this way, comparisons become more in- 69 dependent of the chosen image dataset. However, the general 70 conclusion is that more research should be directed toward a 71 combination of different methods, since the performance of 72 a method usually depends on the type of scene with which it 73 deals.¹¹ Recently, some interesting studies have pointed to- 74 ward this direction,¹² i.e., trying to find which statistical 75 properties of the scenes determine the best color constancy 76 method to use. In all these approaches, the evaluation of the 77 performance of the algorithms has been based on comput- 78 ing the angular error between the selected solution and the 79 actual solution that is provided by the acquisition method. 80

Other recent proposals^{13,14} turn away from the usual **81** approach and deal instead with multiple solutions delegating **82** the selection of a unique solution to a subsequent step that **83** depends on high-level, task-related interpretations, such as **84** the ability to annotate the image content. In this example, **85** the best solution would be the one giving the best semantic **86** annotation of the image content. It is in this kind of ap-**87** proach where the need for a different evaluation emerges, **88** since the performance depends on the visual task and this **89** can lead to an inability to compare different methods. **90** Hence, to be able to evaluate this performance and to com-**91** pare it with other high-level methods, we propose to explore **92** a new evaluation procedure.

[▲]IS&T Member.

Received Aug. 25, 2008; accepted for publication Dec. xx, xxxx; published online Dec. xx, xxxx.



Figure 1. Images regularly selected in the experiment as natural (left) vs images hardly ever selected (right).

94 In summary, the goal of this paper is to show the results 95 of a new psychophysical experiment following the lines of 96 that presented by Vazquez et al.¹⁵ The previous results were 97 confirmed, that is, humans do not choose the minimum 98 angular error solution as the more natural one. Further-99 more, in this paper we propose a new measure to reduce the 100 gap between the error measure and the human preference. 101 Our new experiment represents an improvement over the 102 old one in that it considers the uncertainty level of the ob-103 server responses and it uses a new, improved image dataset. 104 This new dataset has been built by using a neutral gray 105 sphere attached to the calibrated camera to better estimate 106 the illuminant of the scene. We have worked with the 107 Shades-of-Gray¹⁶ algorithm instead of CRule.¹⁷ This deci-108 sion was made on the basis that CRule is calibrated whereas 109 the other algorithms are not.

110 EXPERIMENTAL SETUP

AQ:

#2

AQ:

#3

111 Subjects were presented with a pair of images (each one a 112 different color constancy solution) on a CRT monitor and 113 asked to select the image that seems "most natural." The term "natural" was chosen not because it refers to natural ¹¹⁴ objects but because it refers to natural viewing conditions, ¹¹⁵ implying the least amount of digital manipulation or global ¹¹⁶ perception of an illuminant. Figure 1 shows some exemplary ¹¹⁷ pictures from the database. The pictures on the left are ex- ¹¹⁸ amples of images selected as natural most of the time, while ¹¹⁹ those on the right are examples of images hardly ever se- ¹²⁰ lected as natural. ¹²¹

The global schematics of the experiment are shown in 122 Figure 2. We used a set of 83 images from a new image 123 dataset that was built for this experiment (the image gather-124 ing details are explained below). The camera calibration al-125 lows us to obtain the Commission Internationale de 126 l'Eclairage (CIE) 1931 XYZ values for each pixel and conse-127 quently, we converted 83 images from CIE XYZ space to CIE 128 standard red, green blue (sRGB). Following this, we replaced 129 the original illuminant by D65 using the chromaticity values 130 of the gray sphere that was present in all image scenes. 131

From the original images, five new pictures were created 132 by reilluminating the scene with five different illuminants. 133 To this end we have used the chromatic values of each il- 134 luminant (three Plankians: 4000, 7000, and 10,000 K, and 135 two arbitrary illuminants: greenish (x=0.3026, y=0.3547) 136 and purplish (x=0.2724, y=0.2458), totaling 415 images. 137 Afterward, the three color constancy algorithms 138 (Gray-World,² Shades-of-Gray,¹⁶ and MaxName¹⁵) explained 139 below were applied to the newly created images. Conse- 140 quently, we obtain one solution per test image per algorithm, 141 totaling 1245 different solutions. These solutions were con- 142 verted back to CIE XYZ to be displayed on a calibrated CRT 143 monitor (Viewsonic P227f, which was tested to confirm its 144 uniformity across the screen surface) using a visual stimulus 145 generator (Cambridge Research Systems ViSaGe). The 146 monitor's white point chromaticity was (x=0.315, 147)y=0.341), and its maximum luminance was 123.78 Cd/m². 148 The experiment was conducted in a dark room in which the 149 only light present in the room came from the monitor itself. 150



Figure 2. Experiment schedule.



