

Towards a Perceptual Evaluation Framework for Lighting Estimation

Supplementary Material

Anonymous CVPR submission

Paper ID 5266

001 This document complements our main paper, providing
 002 the following supplementary information regarding the con-
 003 ducted experiments and analysis:

- 004 • Additional description of the lighting estimation methods
 005 (see sec 3.2. of the main paper) used in sec. 1.1;
- 006 • A more in-depth description of the scene selection (see
 007 sec 3.2. of the main paper) given as input to the lighting
 008 estimation methods in sec. 1.2;
- 009 • Details on the geometry, materials, and rendering of the
 010 stimuli used in the psychophysical study (see sec 3.2. of
 011 the main paper), in sec. 1.3;
- 012 • Description of the hardware used for the psychological
 013 experiment (see sec 3.3. of the main paper) in sec. 1.4;
- 014 • Additional information on the procedure during the psy-
 015 chophysical experiments (see sec 3.3. of the main paper),
 016 in sec. 1.5;
- 017 • A more in-depth analysis of the participants in the psy-
 018 chological experiment (see sec 3.3. of the main paper), in
 019 sec. 1.6;
- 020 • Additional analysis of the scores obtained in the psy-
 021 chophysical experiments (see sec 4.2. of the main paper),
 022 with examples (in sec. 2.1), per image score (in sec. 2.2),
 023 and agreement for individual observers (in sec. 2.3);
- 024 • Additional analysis of the scores of the various metrics
 025 (see sec 5. of the main paper), in sec. 3;
- 026 • Additional comparisons between different network archi-
 027 tectures and an analysis of the selected network (see sec
 028 6.1. of the main paper), in sec. 4.1;
- 029 • Additional content regarding the generalisation study con-
 030 ducted with the selected architecture (see sec 6.2. of the
 031 main paper), in sec. 4.2.

032 1. Psychophysical experiment

033 The various steps of our psychophysical study are described
 034 in the following sections. The lighting estimation methods
 035 used are described in sec. 1.1, and the selected scene given as
 036 input to them is detailed in sec. 1.2. The design of the stimuli
 037 is explained in sec. 1.3. The hardware and the procedure
 038 of the experiment are discussed in sec. 1.4 and sec. 1.5,

039 respectively. The study participants are described in more
 040 detail in sec. 1.6.

041 1.1. Lighting estimation methods

042 In the following sections, we describe the lighting estimation
 043 methods used in the indoor and outdoor psychophysical
 044 studies (see sec 3.2. of the main paper).

045 **Environment map.** We consider three state-of-the-art
 046 non-parametric lighting estimation methods to light our vir-
 047 tual scene. Weber *et al.* [23] proposes a two-stage indoor-
 048 only approach, where a dominant light source and scene
 049 layout are estimated and given as input to a texture network
 050 to predict the entire environmental texture based on the in-
 051 put image. EverLight [6] proposes a different two-stage
 052 method working simultaneously indoors and outdoors. It
 053 first estimates the lighting parameters as spherical gaussians,
 054 and then integrates them into environment map generation
 055 via guided co-modulation [5]. StyleLight [21] is a recent
 056 method that leverages the training of StyleGAN [12] with
 057 GAN inversion [17] to predict complete 360° environment
 058 maps from input images.

059 **Parametric.** We consider Gardner *et al.* [8] to provide
 060 parametric indoor lighting estimations. This method predicts
 061 three light sources parameterised by their direction, distance,
 062 angular size, and colour. For the outdoor parametric lighting
 063 model, we chose Zhang *et al.* [24], which trains a network to
 064 directly estimate the Lalonde-Matthews sky parameters [14]
 065 from a given outdoor image.

066 **Classical.** In contrast to these recent sophisticated
 067 learning-based methods, we also include Khan *et al.* [13].
 068 This technique lacks any learning components and instead
 069 determines the lighting conditions by projecting the back-
 070 ground image onto a sphere and then mirroring it to generate
 071 a complete LDR environment map.

072 All the environment maps generated by the indoor and
 073 outdoor lighting estimation methods for the user study are
 074 presented in fig. 1. The first and ninth columns display
 075 the input given to the lighting estimation methods (more
 076 details in sec. 1.2), extracted from the ground-truth panora-
 077 mas, shown in the third and eleventh columns. The second

078 and tenth columns show the reconstructed first-order spher-
079 ical harmonics, used to select the scenes (more details in
080 sec. 1.2).

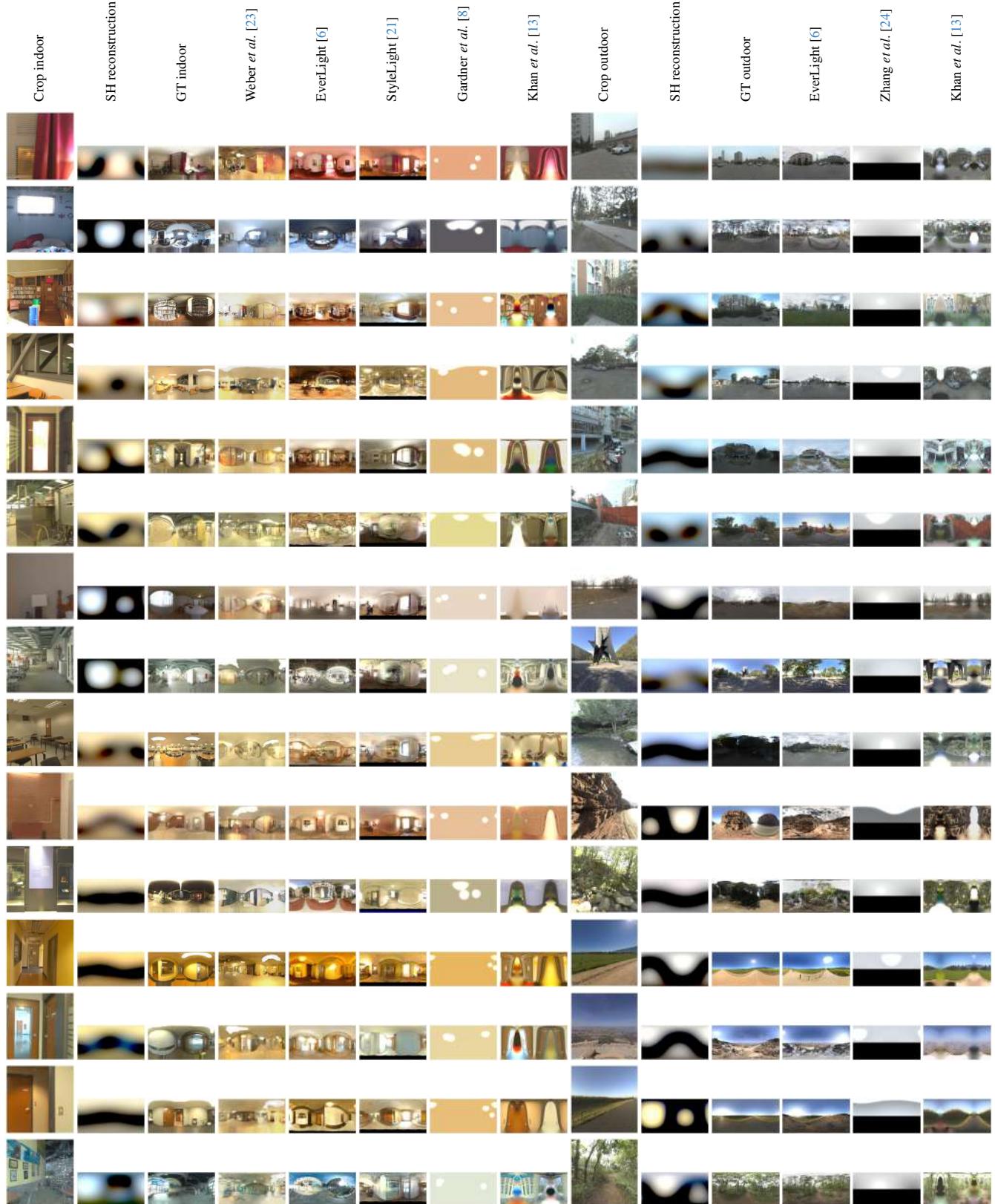


Figure 1. IBLs generated by the different indoor and outdoor lighting methods (columns) for each scene (rows). The first and ninth columns correspond to the region extracted from the indoor/outdoor scene, corresponding to a $50^\circ/90^\circ$ FoV. This region is taken from the centre of the full GT panorama (for most scenes), shown in the third and eleventh columns. The second and tenth column correspond to the reconstruction of the first-order spherical harmonics, showing the variety of the lighting in the selected scenes. The IBLs are reexposed and tonemapped with $\gamma = 2.2$ for display.

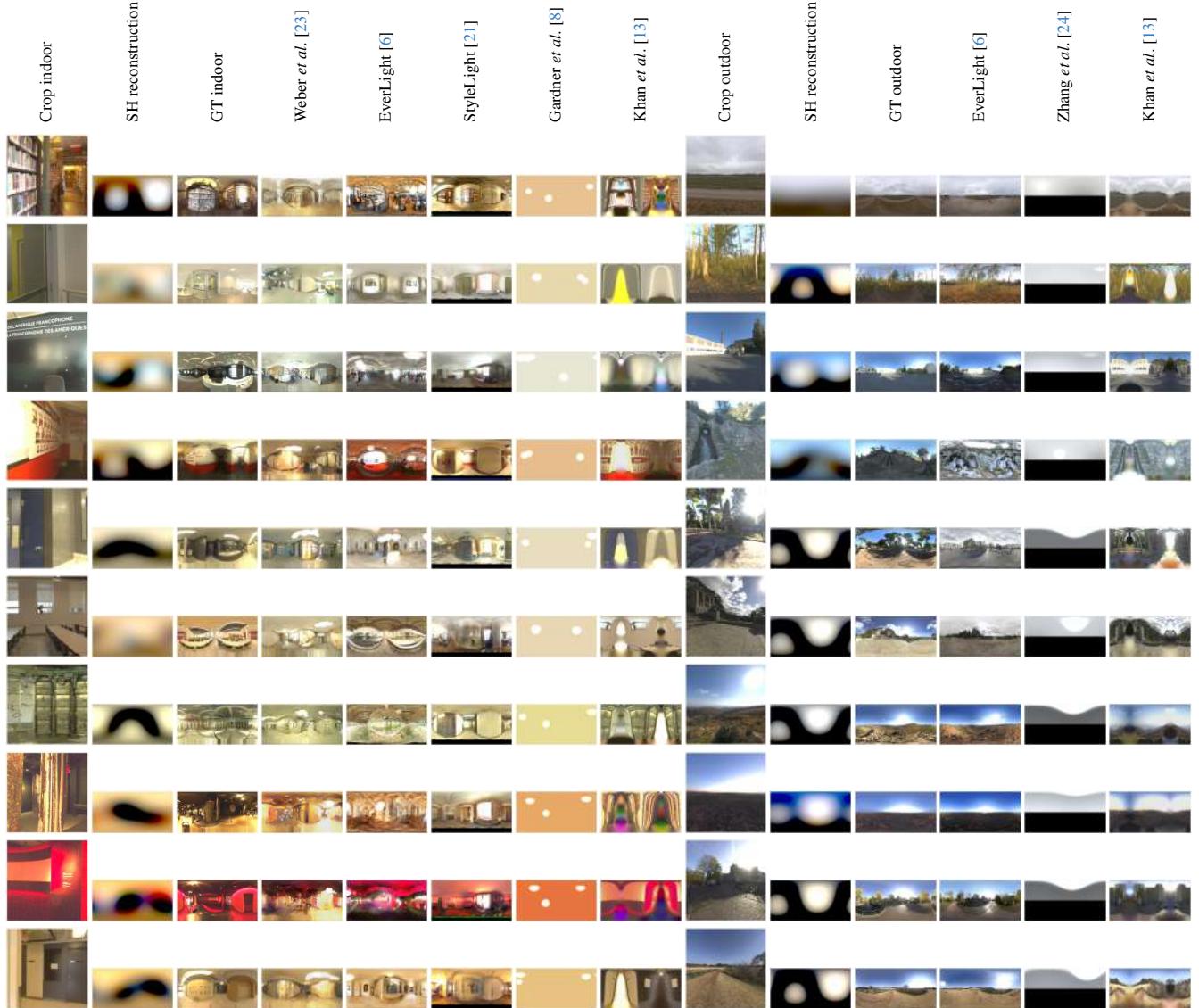


Figure 1. (contd) IBLs generated by the different indoor and outdoor lighting methods (columns) for each scene (rows). The first and ninth columns correspond to the region extracted from the indoor/outdoor scene, corresponding to a $50^\circ/90^\circ$ FoV. This region is taken from the centre of the full GT panorama (for most scenes), shown in the third and eleventh columns. The second and tenth column correspond to the reconstruction of the first-order spherical harmonics, showing the variety of the lighting in the selected scenes. The IBLs are reexposed and tonemapped with $\gamma = 2.2$ for display.