Figure 3. Image dataset under D65 illuminant.

The experiment was conducted on ten naive observers recruited among university students and staff (none of the observers had previously seen the picture database). All observers were tested for normal color vision using the Ishihara of the Farnsworth dichotomous tests (D-15). Pairs of picties tures, each obtained using one of two different color constancy algorithms, were presented one on top of the other on sa gray background (31 Cd/m²). The order and position of picture pairs were random. Each picture subtended 160 10.5° × 5.5° to the observer and was viewed from a distance 161 of 146 cm. This brings us to 1245 pairs of observations per 162 observer. No influence on picture (top or bottom) position 163 in the observers' decision was found.

For each presentation, observers were asked to select the for each presentation, observers were asked to select the picture that seemed most natural and to rate their selection for by pressing a button on an IR button box. The setup (six for buttons) allowed observers to register how convinced they marginally convinced). For example, observers who were ro strongly convinced that the top image was more natural than ro strongly convinced that the top image was more natural than ro strongly convinced that the bottom picture was the ro marginally convinced that the bottom picture was the ro marginally convinced that the bottom picture was the ro most natural they would press button 4 and so on. There respond to each choice. The total experiment lasted approxiro mately 90 min (divided in three sessions of 30 min each).

177 A New Image Dataset

178 To test the models we need a large image dataset of good 179 quality natural scenes. From a colorimetric point of view, the 180 obvious choice is to produce hyperspectral imagery in order 181 to reduce metameric effects. However, hyperspectral outdoor 182 natural scenes are difficult to acquire since the exposure 183 times needed are long, and their capture implies control over 184 small movements or changes in the scene, (not to mention 185 the financial cost of the equipment). There are currently 186 good quality image databases available (such as the 187 hyperspectral dataset built by Foster et al.¹⁸ and Brelstaff et **188** al.¹⁹), but they either contain specialized (i.e., nongeneral) 189 imagery, or the number of scenes is not large enough for our **190** purposes. For this reason, and because metamerism is rela-**191** tively rare in natural scenes, 20,21 we decided to acquire our 192 own dataset of 83 images (see Figure 3) using a trichromatic 193 digital color camera (Sigma Foveon D10) calibrated to pro-194 duce CIE XYZ pixel representations.

The camera was calibrated at Bristol University (UK)Experimental Psychology laboratary by measuring its colorsensors' spectral sensitivities using a set of 31 narrow band



Figure 4. Camera and gray sphere setup.



Revell RAL4012-Matt

Figure 5. Reflectance of the paint used on the ball.

interference filters, a constant-current incandescent light ¹⁹⁸ source, and a TopCon SR1 telespectroradiometer (a process 199 similar to that used by others^{22,23}). The calibrated camera 200 allows us to obtain a measure of the CIE XYZ values for 201 every pixel in the image. Images were acquired around the 202 city of Barcelona at different times of the day and on three 203 different days in July 2008. The weather was mostly sunny 204 with a few clouds. We mounted a gray ball in front of the 205 camera (see Figure 4) following the ideas of Ciurea and 206 Funt.²⁴ The ball was uniformly painted using several thin 207 layers of spray paint (Revell RAL7012-Matt, whose reflec- 208 tance was approximately constant across the camera's re- 209 sponse spectrum, and its reflective properties were nearly 210 Lambertian—see Figure 5). The presence of the gray ball 211 (originally located at the bottom-left corner of every picture 212 and subsequently cropped out) allows us to measure and 213 manipulate the color of the illuminant. Images whose chro- 214 maticity distribution was not spatially uniform (as measured 215 on the gray ball) were discarded. 216

Selected Color Constancy Algorithms

In this section we briefly summarize the three methods we **218** have selected for our analysis. We have chosen two well- **219** known methods, Gray-World² and Shades-of-Gray,¹⁶ and a **220**

217

more recent method, the MaxName algorithm.¹⁵ The GrayWorld algorithm (an uncalibrated method based on a strong
assumption about the scene) was selected because of its
popularity in literature. The Shades-of-Gray algorithm (another uncalibrated algorithm) was selected because it considerably improves performance with respect to Gray-World
(another uncalibrated algorithm such as Gray-Edge²⁵ could
also have been used). Finally, MaxName¹⁵ was selected because it uses high-level knowledge to correct the illuminant.
We give a brief outline of these methods below.

(1) Gray-World. It was proposed by Buchsbaum,² and 231 it is based on the hypothesis that mean chromatic-232 ity of the scene corresponds to gray. Given an im-233 AQ: #5 age $f = (R, G, B)^T$ as a function of RGB values, and 234 adopting the diagonal model of illuminant 235 change,²⁶ then an illuminant (α, β, γ) accomplishes 236 the Gray-World hypothesis if 237

$$\frac{\int f \partial x}{\int \partial x} = k \cdot (\alpha, \beta, \gamma), \tag{1}$$

where k is a constant.

238

239

240

241

242

243

244

245

246

247

248

249

250

251

259

AQ.

AQ:

#6

(2) Shades-of-Gray. It was proposed by Finlayson and Trezzi.¹⁶ This algorithm is a statistical extension of the Gray-World and MaxRGB²⁷ algorithms. It is based on the Minkowski norm of images. An illuminant (α, β, γ) is considered as the scene illuminant if it accomplishes

$$\left(\frac{\int f^p \partial x}{\int \partial x}\right)^{1/p} = k \cdot (\alpha, \beta, \gamma), \qquad (2)$$

where *k* is a constant. Actually, this is a family of methods where p=1 is the Gray-World method and $p=\infty$ is the MaxRGB algorithm. In this case we have used p=12, since it is the best solution for our dataset.

(3) MaxName. This algorithm is a particular case of the one presented by Vazquez et al.¹⁵ It is based on giving more weight to those illuminants that maximize the number of color names in the scene. That is, MaxName builds a weighted feasible set by considering nameable colors; this is prior knowledge given by

 $\mu_k = \int_{\omega} S(\lambda) E(\lambda) R_k(\lambda) \partial \lambda, \quad k = R, G, B, \quad (3)$

where, $S(\lambda)$ are the surface reflectances having 260 maximum probability of being labeled with a basic 261 color term, also called focal reflectances from the 262 work of Benavente et al.²⁸ In addition to the basic 263 color terms, we added a set of skin colored 264 reflectances. In Eq. (3), $E(\lambda)$ is the power distribu-265 tion of a D65 illuminant, and $R_k(\lambda)$ are the CIE 266 RGB 1955 color matching functions. 267

We define μ as the set of all *k*-dimensional nameable ²⁶⁸ colors obtained from Eq. (3). The number of elements of μ ²⁶⁹ depends on the number of reflectances used. Following this, ²⁷⁰ we compute the *semantic matrix*, denoted as (SM), which is ²⁷¹ a binary representation of the color space as a matrix, where ²⁷² a point is set to 1 if it represents a nameable color, that is, it ²⁷³ belongs to μ and 0 otherwise. Then, for a given input image, ²⁷⁴ *I*, we compute all possible illuminant changes $I_{\alpha,\beta,\gamma}$. For each ²⁷⁵ one, we calculate its nameability value. This is done by ²⁷⁶ counting how many points of the mapped image are name- ²⁷⁷ able colors in SM and can be computed by a correlation in ²⁷⁸ log space: ²⁷⁹

$$Nval_{\alpha,\beta,\gamma} = \log(H_{bin}(I)) * \log(SM).$$
(4) 280

In the previous equation, H_{bin} is the binarized histo- 281 gram of the image, *Nval* at the position (α, β, γ) is the num- 282 ber of coincidences between the SM and $I_{\alpha,\beta,\gamma}$. *Nval* is a 283 three-dimensional matrix, depending on all the feasible 284 maps, (α, β, γ) . From this matrix, we select the most feasible 285 illuminant as the one that accomplishes 286

$$(\alpha, \beta, \gamma) = \arg \max_{(\alpha, \beta, \gamma)} Nval, \tag{5}$$

290

that is, the one giving the maximum number of nameable 288 colors. 289

RESULTS

The results of the experiment validate those presented by 291 Vazquez et al.¹⁵ with a different image dataset and a different 292 set of algorithms. The main finding is that preferred solu- 293 tions, namely, the more natural in the psychophysical experi- 294 ment, do not always coincide with solutions of minimum 295 angular error. In fact, this agreement only happened in 43% 296 of the observations, independently of the degree of certainty 297 of the observers when making the decision. 298

Since the experimental procedure allows us to define a 299 partition in the interval [0,1] to encode the subject selection 300 and each observation represents a decision between two im- 301 ages, then for each observation we label one image as the 302 result from Method A and the other as the result from 303 Method B (Methods A and B are labeled as 1 and 0, respec- 304 tively). The confidence of the decision is considered at three 305 different levels (the three buttons that the subject was al- 306 lowed to press yield an ordinal paired comparison²⁹). For 307 example, suppose that a scene processed by Method A is 308 presented on top of the screen and a second scene processed 309 by Method B is presented at the bottom (the physical posi- 310 tion of the scenes was randomized in each trial, but let us 311 consider an exemplary layout). If subjects think that the top 312 picture is more natural they will press one of the top buttons 313 in Fig. 2, according to how strongly they are convinced. Sup- 314 pose the subject presses button 3 (top-right: definitely more 315 natural), then the response is coded as 1. If the choice is 316 button 2 (top-center: sufficiently more natural) the response 317 is coded as 0.8, etc. (see Table I). If, on the contrary subjects 318 think the bottom picture (Method B) is more natural, then 319

J. Imaging Sci. Technol.

Vazquez et al.: Color constancy algorithms: Psychophysical evaluation on a new dataset



Table I. Button codification.