081 We utilise the output generated by these models as environment maps to illuminate the synthetic scene (sec. 1.3).
082 For the parametric models, we initially convert the output
083 parameters into environment maps and employ them in ren-
084 dering stimuli.
085

086 1.2. Lighting estimation input scenes

087 High dynamic range (HDR) panorama images are used to
088 extract limited FoV low dynamic range (LDR) regions to
089 give as input to the lighting estimation methods (see sec 3.2.
090 of the main paper). In our assessment of indoor lighting
091 estimation techniques, we adhered to the procedures out-
092 lined in Weber *et al.* [23]. Our evaluation was conducted
093 on the test set of Laval indoor dataset [7], comprising 224
094 high-resolution HDR panoramas. Within this test set, we
095 systematically extracted 10 LDR images from each of the
096 panoramas using the sampling distribution identical to We-
097 ber *et al.* [23]. This process yielded a grand total of 2240
098 images for our evaluation. We adopted the approach detailed
099 in [6] to assess outdoor lighting. This method leverages
100 839 distinct outdoor HDR panoramas sourced from SHLight
101 dataset [3]. From these, it derives three LDR images ac-
102 cording to the sampling distribution of [6], resulting in an
103 evaluation set comprising 2517 images.

104 **Region extraction.** To obtain the images given as input
105 to the lighting estimation methods, an FoV of $50^\circ/90^\circ$ is ex-
106 tracted from the centre of the indoor/outdoor HDR panorama,
107 which is tonemapped with $\gamma = 2.2$ and reexposed. The ex-
108 tracted regions have a resolution of 512×512 for the indoor
109 images and 256×256 for the outdoor images. Examples of
110 the extracted regions used in this study are shown in the first
111 and ninth columns of fig. 1.

112 **Scene selection.** 25 scenes are selected from the indoor
113 and outdoor datasets. We limited the number of scenes in our
114 study to keep the experiment time below ~ 30 min, in order
115 to avoid errors caused by observers' fatigue. 25 scenes are
116 considered sufficient to represent different types of environ-
117 ments with diverse lighting. In order to have a great variety
118 of lighting environments, the coefficients of the first-order
119 spherical harmonics are extracted from the HDR panora-
120 mas, using the `skylibs` Python library. The scenes are
121 selected by taking the medoid of the clusters obtained using
122 the k-means algorithm, where $k = 25$, from the `sklearn`
123 Python library. The resulting clusters for the indoor (left)
124 and outdoor (right) are shown in fig. 2. Examples of the
125 first-order spherical harmonics reconstruction of the selected
126 scenes are shown in the second and tenth columns of fig. 1,
127 which indeed demonstrates a great variety. The scenes that
128 contain too much noise are removed from the dataset, as
129 they would potentially distract the observers from judging
130 the realism of the inserted virtual object.

131 1.3. Stimuli

132 The stimuli for tasks 1 and 2 (see sec 3.2. of the main paper)
133 have different geometries (described below), and each task
134 has a diffuse and glossy variation, with details regarding the
135 materials given subsequently. The rendering details are also
136 indicated.

137 For reference, the stimuli used in the indoor and outdoor
138 psychophysical experiments, for the diffuse/glossy experi-
139 ments are shown in fig. 3/fig. 4 for task 1 and fig. 5/fig. 6 for
140 task 2.

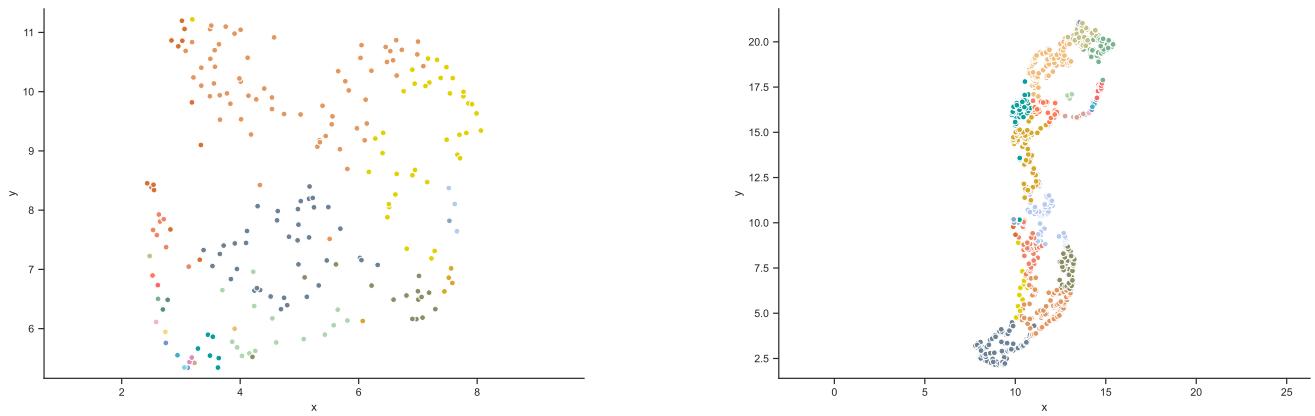


Figure 2. Clusters obtained k -means algorithm, where $k = 25$, for the indoor (left) and outdoor (right) high dynamic range (HDR) panorama datasets. The projection to \mathbb{R}^2 is done using UMAP [15]. The different colours indicate the different clusters. The axes are in arbitrary units.

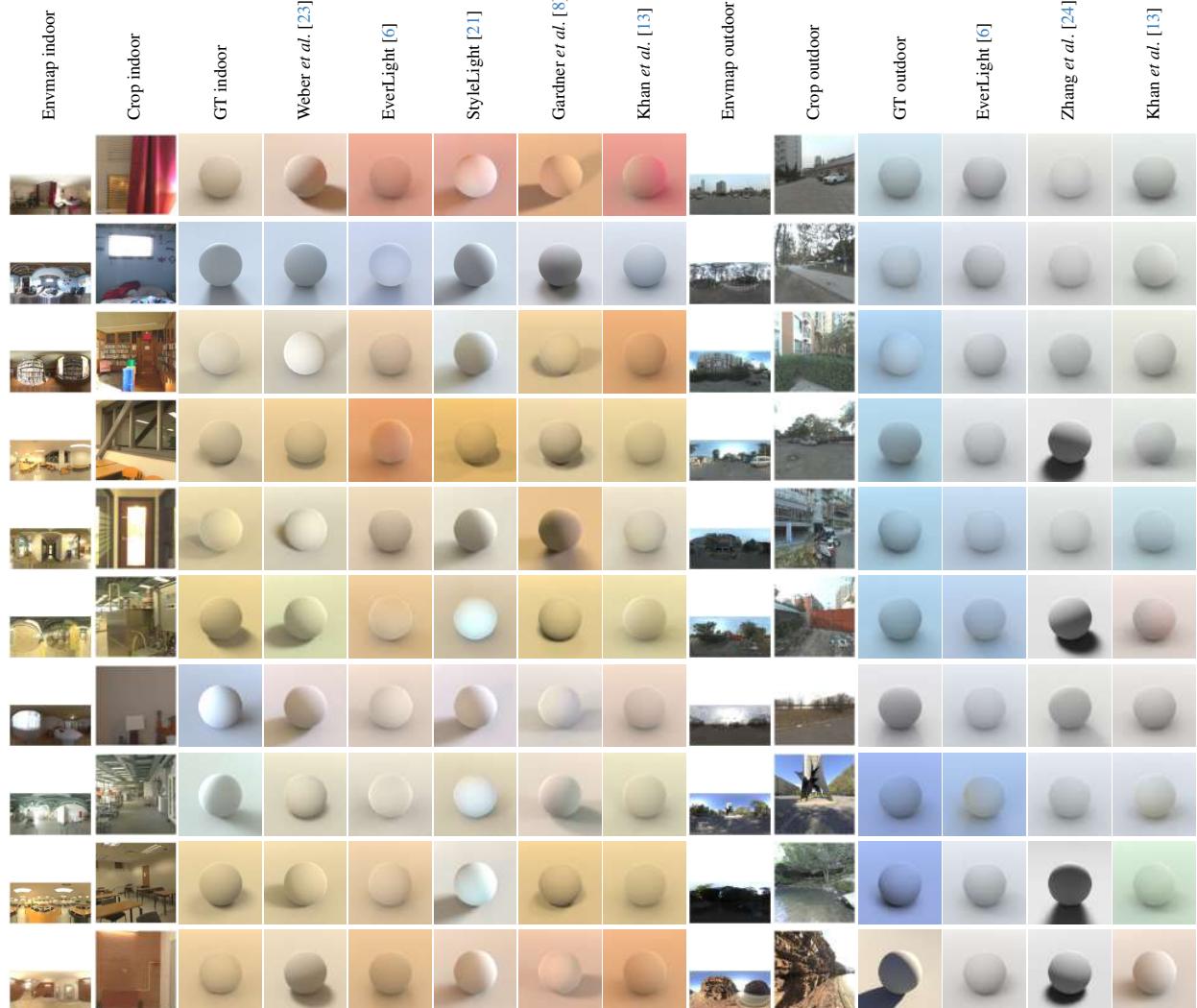


Figure 3. Stimuli used for the task 1 experiment with the diffuse sphere. The full HDR panorama (first and ninth columns) is reexposed and tonemapped with $\gamma = 2.2$ for display, for the indoor and outdoor cases. The region extracted from the scene (second and tenth columns), corresponding to a $50^\circ/90^\circ$ FoV, taken from the centre of the full indoor/outdoor panorama (for most scenes). The rendered stimuli using the ground truth IBLs (first/ninth columns) are shown in the third/eleventh columns. The other columns are the rendered stimuli using the IBLs (shown in fig. 1) produced by the different lighting estimation methods.