Figure 6. Comparison to the mean observer (black line).

Table II. Correlation between each observer and mean observer.

Observer	1	2	3	4	5	6	7	8	9	10
Correlation	0.54	0.57	0.59	0.55	0.52	0.23	0.48	0.63	0.61	0.55
CV	52,49%	57,96%	37,65%	52,28 %	52,69%	59,85%	47,12%	51,13%	25,36%	42,8 1%

³²⁰ they will press a button from the lower row (Fig. 2). If they ³²¹ are marginally convinced, they will pick button 4 (bottom-³²² left), and the response will be coded as 0.4 according to ³²³ Table I. Similarly if they are strongly convinced, they will ³²⁴ press button 6 (bottom-right), and the response will be ³²⁵ coded as 0. In this way we collect not only the direction of ³²⁶ the response but its certainty. Observers' certainty was found ³²⁷ to be correlated (corr. coef. 0.726) to a simple measure of ³²⁸ image difference (the angular error between each image ³²⁹ pair). This technique is similar to that used by other ³³⁰ researchers.^{30–33}

We have computed two different measures of observer variability. The first measure is the correlation coefficient between individual subjects and the average (in black in Figure 6). Table II shows this measure. The idea behind this ³³⁴ analysis is to detect outliers (subjects with a distribution of ³³⁵ results significantly different from the rest of the observers, ³³⁶ i.e., low correlation). Our second measure is the coefficient ³³⁷ of variation (CV),^{^{34,35}} which computes the difference be- ³³⁸ tween two statistical samples (see Table II). Both measures ³³⁹ were calculated for the whole 1245 observations (three com- ³⁴⁰ binations of color constancy solutions × 415 observations ³⁴¹ per combination). From the table, and from the distribution ³⁴² of the plots in Fig. 6, we decided to omit data from observer ³⁴³ 6 (very low correlation coefficient and highest coefficient of ³⁴⁴ variation) in all subsequent analyses. ³⁴⁵

As a first approach to analyze our results we computed **346** the mean of the observers' responses for each pairwise com- **347**

Table III. Results of the experiment in the one-to-one comparison.

Selected method vs method	Shades-of-Gray (%)	Gray-World (%)	MaxName (%)
Shades-of-Gray	_	68.1	50.6
Gray-World	31.9	_	37.6
MaxName	49.4	62.4	—

Table IV. Experiment results in a general comparison.

Method	Wins (%)
Shades-of-Gray	35.18
Gray-World	16.63
MaxName	39.28
Three-equally selected	8.92

TABLE V. RESULTS USING THUISIONES INW OF COMPANIANCE [DAGING	Table	: V.	Results	s using	Thurstone's	law of	comparative	judgment
--	-------	------	---------	---------	-------------	--------	-------------	----------

Method	Wins (%)
Shades-of-Gray	42.65
MaxName	36.39
Gray-World	20.96

³⁴⁸ parison. We considered that a method was selected if the ³⁴⁹ mean of the encoded decisions, computed for all nine ob-³⁵⁰ servers, is greater than 0.5 (when the method was encoded as ³⁵¹ 1) or lower than 0.5 (when the method was encoded as 0). ³⁵² The performance does not vary significantly if we do not ³⁵³ consider the cases where the average value is too close to the ³⁵⁴ chance rate (e.g., averages between 0.45 and 0.55). The re-³⁵⁵ sults of these pairwise comparisons are given in Table III. ³⁵⁶ For each pair of methods, we show the percentage of cases ³⁵⁷ where it has been selected against the others. Thus, results in ³⁵⁸ Table III can be interpreted as follows: each method (in ³⁵⁹ rows) is preferred to a certain percentage of trials over the ³⁶⁰ method in the columns. For example, Shades-of-Gray is pre-³⁶¹ ferred in 68.1% of the trials against Gray-World. The percentages in Table III show that the images produced by Shades-of-Gray and MaxName are preferred to those produced by Gray-World (68.1% and 62.4%). However, there is no clear preference when compared against each other (50.6% Shades-of-Gray preference versus MaxName).

In Table IV we show a global comparison of all algo- 368 rithms (the percentages are computed for all 415 images). A 369 method was considered a "winner" for a given image if it 370 was selected in two of the three comparisons. Methods were 371 evaluated in the same way as we did for results in Table III 372 (that is, a greater than 0.5 mean value from all observers is 373 encoded as 1). Evaluating this way, there are some cases 374 where the three methods are equally selected (this happens 375 in 8.92% of the images). This analysis was formulated in 376 order to remove nontransitive comparisons (e.g., Method A 377 beats Method B, Method B beats Method C, and Method C 378 beats Method A). Hence, we can conclude from these 379 straightforward analyses that solutions from MaxName are 380 preferred in general but are closely followed by Shades-of- 381 Gray (39.28% and 35.18%, respectively). We can also state 382 that Gray-World solutions are the least preferred in general 383 (with a low percentage of 16.63%). Moreover, the best an- 384 gular error solution is selected in 42.96% of the cases. 385

We have also calculated Thurstone's law of comparative judgment³⁶ coefficients from our data (Table V), obtained from the ordinal pairwise comparisons. Using this measure, results are not very different (Shades-of-Gray and MaxName are clearly better than Gray-World although the ranking changes), and images with minimal angular error are only selected in 45% of the cases.

Finally, we have computed two overall analyses (consid- 393 ering all scenes as one) in order to extract a global ranking 394 for our color constancy methods: Thurstone's law of com- 395 parative judgment³⁶ and the Bradley–Terry³⁷ analysis. Table 396 VI shows the results of Bradley and Terry's cumulative logit 397 model for pairwise evaluations extended to ordinal 398 comparisons.²⁹ These results are shown in the "estimate" 399 column where the estimate reference has been set to 0 for 400 the smallest value (Gray-World model). The standard error 401 of this ranking measure shows that the two best models 402 (Shades-of-Gray and MaxName) are better than Gray-World 403 and arguably close to each other. Table VII shows a similar 404 analysis using Thurstone's law of comparative judgment³⁶ 405 and considering all scenes as one. 406

AQ:

Table VI. Results using Bradley–Terry ordinal pairwise comparison analysis.

Parameter	DF	Estimate	Standard error	Wald 95% confidence li	mits	Chi-square	Pr > Chisq
Shades-of-Gray	1	1.609	1.2231	-0.7882	4.0063	1.73	0.1883
MaxName	1	1.0256	0.8435	-0.6278	2.6789	1.48	0.2241
Gray-World	0	0	0	0	0		

Table VII. Results using Thurston's law of comparative judgment binary pairwise comparison analysis.

Parameter	DF	Estimate	Standard error	Wald 95% confidence	limits	Chi-square	Pr > Chisq
Shades-of-Gray	1	0.196	0.0031	0.19	0.2021	4040.2	< 0.0001
MaxName	1	0.1283	0.0031	0.1223	0.1343	1743.22	< 0.0001
Gray-World	0	0	0	0	0		

As we mentioned above, our experiment shows that im-408 ages having minimum angular error with respect to the ca-409 nonical solution are selected in less than half the observa-410 tions (when we ask people for the most natural image, the 411 response does not always correspond to the optimal physical 412 solution). Moreover, this result is maintained even if we dis-413 card responses with low levels of certainty. In order to quan-414 tify this fact, in the next section we will introduce a new 415 measure to complement the current performance evaluation 416 of color constancy algorithms.

417 PERCEPTUAL PERFORMANCE EVALUATION

 Assuming the ill-posed nature of the problem, the difficulty of finding an optimal solution and the results of the present experiment, we propose an approach to color constancy al- gorithms that involves human color constancy by trying to match computational solutions to perceived solutions. Hence, we propose a new evaluation measurement, the *per- ceptual angular error*, which is based on perceptual judg- ments of adequacy of a solution instead of the physical so- lution. The approach that we propose in this work does not try to give a line of research alternative to the current trends, which focus on classifying scene contents to efficiently com- bine different methods. Here we try to complement these efforts from a different point of view that we could consider as more "top-down," instead of the "bottom-up" nature of the usual research.