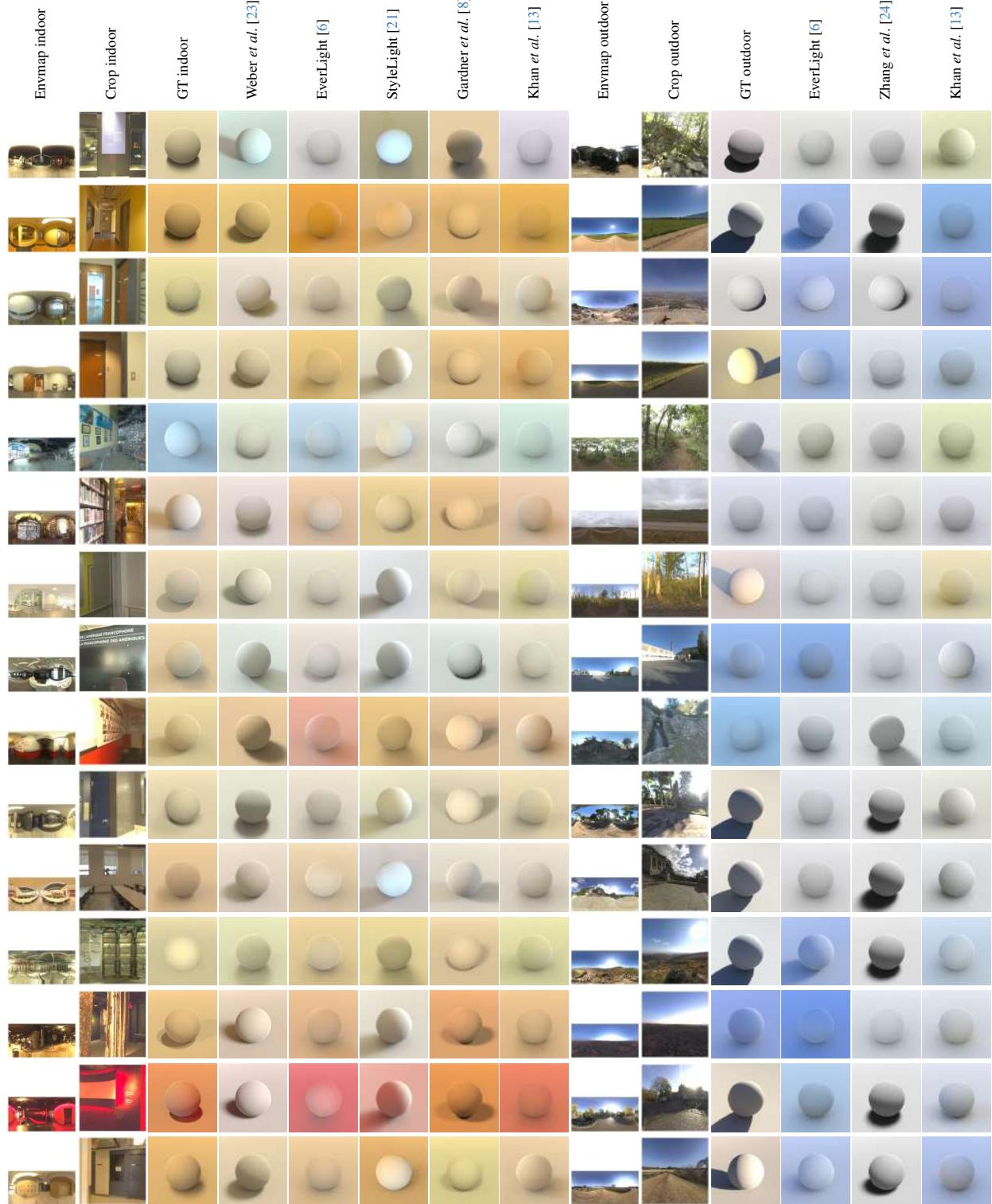


Figure 3. (contd) Stimuli used for the task 1 experiment with the diffuse sphere. The full HDR panorama (first and ninth columns) is reexposed and tonemapped with $\gamma = 2.2$ for display, for the indoor and outdoor cases. The region extracted from the scene (second and tenth columns), corresponding to a $50^\circ/90^\circ$ FoV, taken from the centre of the full indoor/outdoor panorama (for most scenes). The rendered stimuli using the ground truth IBLs (first/ninth columns) are shown in the third/eleventh columns. The other columns are the rendered stimuli using the IBLs (shown in fig. 1) produced by the different lighting estimation methods.

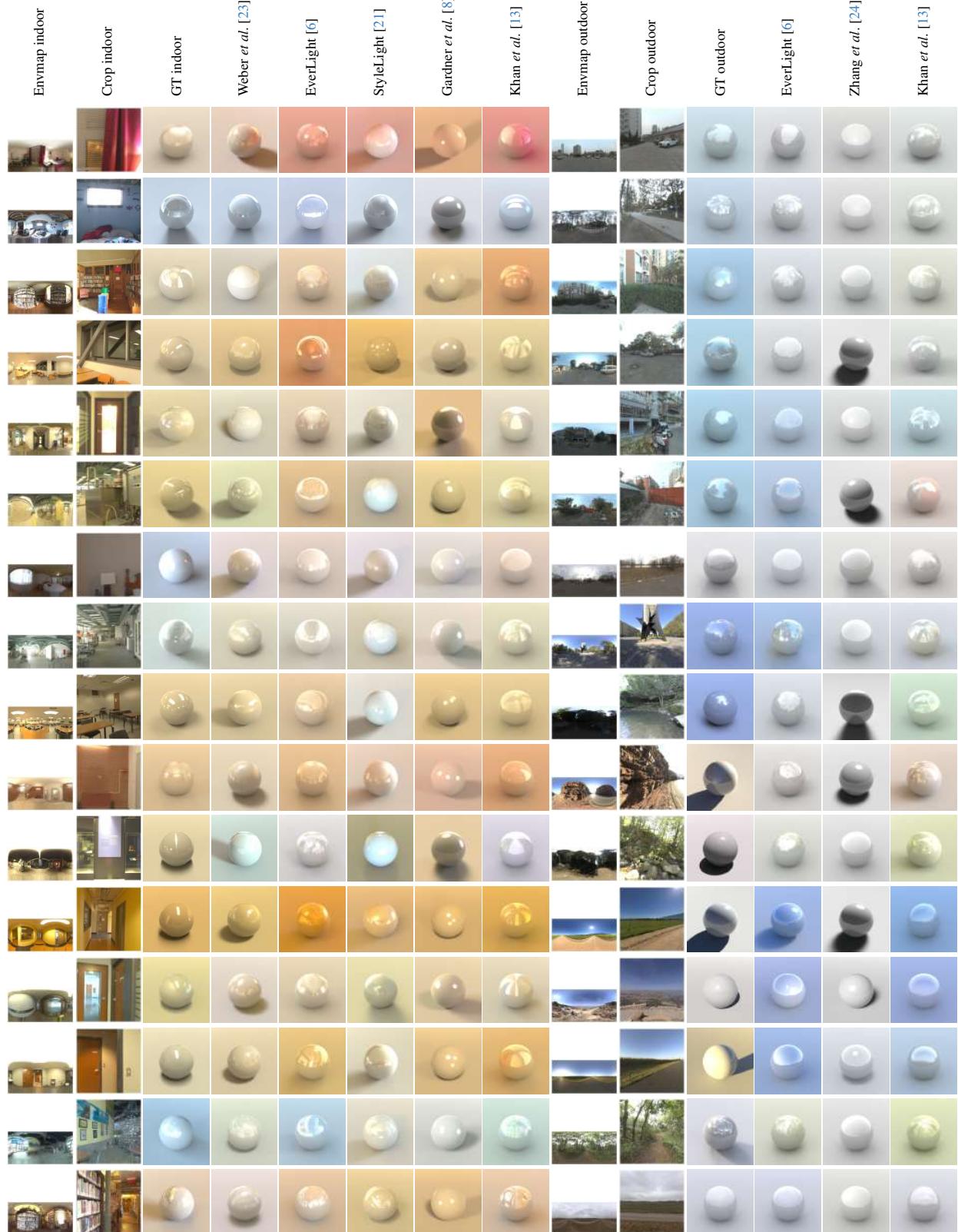


Figure 4. Stimuli used for the task 1 experiment with the glossy sphere. The full HDR panorama (first and ninth columns) is reexposed and tonemapped with $\gamma = 2.2$ for display, for the indoor and outdoor cases. The region extracted from the scene (second and tenth columns), corresponding to a $50^\circ/90^\circ$ FoV, taken from the centre of the full indoor/outdoor panorama (for most scenes). The rendered stimuli using the ground truth IBLs (first/ninth columns) are shown in the third/eleventh columns. The other columns are the rendered stimuli using the IBLs (shown in fig. 1) produced by the different lighting estimation methods.

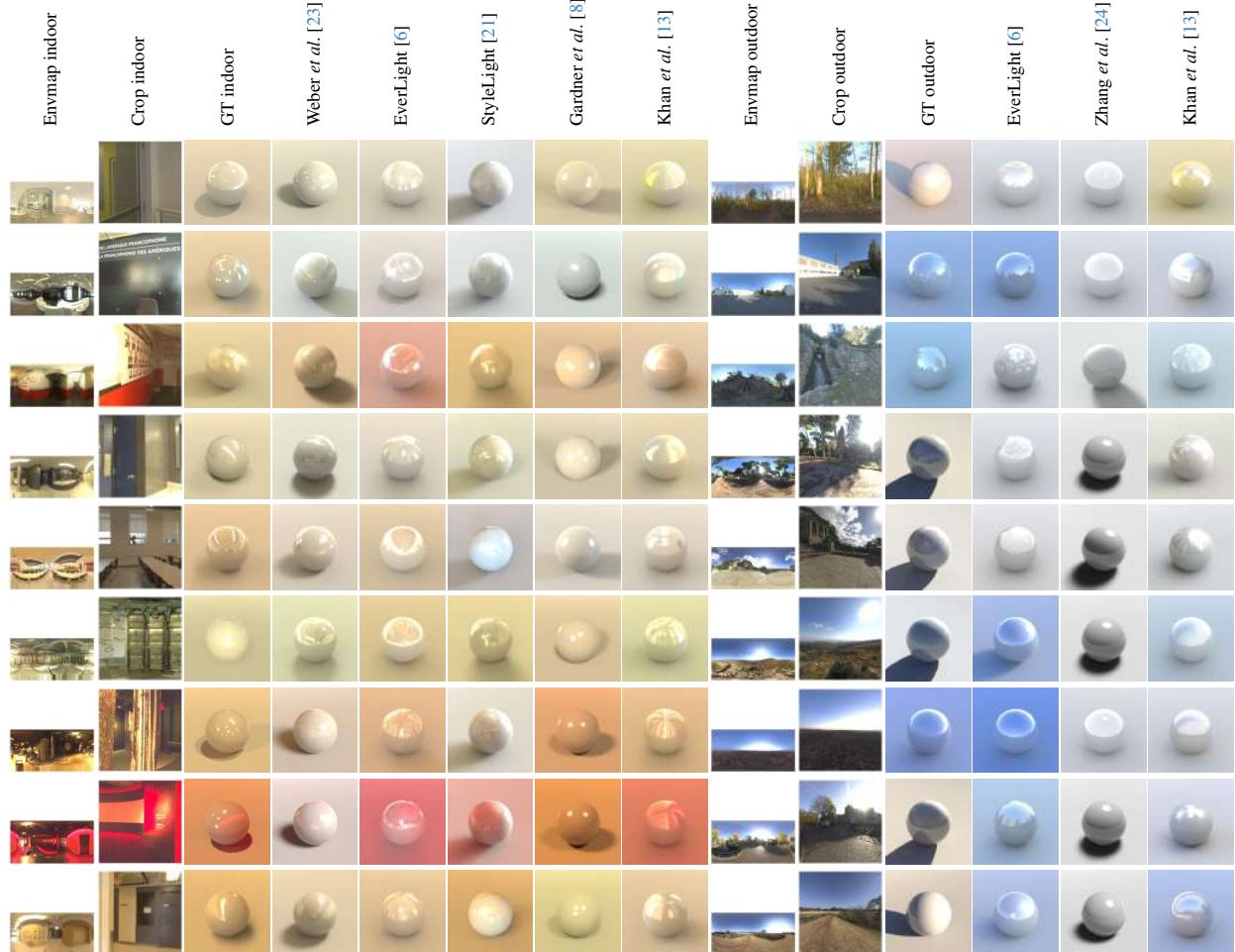


Figure 4. (contd) Stimuli used for the task 1 experiment with the glossy sphere. The full HDR panorama (first and ninth columns) is reexposed and tonemapped with $\gamma = 2.2$ for display, for the indoor and outdoor cases. The region extracted from the scene (second and tenth columns), corresponding to a $50^\circ/90^\circ$ FoV, taken from the centre of the full indoor/outdoor panorama (for most scenes). The rendered stimuli using the ground truth IBLs (first/ninth columns) are shown in the third/eleventh columns. The other columns are the rendered stimuli using the IBLs (shown in fig. 1) produced by the different lighting estimation methods.