As mentioned above, the most common performance
evaluation for color constancy algorithms consists of measuring how close their proposed solution is to the physical
solution, independently of the other concerns. This has been
computed as

$$e_{ang} = a \cos\left(\frac{\rho_w \hat{\rho}_w}{\|\rho_w\|\|\hat{\rho}_w\|}\right),\tag{6}$$

 which represents the angle between the actual white point of the scene illuminant, ρ_w , and the estimation of this point given by the color constancy method, $\hat{\rho}_w$, which can be un- derstood as a chromaticity distance between the physical so- lution and the estimate. The current consensus is that none of the current algorithms present a good performance on all the images,³⁸ and a combination of different algorithms of- fers a promising option for further research. Our proposal here is to introduce a new measure, the *perceptual angular error*, e_{ane}^{p} , that would be computed in a similar way:



Figure 7. Estimated perceptual angular error (between method estimations and preferred illuminants).

$$e_{ang}^{p} = a \cos\left(\frac{\rho_{w}^{p} \hat{\rho}_{w}}{\|\rho_{w}^{p}\|}\right), \qquad (7)$$

where ρ_w^p is the perceived white point of the scene (which 450 should be measured psychophysically) and $\hat{\rho}_w$ is an estima- 451 tion of this point, that is the result of any color constancy 452 method, as in Eq. (6). The difficulty of this new measure- 453 ment arises from the complexity of building a large image 454 dataset, where ρ_w^p , the perceived white point of the images 455 has been measured.

In this work we propose a simple estimation of this **457** perceived white point by considering the images preferred in **458** the previous experiment. Hence, the perceived white point is **459** given by the images coming from the color constancy solu- **460** tions that have been preferred by the observers. The pre- **461** ferred solutions, that is, the most natural solutions, can give **462** us an approximation to the perceived image white point. **463**

Making the above consideration, in Figure 7 we can see 464 how the estimation of the perceptual angular error works for 465 the three tested algorithms. In the abscissa we plot a ranking 466 of the observations in order to get the perceptual errors in 467 descending order. In the ordinate we show the estimated 468 perceptual angular error for each created image (that is, 415 469 different inputs to the algorithms). A numerical estimation 470 of the perceptual angular error could be the area under the 471 curves plotted in Fig. 7. In the figure we can see that both 472 Shades-of-Gray and MaxName work quite similarly, while 473 Gray-World presents the highest perceptual error. This new 474 measurement agrees with the conclusion we summarized in 475 the previous section and provides a complementary measure 476

438



Figure 8. Angular error between methods estimations and canonical illuminant.

Table VIII. Angular error for the different methods on 415 images of the dataset.

Method	Mean	rms	Median
MaxName	7.64°	8.84°	6.78°
Shades-of-Gray	7.84°	9.70°	5.95°
Gray-World	10.05°	12.70°	7.75°

⁴⁷⁷ to evaluate color constancy algorithms. In Figure 8 we show 478 a similar plot for the usual angular error.

In Tables VIII and IX we show the different statistics on the computed angular errors. In Table VIII, the angular error between the estimated illuminant and the canonical illuminant are shown. In this case, MaxName and Shades-of-Gray present better results than Gray-World. In Table IX equal statistics are computed for the estimated perceptual angular error. The results in this table confirm the conclusions we obtained from Fig. 7.

487 CONCLUSIONS

488 This paper explores a new research line, the psychophysical 489 evaluation of color constancy algorithms. Previous research 490 points to the need to further explore the behavior of high-491 level constraints needed for the selection of a feasible solu-492 tion (to avoid the dependency of current evaluations on the 493 statistics of the image dataset). With this aim in mind, we 494 have performed a psychophysical experiment in order to 495 compare three computational color constancy algorithms: 496 Shades-of-Gray, Gray-World, and MaxName. The results of 497 the experiment show Shades-of-Gray and MaxName meth-498 ods have quite similar results, which are better than those 499 obtained by the Gray-World method and that in almost half 500 the judgments, subjects have preferred solutions that are not 501 the closest ones to the optimal solutions.

502 Considering that subjects do not prefer the optimal so-503 lutions in a large percentage of judgments, we have intro-504 duced a new measure based on the perceptual solutions to 505 complement current evaluations: the perceptual angular er-

 Table IX.
 Estimated perceptual angular error for the different methods on 415 images of the dataset.

Method	Mean	rms	Median
MaxName	3.86°	6.02°	2.61°
Shades-of-Gray	3.79°	5.66°	2.86°
Gray-World	6.70°	9.01°	5.85°

ror. It tries to measure the proximity of the computational ⁵⁰⁶ solutions versus the human color constancy solutions. The ⁵⁰⁷ current experiment allows computing an estimation of the ⁵⁰⁸ perceptual angular error for the three explored algorithms. ⁵⁰⁹ However, our main conclusion is that further work should ⁵¹⁰ be done in the line of building a large dataset of images ⁵¹¹ linked to the perceptually preferred judgments. ⁵¹²

To this end a new, more complex experiment, perhaps **513** related to the one proposed in Ref. 39, must be done in **514** order to obtain the perceptual solution of the images inde- **515** pendently of the algorithms being judged. **516**

Acknowledgments

This work has been partially supported by projects **518** TIN2004-02970, TIN2007-64577, and Consolider-Ingenio **519** 2010 CSD2007-00018 of Spanish MEC (Ministry of Sci- **520** ence). C.A.P. was funded by the Ramon y Cajal research **521** programme of the MEC (RYC-2007-00484). The authors **522** thank J. van de Weijer for his insightful comments. **523**

REFERENCES

- ¹S. Hordley, "Scene illuminant estimation: Past, present, and future," 525 Color Res. Appl. 31, 303 (2006). 526
- ²G. Buchsbaum, "A spatial processor model for object colour **527** perception," J. Franklin Inst. **310**, 1 (1980). **528**
- ³V. C. Cardei, B. Funt, and K. Barnard, "Estimating the scene **529** illumination chromaticity by using a neural network," J. Opt. Soc. Am. **530** A **19**, 2374 (2002). **531**
- ⁴G. Finlayson, S. Hordley, and R. Xu, "Convex programming colour 532 constancy with a diagonal-offset model," *Proc. International Conference* 533 *on Image Processing (ICIP)* (IEEE Computer Society, Los Alamitos, CA, 534 2005) pp. 2617–2620. 535
- ⁵K. Barnard, "Improvements to gamut mapping colour constancy 536 algorithms," *Proc. European Conference on Computer Vision (ECCV)* 537 (Springer, Berlin, 2000) pp. 390–403. 538
- ⁶G. Finlayson, P. Hubel, and S. Hordley, "Color by correlation," Proc. 539 IS&T/SID 5th Color Imaging Conference (IS&T, Springfield, VA, 1997) 540 pp. 6–11. 541
- ⁷B. Funt, M. Drew, and J. Ho, "Color constancy from mutual reflection," 542 Int. J. Comput. Vis. 6, 5 (1991).
- ⁸K. Barnard, V. Cardei, and B. Funt, "A comparison of computational 544 color constancy algorithms. I: Methodology and experiments with 545 synthesized data," IEEE Trans. Image Process. 11, 972 (2002). 546
- ⁹K. Barnard, L. Martin, A. Coath, and B. Funt, "A comparison of 547 computational color constancy algorithms. II: Experiments with image 548 data," IEEE Trans. Image Process. 11, 985 (2002). 549
- ¹⁰ S. Hordley and G. Finlayson, "Re-evaluating colour constancy 550 algorithm," *Proc. 17th International Conference on Pattern Recognition* 551 (IEEE Computer Society, Los Alamitos, CA, 2004) pp. 76–79. 552
- ¹¹V. Cardei and B. Funt, "Committee-based color constancy," *Proc.* 553
 IS&T/SID 7th Color Imaging Conference (IS&T, Springfield, VA, 1999) 554
 ¹²A. Giisenii and T. Course, "Col.
- ¹²A. Gijsenij and T. Gevers, "Color constancy using natural image 556 statistics," *Proc. 2007 IEEE Conference on Computer Vision and Pattern* 557 *Recognition*, Vols. 1–8 (IEEE Computer Society, Los Alamitos, CA, 2007) 558 pp. 1806–1813.
- ¹³F. Tous, "Computational framework for the white point interpretation 560

524

517

- 561 base on color matching," Ph.D. thesis, Universitat Autònoma de 562 Barcelona, Barcelona (2006) (unpublished).
- ¹⁴ J. V. van de Weijer, C. Schmid, and J. Verbeek, "Using high-level visual 563 information for color constancy," Proc. International Conference on 564 565 Computer Vision (IEEE Computer Society, Los Alamitos, CA, 2007).
- 15 J. Vazquez, M. Vanrell, R. Baldrich, and C. A. Párraga, "Towards a 566
- 567 psychophysical evaluation of colour constancy algorithms," Proc. 568 IS&TsCGIV 2008/MCS/08-4th European Conference on Colour in
- Graphics, Imaging, and Vision (IS&T, Springfield, VA, 2008) pp. 372-377. 569 570 G. Finlayson and E. Trezzi, "Shades of gray and colour constancy," Proc. 571 IS&T/SID 12th Color Imaging Conference (IS&T, Springfield, VA, 2004)
- pp. 37-41. 572 573 ⁷D. A. Forsyth, "A novel algorithm for color constancy," Int. J. Comput. Vis. 5, 5 (1990). 574
- ¹⁸D. H. Foster, S. M. C. Nascimento, and K. Amano, "Information limits 575 on neural identification of colored surfaces in natural scenes," Visual 576
- 577 Neurosci. 21, 331 (2004). ⁹G. J. Brelstaff, C. A. Parraga, T. Troscianko, and D. Carr, "Hyperspectral camera system: Acquisition and analysis," Proc. SPIE **2587**, 150–159 578 579
- 580 (1995).²⁰D. H. Foster, K. Amano, S. M. C. Nascimento, and M. J. Foster, 581 "Frequency of metamerism in natural scenes," J. Opt. Soc. Am. A 23, 582
- 583 2359 (2006). ²¹M. G. A. Thomson, S. Westland, and J. Shaw, "Spatial resolution and 584
- 585 metamerism in coloured natural scenes," Perception 29, 123 (2000).
- ²²A. Olmos and F. A. A. Kingdom, "A biologically inspired algorithm for 586 the recovery of shading and reflectance images," Perception 33, 1463 587 588 (2004).
- ²³C. A. Párraga, T. Troscianko, and D. J. Tolhurst, "Spatiochromatic 589 properties of natural images and human vision," Curr. Biol. 12, 483 590 591 (2002)
- ²⁴ F. Ciurea and B. Funt, "A large image database for color constancy 592 research," Proc. IS&T/SID 11th Color Imaging Conference (IS&T, 593
- 594 Springfield, VA, 2003) pp. 160-164.
- ²⁵ J. V. van de Weijer, T. Gevers, and A. Gijsenij, "Edge-based color 595