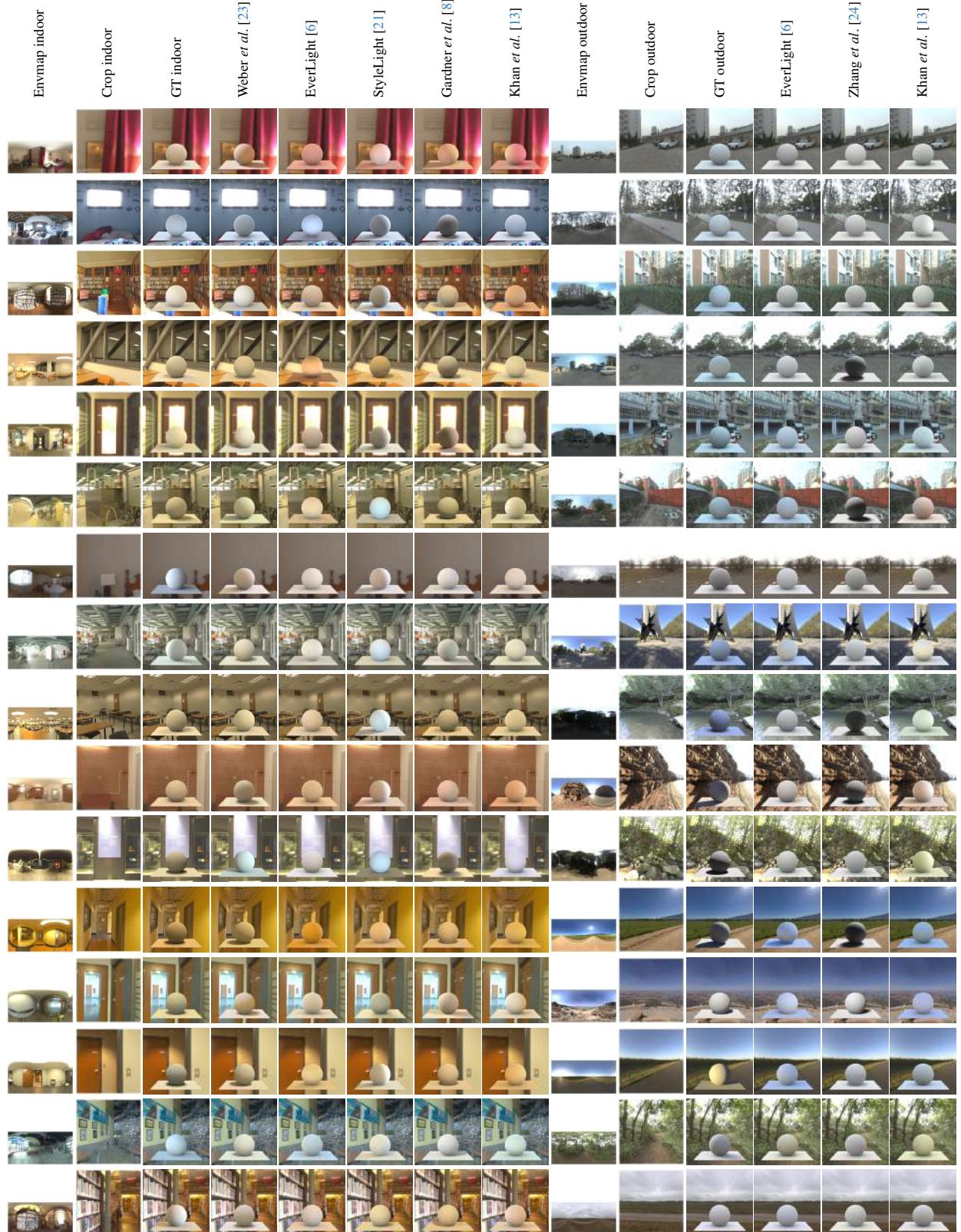


Figure 5. Stimuli used for the task 2 experiment with the diffuse sphere. The full HDR panorama (first and ninth columns) is reexposed and tonemapped with $\gamma = 2.2$ for display, for the indoor and outdoor cases. The region extracted from the scene (second and tenth columns), corresponding to a $50^\circ/90^\circ$ FoV, taken from the centre of the full indoor/outdoor panorama (for most scenes). The rendered stimuli using the ground truth IBLs (first/ninth columns) are shown in the third/eleventh columns. The other columns are the rendered stimuli using the IBLs (shown in fig. 1) produced by the different lighting estimation methods.

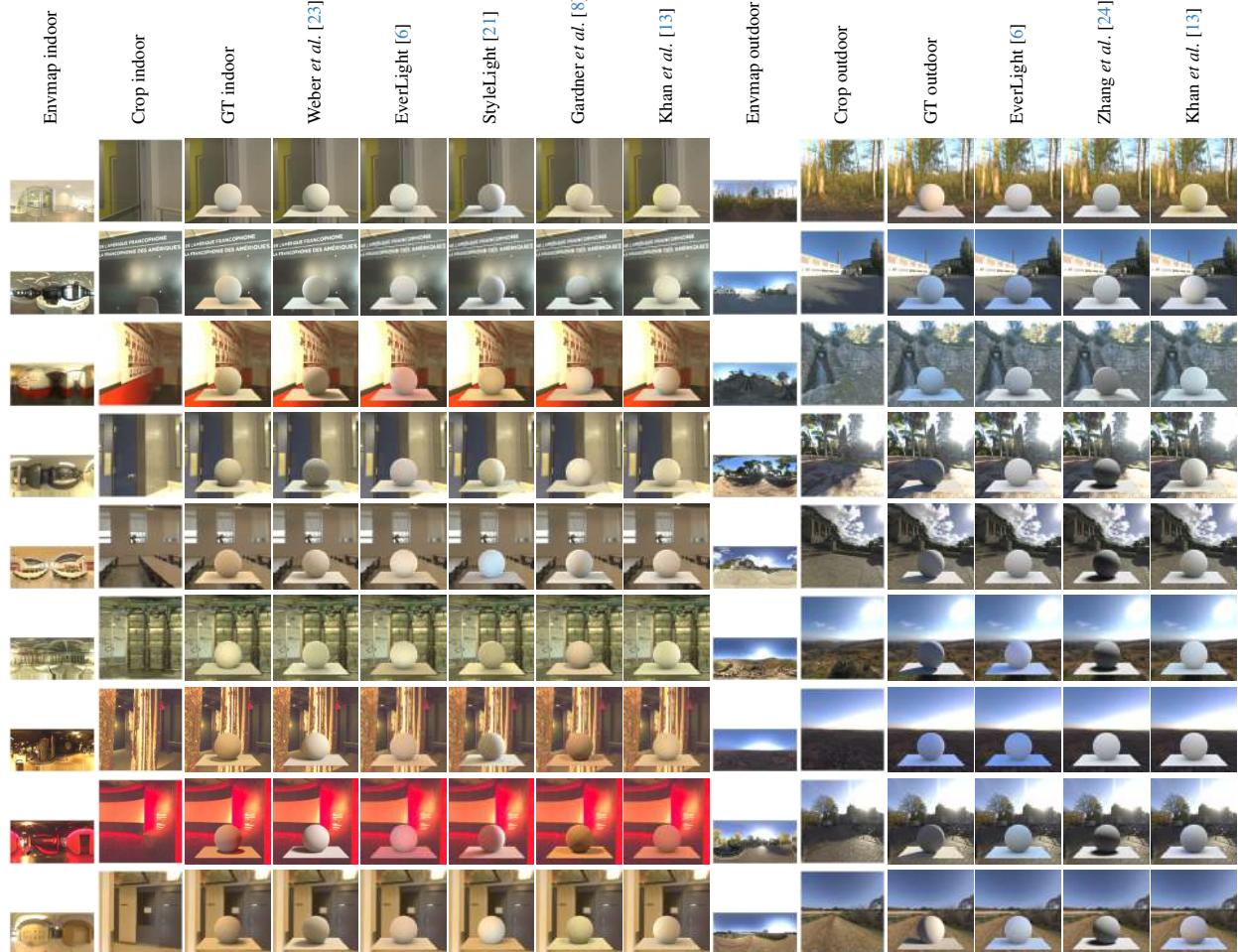


Figure 5. (contd) Stimuli used for the task 2 experiment with the diffuse sphere. The full HDR panorama (first and ninth columns) is reexposed and tonemapped with $\gamma = 2.2$ for display, for the indoor and outdoor cases. The region extracted from the scene (second and tenth columns), corresponding to a $50^\circ/90^\circ$ FoV, taken from the centre of the full indoor/outdoor panorama (for most scenes). The rendered stimuli using the ground truth IBLs (first/ninth columns) are shown in the third/eleventh columns. The other columns are the rendered stimuli using the IBLs (shown in fig. 1) produced by the different lighting estimation methods.