- 596 constancy," IEEE Trans. Image Process. 16, 2207-2214 (2007).
- ²⁶G. Finlayson, M. Drew, and B. Funt, "Diagonal transforms suffice for **597** color constancy," Proc. 4th International Conference on Computer Vision 598 (IEEE Computer Society, Los Alamitos, CA, 1993) pp. 164-171. 599
- ²⁷ E. Land, "Retinex theory of color-vision," Sci. Am. 237, 108 (1977).
 ⁶⁰⁰
 ²⁸ R. Benavente, M. Vanrell, and R. Baldrich, "Estimation of fuzzy sets for 601
- computational colour categorization," Color Res. Appl. **29**, 342 (2004). **602** A. Agresti, An Introduction to Categorical Data Analysis (Wiley, New 603
- York and Chichester, 1996) pp. 436-439. 604 ³⁰ P. Courcoux and M. Semenou, "Preference data analysis using a paired 605
- comparison model," Food Qual. Preference 8, 353 (1997). 606 ³¹G. Gabrielsen, "Paired comparisons and designed experiments," Food 607 608
- Qual. Preference 11, 55 (2000). ³² J. Fleckenstein, R. A. Freund, and J. E. Jackson, "A paired comparison **609** test of typewriter carbon papers," Tappi J. 41, 128 (1958). 610
- A. Agresti, "Analysis of ordinal paired comparison data," J. R. Stat. Soc., 611 Ser. C, Appl. Stat. 41, 287 (1992). 612
- ³⁴M. Luo, A. Clarke, P. Rhodes, A. Schappo, S. Scrivener, and C. Tait, 613 'Quantifying color appearance I. LUTCHI color appearance data," Color 614 Res. Appl. 16, 166 (1991). 615
- ³⁵ M. Luo, A. Clarke, P. Rhodes, A. Schappo, S. Scrivener, and C. Tait, 616 "Quantifying color appearance II. Testing color models performance 617
- using LUTCHI color appearance data," Color Res. Appl. **16**, 181 (1991). **618** ³⁶L. Thurstone, "A law of comparative judgment," Psychol. Rev. **34**, 273 **619** (1927). 620
- $^{37}\mbox{R}.$ A. Bradley and M. B. Terry, "Rank analysis of incomplete block 621 designs: I. the method of paired comparisons," Biometrika 39, 22 622 (1952).623
- ³⁸ B. Funt, K. Barnard, and L. Martin, "Is machine colour constancy good 624 enough?" Proc. 5th European Conference on Computer Vision (Springer, 625 Berlin, 1998) pp. 445-459. 626
- ³⁹ P. D. Pinto, J. M. Linhares, and S. M. Nascimento, "Correlated color 627 temperature preferred by observers for illumination of artistic 628 paintings," J. Opt. Soc. Am. A 25, 623 (2008). 629

AUTHOR QUERIES — 009903IST

- #1 Au: Please supply accepted date.
- #2 Au: Please check change to sentence beginning with "The previous..." to ensure meaning is preserved.
- #3 Au: Please check change to sentence beginning with "This decision..." to ensure meaning is preserved.
- #4 Au: Please verify definitions for CIE and sRGB
- #5 Au: Please verify spelling of Buchsbaum 2, Please see Reference list.
- #6 Au: Please verify Finlayson and Trezzi. 16
- #7 Au: Please verify Benavente et al.28.
- #8 Au: Please check change to sentence begning with "Thus, results..." to ensure meaning is preserved.
- #9 Au.: Please verify all information for all references
- #10 Au: Please define DF.