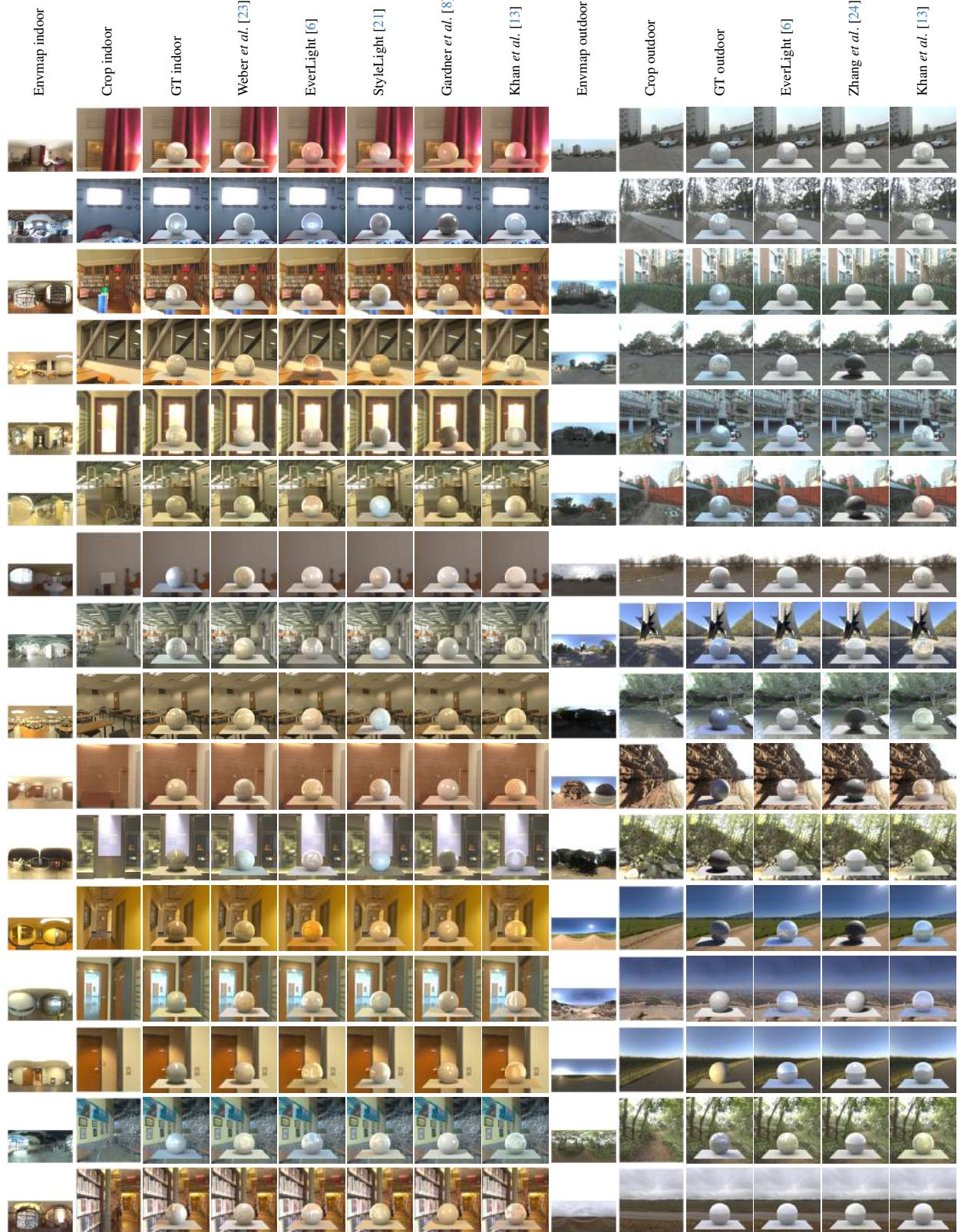


Figure 6. Stimuli used for the task 2 experiment with the glossy sphere. The full HDR panorama (first and ninth columns) is reexposed and tonemapped with $\gamma = 2.2$ for display, for the indoor and outdoor cases. The region extracted from the scene (second and tenth columns), corresponding to a $50^\circ/90^\circ$ FoV, taken from the centre of the full indoor/outdoor panorama (for most scenes). The rendered stimuli using the ground truth IBLs (first/ninth columns) are shown in the third/eleventh columns. The other columns are the rendered stimuli using the IBLs (shown in fig. 1) produced by the different lighting estimation methods.

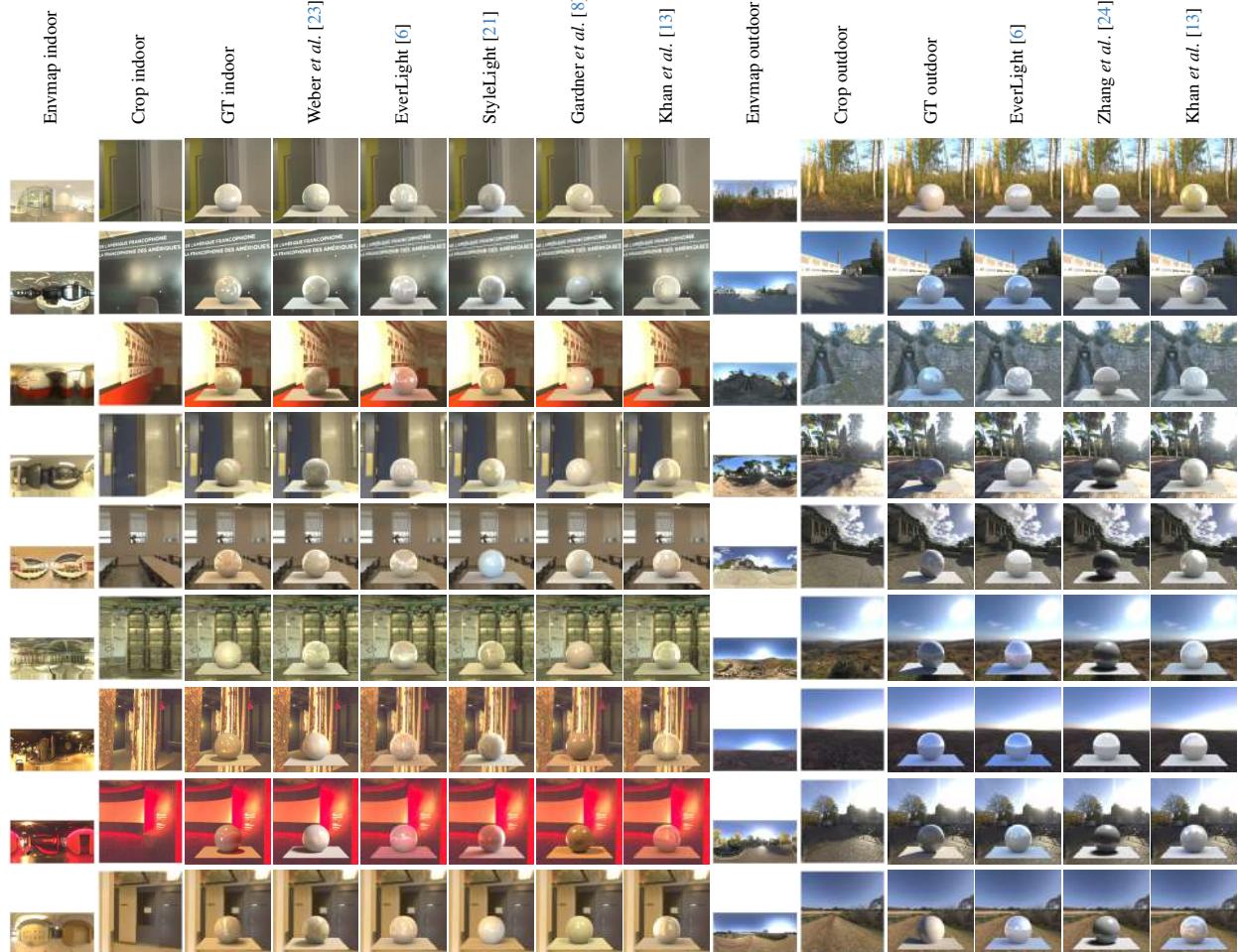


Figure 6. (contd) Stimuli used for the task 2 experiment with the glossy sphere. The full HDR panorama (first and ninth columns) is reexposed and tonemapped with $\gamma = 2.2$ for display, for the indoor and outdoor cases. The region extracted from the scene (second and tenth columns), corresponding to a $50^\circ/90^\circ$ FoV, taken from the centre of the full indoor/outdoor panorama (for most scenes). The rendered stimuli using the ground truth IBLs (first/ninth columns) are shown in the third/eleventh columns. The other columns are the rendered stimuli using the IBLs (shown in fig. 1) produced by the different lighting estimation methods.

141 **Geometry.** For both tasks, the stimulus corresponds to a
142 sphere, with a radius of 1.5 m, on a plane, to act as a shadow
143 catcher. In task 1, the plane is 30 m × 30 m to cover the entire FoV and in task 2, the plane is 2.5 m × 3 m to allow
144 the composited background to be seen. The background
145 image used in task 2 corresponds to the extracted region
146 from scenes given as input to the indoor lighting estimation
147 methods (sec. 1.2), to give the virtual objects context.
148

149 A virtual camera is positioned parallel to the $x - y$ plane,
150 facing the $+x$ axis, thus capturing in its FoV the $y - z$ plane,
151 where y points towards the left and z upwards. The virtual
152 camera is raised by 1.6 m from the perpendicular axis of the
153 horizontal plane (z) with regard to the origin, to simulate the
154 standard height of humans. In task 1, the FoV is inclined
155 by 30° with regards to the horizon, to include the shadows
156 produced by the sphere on the plane and for task 2, the FoV
157 is inclined by 90°, to align realistically with the composited
158 background image.

159 **Materials.** For both tasks, two separate experiments are
160 done on spheres with two different materials (diffuse and
161 glossy), which use the Disney Principled BRDF [2]. The
162 diffuse material has a roughness of 1.0 and a specularity
163 of 0.0, whilst the glossy material has a roughness of 0.1
164 and a specularity of 1.0. The Lambertian sphere allows the
165 observers to evaluate the lower frequency light cues, such as
166 the colour, intensity, and degree of collimation. The opaque
167 glossy sphere includes the high frequency light cues and the
168 texture to be judged by the observer. For both versions, the
169 plane maintains a grey Lambertian material with the same
170 parameters as the diffuse sphere. All the objects have an
171 albedo of 0.18.

172 **Rendering.** The synthetic objects are rendered using
173 the physically based rendering engine Cycles in
174 Blender [4]. The rendered stimuli have a resolution
175 of 256 px × 256 px, as the extracted regions from scenes
176 given as input to the indoor lighting estimation methods
177 (sec. 1.2). The renders are saved in the `exr` format and then
178 tonemapped with $\gamma = 2.2$ and reexposed to be displayed on
179 the monitor used during the experiment (sec. 1.4).

180 1.4. Hardware

181 The experiment (see sec 3.3. of the main paper) is conducted
182 in a controlled lab setting to ensure the data collected is
183 uniform. The experiment is carried out in a matte black
184 room (painted walls and ceilings, with black rug flooring)
185 with a standard keyboard placed on a desk and the monitor
186 set to sRGB. The monitor was the only light source. The
187 observers are seated at ~ 70 cm from the monitor, which
188 gives a 11.5°/17° visual angle for task 1/2. The experimental
189 setup is shown in 7.



Figure 7. Photograph of the experimental setup of the psychophysical experiment.

The experiment runs on MATLAB (version R2023a) and uses the Psychophysics Toolbox. An example of the screen displayed to the observers is shown in 8 for all four experiments.

190 1.5. Procedure

During the experiment (see sec 3.3. of the main paper), the images are selected using the arrows on the keyboard. The background is middle grey.

Observers are asked to participate in two of the four experiments (task 1 and task 2), with randomly assigned material for the first task and the opposite material for the second task, to avoid potentially causing bias. A break is offered between the experiments to avoid fatigue. Each experiment takes 10–35 min to complete. No time restriction is imposed on the observers to avoid inducing stress and bias. The observers are advised to follow their intuition to determine their preference and that each combination of stimuli shown should be analysed in around than 5 s, so the experiment would not last too long. This is done to avoid the fatigue or boredom they experience when doing the task for too long.

At the beginning of each experiment, a short tutorial is shown to the observer with an example not included in the dataset. The observers are informed that the images always contain the same sphere (same geometry) made of the same material, for all the stimuli they see during that specific experiment, and that only the lighting has changed. They are also informed that there is no right answer, and that we are only trying to measure their preferences. The participants are unaware that different lighting estimation methods have been used to produce the stimuli. To confirm that the observers are not colourblind, an Ishihara test is conducted for each participant before starting the experiments.

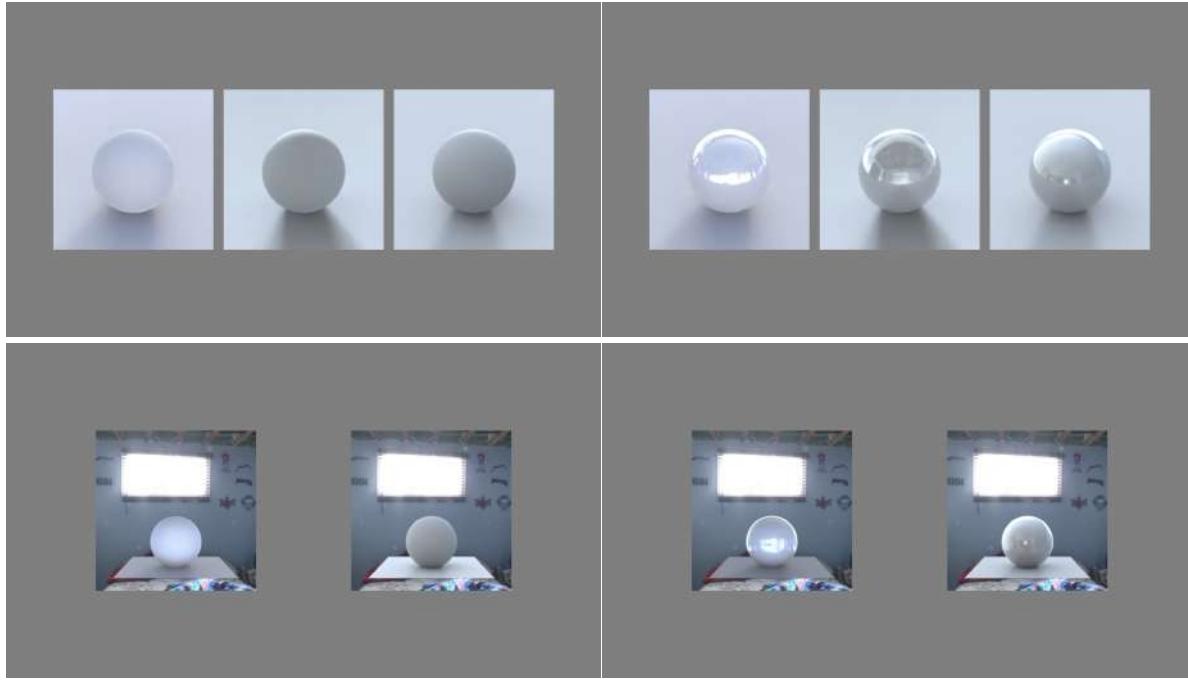


Figure 8. Examples of the stimuli presented during the experiments. Top row are the stimuli fr the comparison experiments for the diffuse (left) and glossy (right) spheres and the bottom row are the stimuli for the realism experiments.

1.6. Participants

A total of 49 unique observers (33M/16F) participated in the study (see sec 3.3. of the main paper). Some observers were asked to participate twice (in different sessions) and were assigned the two remaining experiments they had not previously done. 12 observers participated in all four outdoor experiments.

Fig. 9 shows the temporal evolution of the score based on the number of participants, to validate the convergence of the preferred methods by the observers. The curves show little variance as more participants are added after ~ 15 for the indoor experiment (left) and ~ 8 for the outdoor experiment (right), which confirms the number of participants included in the study is sufficient to describe a general tendency.

None of the participants were authors. 11 observers were students from the computer vision department (labelled *experts*), who were unaware of the project.

The scores obtained for the indoor lighting estimation methods for the expert (yellow) and naive (teal) observers are shown in fig. 10 (the procedure for computing score is described in sec. 4.1. of the main paper). It is possible to see that the scores for all the experiments are similar, which confirms that both samples of observers have the same trends and do not use different light cues in the stimuli.

The agreement score (described in sec. 5.1. of the main paper) between the expert and naive observers and the metrics, shown in fig. 11, also displays the same trends between

each group, which further confirms that their behaviour is similar.

2. Psychophysical results

Additional analysis of the psychophysical results (see sec 4.2. of the main paper) is done in this section. An example of the observers' ranking of a set of stimuli for an input scene is shown in sec. 2.1, and the trends of the preferred methods per image is shown in sec. 2.2. The agreement score of the individual observers is discussed in sec. 2.3

2.1. Example of stimuli ranking

An example of each indoor lighting estimation method's ranking (in decreasing order) and the associated score for each stimulus for the same input scene is shown in fig. 12, for each experiment. When comparing against the ground truth stimulus (task 1) for the diffuse sphere (first row), observers seem to agree—at least in essence—to what IQA metrics are trying to achieve: having an image as close as possible to the ground truth reference. E.g. the lighting estimation of Gardner *et al.* [8] does not accurately match the ground truth, while the one produced by Weber *et al.* [23] resembles the ground truth very closely. Yet, when the resulting lighting estimations are put into context (task 2, third row), lighting accuracy based on the ground truth does not seem to matter as much to be considered plausible. In this case, the preferred estimated lighting does not seem to match well the ground

222

249

223

250

224

251

225

252

226

253

227

254

228

255

229

256

230

257

231

258

232

259

233

260

234

261

235

262

236

263

237

264

238

265

239

266

240

267

241

268

242

269

243

270

244

271

245

272

246

273

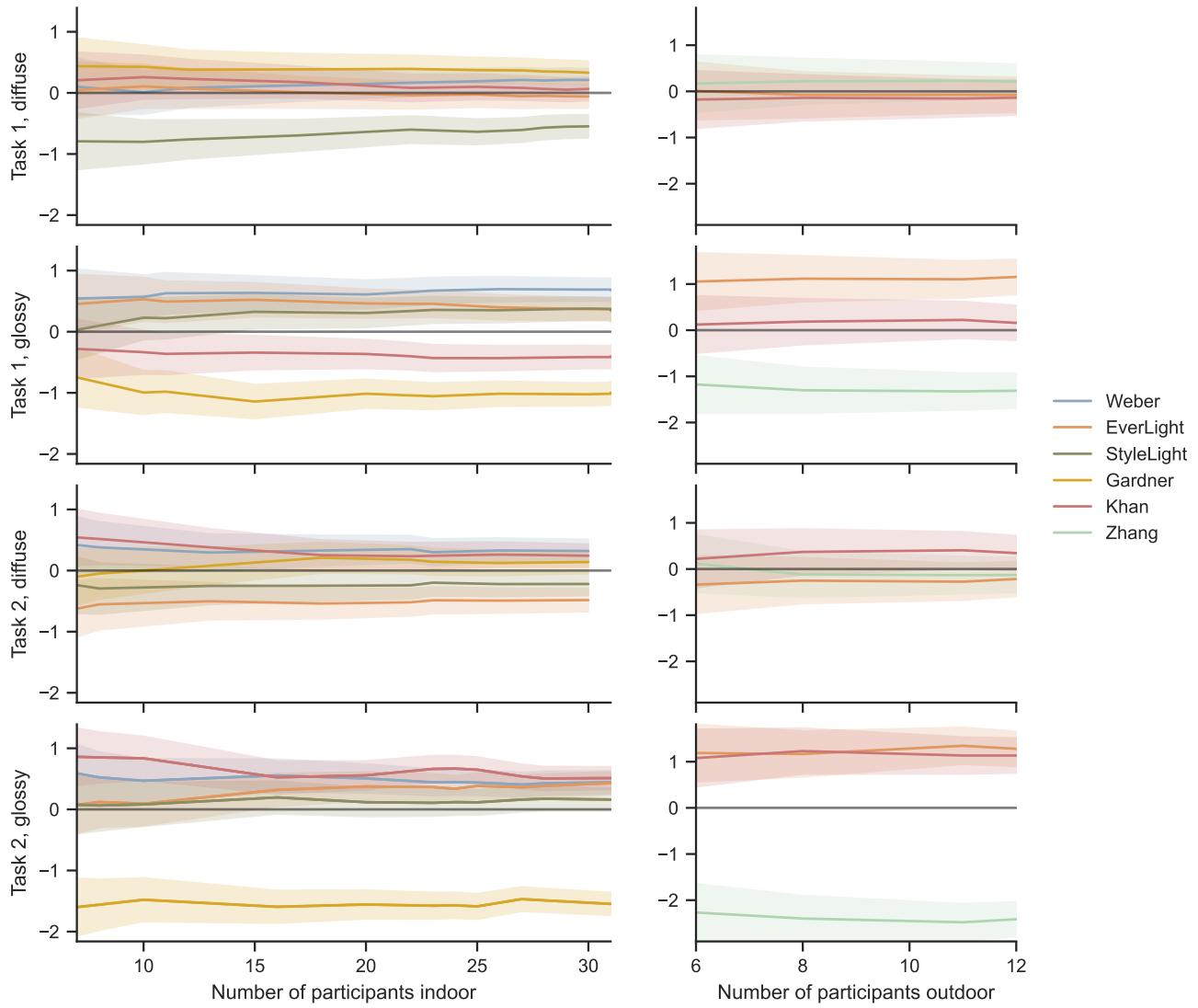


Figure 9. Convergence of the score for the different lighting estimation methods as a function of the number of participants. The line uncertainties correspond to 95 % confidence interval.

274 truth.

275 Fig. 12 also illustrates an example of the trend observed in
276 fig. 4 of the main paper for the glossy experiments produced
277 by Gardner *et al.* [8]. The observers seem to consider texture
278 primordial when observing a more reflective surface. They
279 seemingly consider it more important than having a plausible
280 lighting intensity and direction, when judging the plausibility
281 of an inserted object.

282 2.2. Scores per image

283 The scores from all the observers for each stimulus are shown
284 in fig. 13. This figure clearly shows that glossy stimuli from
285 Gardner *et al.* [8] (indoors) and Zhang *et al.* [24] (outdoor)
286 are, in general, rarely picked. For the other lighting estima-

287 tion methods, we can see a greater variance in the scores,
288 with some images performing very well or very poorly. This
289 indicates some of the limitations of a given method in some
290 specific case.

291 2.3. Individual observer agreement

292 The individual observer agreement scores $\omega^{(i)}$ are shown
293 in fig. 14, for the indoor (left) and outdoor (right) lighting
294 estimation methods, for all experiments (rows). The ob-
295 servers are anonymised by assigning them a random number
296 and a random order between each experiment (i.e. the ob-
297 server labelled “1” in the first row is not necessarily the same
298 observer “1” in the second row), as not all the observers par-
299 participated in the same experiments. The observers labelled as

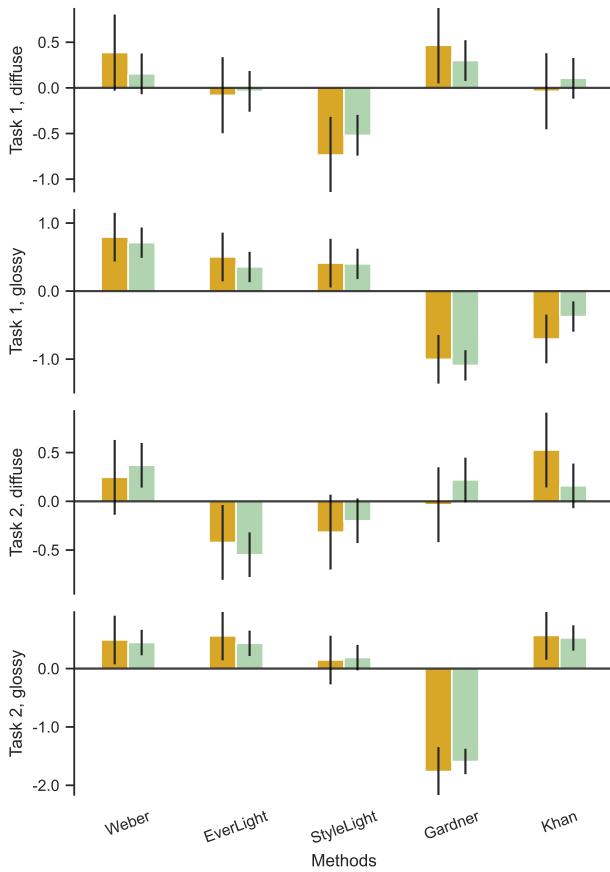


Figure 10. Thurstone Case V Law of Comparative Judgement scores from the expert (yellow) and naive (teal) observers as a function of the different indoor lighting estimation methods (columns), for the different types of sphere materials used in the experiments (rows). Error bars correspond to 95 % confidence interval.

“P” and “R” correspond to the perfect and random observers (defined in sec. 5.1. of the main paper), respectively. This score is how the observers removed from the study are determined (see sec. 5.1. of the main paper for more details). The orange/blue lines are determined by taking the average of the individual observer agreement scores $\omega^{(i)}$ (excluding the perfect and random observers), to obtain the expected observer agreement score (defined in sec. 5.1. of the main paper).

3. Measuring the scores of the IQA metrics

Additional information regarding the scores of the IQA metrics (see sec 5. of the main paper) is provided in this section. In this subsection we aim to correlate in a simpler manner the results of our metrics with our psychophysical experiment. We will consider the same metrics as in the main paper RGB angular error [10], PSNR [11], RMSE, si-RMSE and more recent ones, such as SSIM [22], VIF [19], (LPIPS [25],

PieAPP [16], FLIP [1], and HyperIQA [20].

Broadly speaking, our idea is to consider the metrics as if they were observers in our pair comparison test. In order to do so, we will perform all the same pair-wise comparisons of our experiment, and for each comparison, we will assign a 1 to the image that the metric picks as better and a 0 to the other image. Once we have all the selections, we can compute the scores for each of the metrics. More in detail, the score is then computed for each metric, for each experiment, the same way as described in sec. 4.1. of the main paper. Let us note —as we already did in the main paper— that the results of the metrics for the stimuli in tasks 1 and 2 are the same.

The Thurstone Case V Law of Comparative Judgement scores obtained for each metric computed on the stimuli from the indoor (left) and outdoor (right) lighting estimation methods are shown in fig. 15, for all the experiments. This figure clearly highlights that the different metrics do not agree with each other for a given set of images. This fact reinforces our main arguments in two different ways. Firstly, it shows that current metrics are not feasible enough as by selecting a specific metric the ranking of the methods is also modified. Secondly, it also proves that our approach of considering an ensemble of different metrics to derive our new metric is the avenue to pursue.

4. Learning a metric combination

Additional details on the training of our proposed metrics (see sec 6.1. of the main paper) are given in sec. 4.1, and a supplementary analysis of the generalisation psychophysical study (see sec 6.2. of the main paper) is done in sec. 4.2.

4.1. Formulation and training

In tab. 1, the mean and standard deviation of the training accuracies for the k -fold approach ($k = 10$) are displayed, for various classical learners. The parameters are all the default values. In some cases, *BernoulliNB* performs better than the SVR, however that methods outputs discrete values, and not a continuous ranking like the SVR, thus this method is preferred.

A representation of the learnt metric by the SVR is shown in fig. 16 for the indoor and outdoor validation cases (shown in “Ours” column of fig. 5 in the main paper) and for the holdout and new methods generalisation tests in fig. 17 (which corresponds to the “Ours Holdout” column of fig. 5 in the main paper and sec. 6.2. of the main paper, respectively). The projection of the \mathbb{R}^{10} space of the classical metrics is projected to \mathbb{R}^2 to illustrate the decision boundaries (white lines) of the learnt metrics. The data points with the black contours indicate data points used in validation and the grey contours correspond to the support vectors.

317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340

341
342
343
344
345

346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364

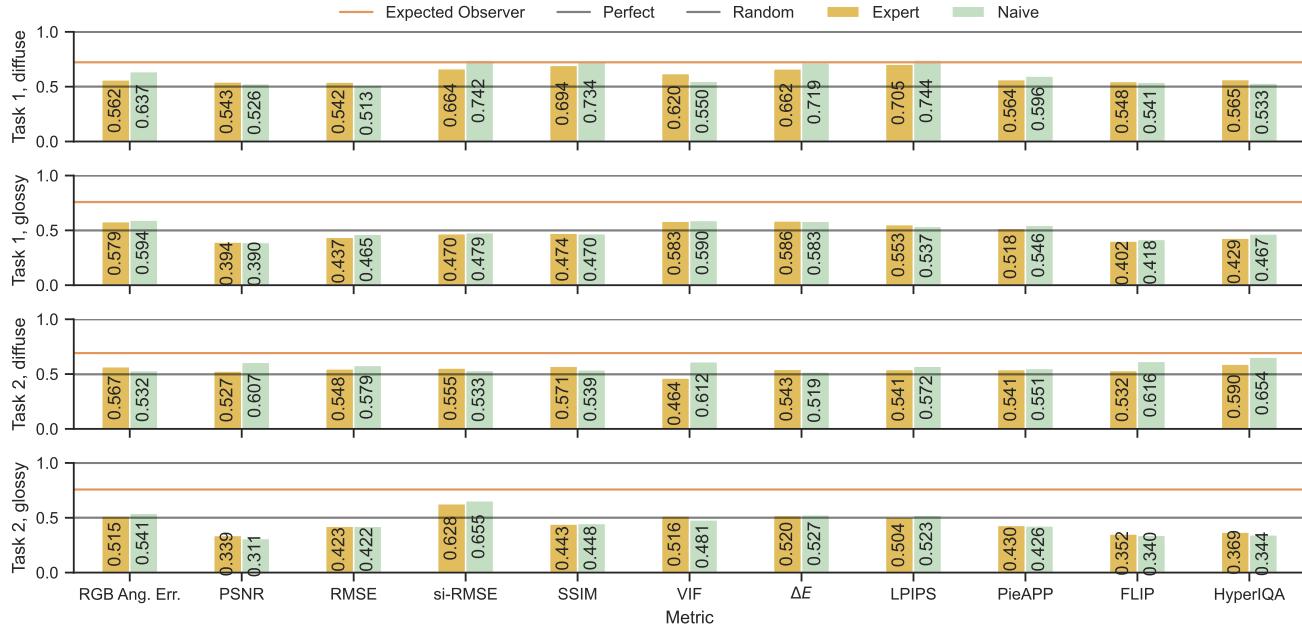


Figure 11. Agreement between the observer scores (expert: yellow left bars; naive: teal right bars) and the metric scores (columns) for all the indoor lighting estimation methods, for the different types of experiments (rows). The lower horizontal grey bar is set at chance level (~ 0.5) and the higher one corresponds to the perfect observer (set at 1.0). The orange line corresponds to the expected observer agreement score for all the observers for the indoor methods (same as the one in fig. 5 of the main paper).

4.2. Generalisation to other lighting estimation methods

All the environment maps generated by the generalisation lighting estimation methods for the user study are presented in fig. 18.

The stimuli used in the indoor and outdoor user study, for the diffuse/glossy experiments, are shown in fig. 19/fig. 20 for task 1 and fig. 21/fig. 22 for task 2.

The scores for all observers for the generalisation lighting estimation methods are shown in fig. 23. Only Weber *et al.* [23] is a method also included in the main psychophysical study; however, it has never been compared to the new methods. The methods are described in sec. 6.2. of the main paper.

Unlike Khan *et al.* [13], which includes texture but lacks a proper HDR lighting estimation, the Average method only includes the average pixel colour of the background image given as input (sec. 1.2). Thus, it includes no texture nor lighting estimation. The fact that for all experiments this method is not preferred demonstrates that excluding texture and lighting estimation cannot produce accurate (task 1) nor realistic (task 2) results. We hypothesize that this lack of texture causes the method to be disliked for the glossy experiments, as it has been empirically shown in the main study with the results from fig. 4 in the main paper for Gardner *et al.* [8] and Zhang *et al.* [24]. However, the naive

approach to lighting estimation from the Average method is similar to Khan *et al.* [13], yet the Average method does not perform well in diffuse spheres, especially on task 2, like Khan *et al.* [13]. When comparing the stimuli from Khan *et al.* [13] (fig. 5) and the Average method (fig. 21), it is possible to see that there are low-frequency lighting variations on the renders done with Khan *et al.* [13] (e.g. spatial colour variations), which we believe contribute to the judgement of plausibility by the observers.

For the same reasons as Gardner *et al.* [8] and Zhang *et al.* [24] in fig. 4 of the main paper, Garon *et al.* [9] seems not to be preferred by observers as it lacks textures, adding empirical evidence that textures are an important part of lighting estimation, especially for the glossy experiments.

However, it is important to remember that this generalisation study has been done on a small number of participants and fewer images, yielding an increased statistical uncertainty. Nevertheless, fig. 9 shows that even with 6 observers, the general trends do not seem to change much.

References

- [1] Pontus Andersson, Jim Nilsson, Tomas Akenine-Möller, Magnus Oskarsson, Kalle Åström, and Mark D. Fairchild. FLIP: a difference evaluator for alternating images. *ACM Comp. Graph. Int. Tech.*, 3(2):15:1–15:23, 2020. 18
- [2] Brent Burley. Physically-based shading at disney. 2012. 15

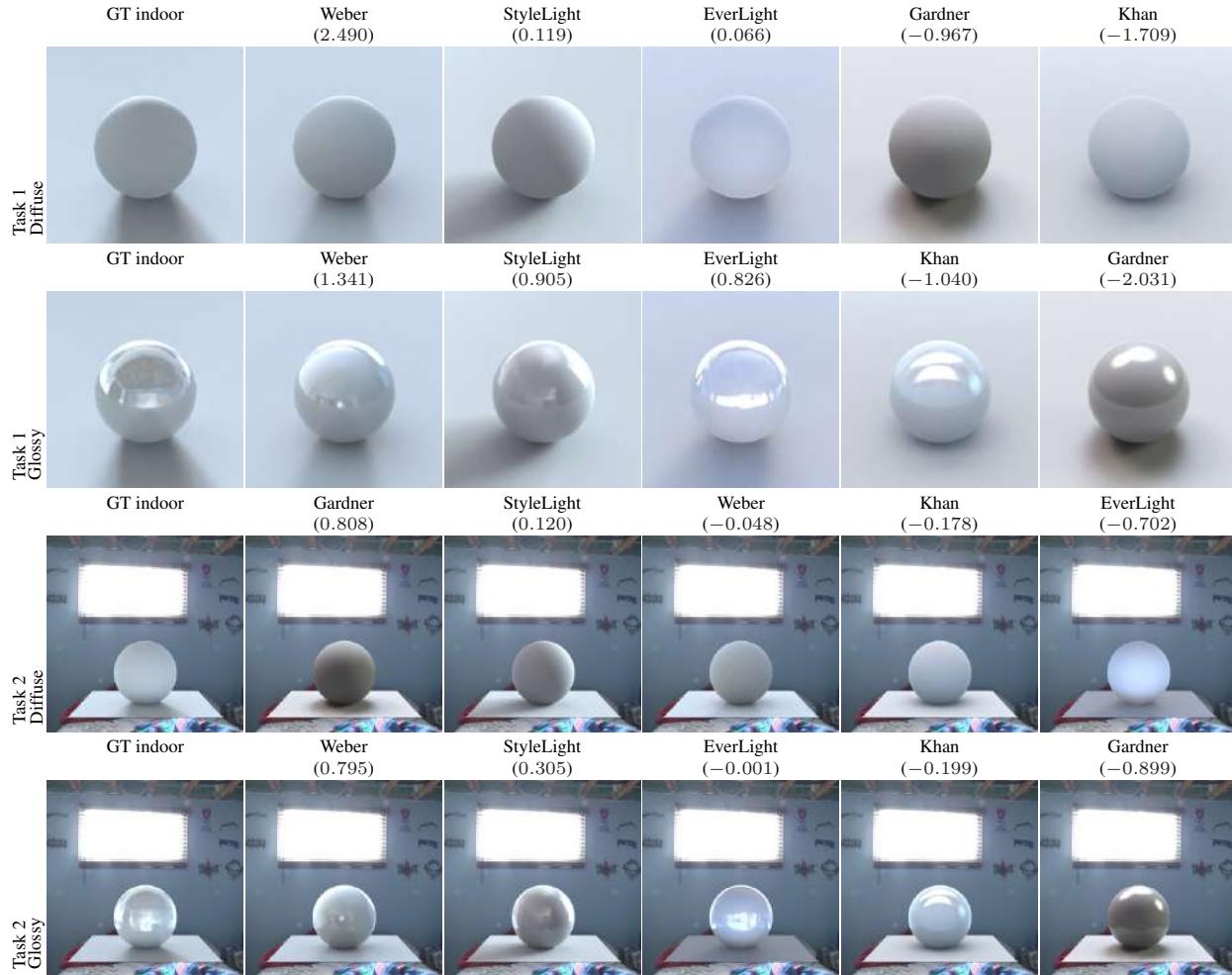


Figure 12. Example of the preferences of all the observers of the stimuli generated by the different lighting estimation methods (columns) ordered decreasingly as a function of their score, for the different types of experiments (rows). The first column is the ground truth (not judged by the observers) associated to the scene for comparison.

- 416 [3] Dachuan Cheng, Jian Shi, Yanyun Chen, Xiaoming Deng, and
- 417 Xiaopeng. Zhang. Learning scene illumination by pairwise
- 418 photos from rear and front mobile cameras. *Comput. Graph.*
- 419 *Forum*, 37(7):213–221, 2018. 5
- 420 [4] Blender Online Community. *Blender - a 3D modelling and*
- 421 *rendering package*. Blender Foundation, Stichting Blender
- 422 Foundation, Amsterdam, 2018. 15
- 423 [5] M. Dastjerdi, Y. Hold-Geoffroy, J. Eisenmann, S. Kho-
- 424 dadadeh, and J. Lalonde. Guided co-modulated GAN for
- 425 360° field of view extrapolation. In *Int. Conf. 3D Vis.*, 2022.
- 426 1
- 427 [6] Mohammad Reza Karimi Dastjerdi, Jonathan Eisenmann,
- 428 Yannick Hold-Geoffroy, and Jean-François Lalonde. Ever-
- 429 light: Indoor-outdoor editable HDR lighting estimation.
- 430 In *Int. Conf. Comput. Vis.*, 2023. 1, 3, 4, 5, 7, 8, 9, 10, 11, 12,
- 431 13, 14, 27, 28, 29, 30, 31
- 432 [7] Marc-André Gardner, Kalyan Sunkavalli, Ersin Yumer, Xiao-
- 433 hui Shen, Emiliano Gambaretto, Christian Gagné, and Jean-

- François Lalonde. Learning to predict indoor illumination from a single image. *ACM Trans. Graph.*, 9(4), 2017. 5 434
- [8] Marc-André Gardner, Yannick Hold-Geoffroy, Kalyan Sunkavalli, Christian Gagné, and Jean-François Lalonde. Deep parametric indoor lighting estimation. In *Int. Conf. Comput. Vis.*, 2019. 1, 3, 4, 7, 8, 9, 10, 11, 12, 13, 14, 16, 17, 19 436
- [9] Mathieu Garon, Kalyan Sunkavalli, Sunil Hadap, Nathan Carr, and Jean-François Lalonde. Fast spatially-varying indoor lighting estimation. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2019. 19, 27, 28, 29, 30, 31 441
- [10] Arjan Gijsenij, Theo Gevers, and Marcel P. Lucassen. Perceptual analysis of distance measures for color constancy algorithms. *J. Opt. Soc. Am. A*, 26(10):2243–2256, 2009. 18 445
- [11] Alain Horé and Djemel Ziou. Image quality metrics: PSNR vs. SSIM. In *Int. Conf. Pattern Recog.*, 2010. 18 446
- [12] Tero Karras, Samuli Laine, and Timo Aila. A style-based 447

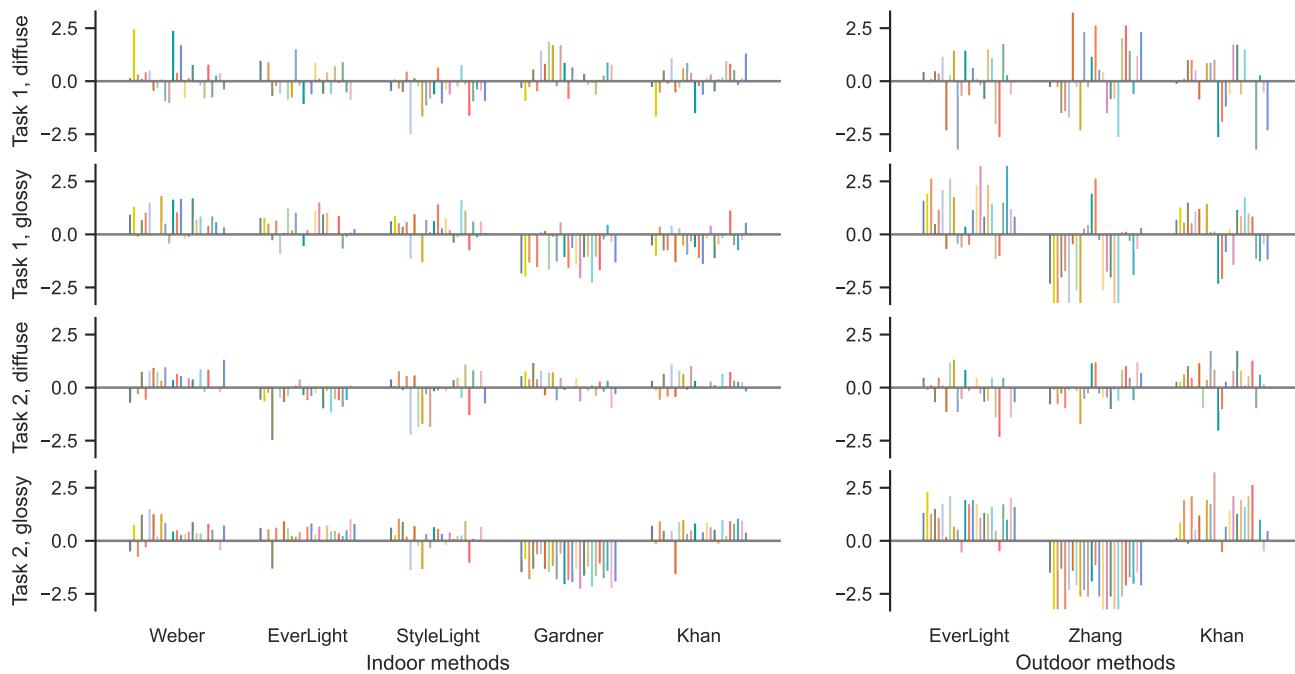


Figure 13. Thurstone Case V Law of Comparative Judgement scores from all the observers for each of the sets of images as a function of the different lighting estimation methods (columns), for the different types of sphere materials used in the experiments (rows).

- generator architecture for generative adversarial networks. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2019. 1
- [13] Erum Arif Khan, Erik Reinhard, Roland W. Fleming, and Heinrich H. Bülthoff. Image-based material editing. *ACM Trans. Graph.*, 25(3):654–663, 2006. 1, 3, 4, 7, 8, 9, 10, 11, 12, 13, 14, 19
- [14] Jean-François Lalonde and Iain Matthews. Lighting estimation in outdoor image collections. In *Int. Conf. 3D Vis.*, 2014. 1
- [15] L. McInnes, J. Healy, and J. Melville. UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction. *ArXiv e-prints*, 2018. 6, 25, 26
- [16] Ekta Prashnani, Hong Cai, Yasamin Mostofi, and Pradeep Sen. Pieapp: Perceptual image-error assessment through pairwise preference. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2018. 18
- [17] Daniel Roich, Ron Mokady, Amit H Bermano, and Daniel Cohen-Or. Pivotal tuning for latent-based editing of real images. *ACM Trans. Graph.*, 2021. 1
- [18] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2022. 27, 28, 29, 30, 31
- [19] Hamid R. Sheikh, Alan C. Bovik, and Gustavo de Veciana. An information fidelity criterion for image quality assessment using natural scene statistics. *IEEE Trans. Image Process.*, 14(12):2117–2128, 2005. 18
- [20] Shaolin Su, Qingsen Yan, Yu Zhu, Cheng Zhang, Xin Ge, Jinqiu Sun, and Yanning Zhang. Blindly assess image quality

- in the wild guided by a self-adaptive hyper network. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2020. 18
- [21] Guangcong Wang, Yinuo Yang, Chen Change Loy, and Ziwei Liu. Stylelight: HDR panorama generation for lighting estimation and editing. In *Eur. Conf. Comput. Vis.*, 2022. 1, 3, 4, 7, 8, 9, 10, 11, 12, 13, 14
- [22] Zhou Wang, Alan C. Bovik, Hamid R. Sheikh, and Eero P. Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE Trans. Image Process.*, 13(4):600–612, 2004. 18
- [23] Henrique Weber, Mathieu Garon, and Jean-François Lalonde. Editable indoor lighting estimation. In *Eur. Conf. Comput. Vis.*, 2022. 1, 3, 4, 5, 7, 8, 9, 10, 11, 12, 13, 14, 16, 19, 27, 28, 29, 30, 31
- [24] Jinsong Zhang, Kalyan Sunkavalli, Yannick Hold-Geoffroy, Sunil Hadap, Jonathan Eisenman, and Jean-François Lalonde. All-weather deep outdoor lighting estimation. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2019. 1, 3, 4, 7, 8, 9, 10, 11, 12, 13, 14, 17, 19
- [25] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *IEEE Conf. Comput. Vis. Pattern Recog.*, 2018. 18

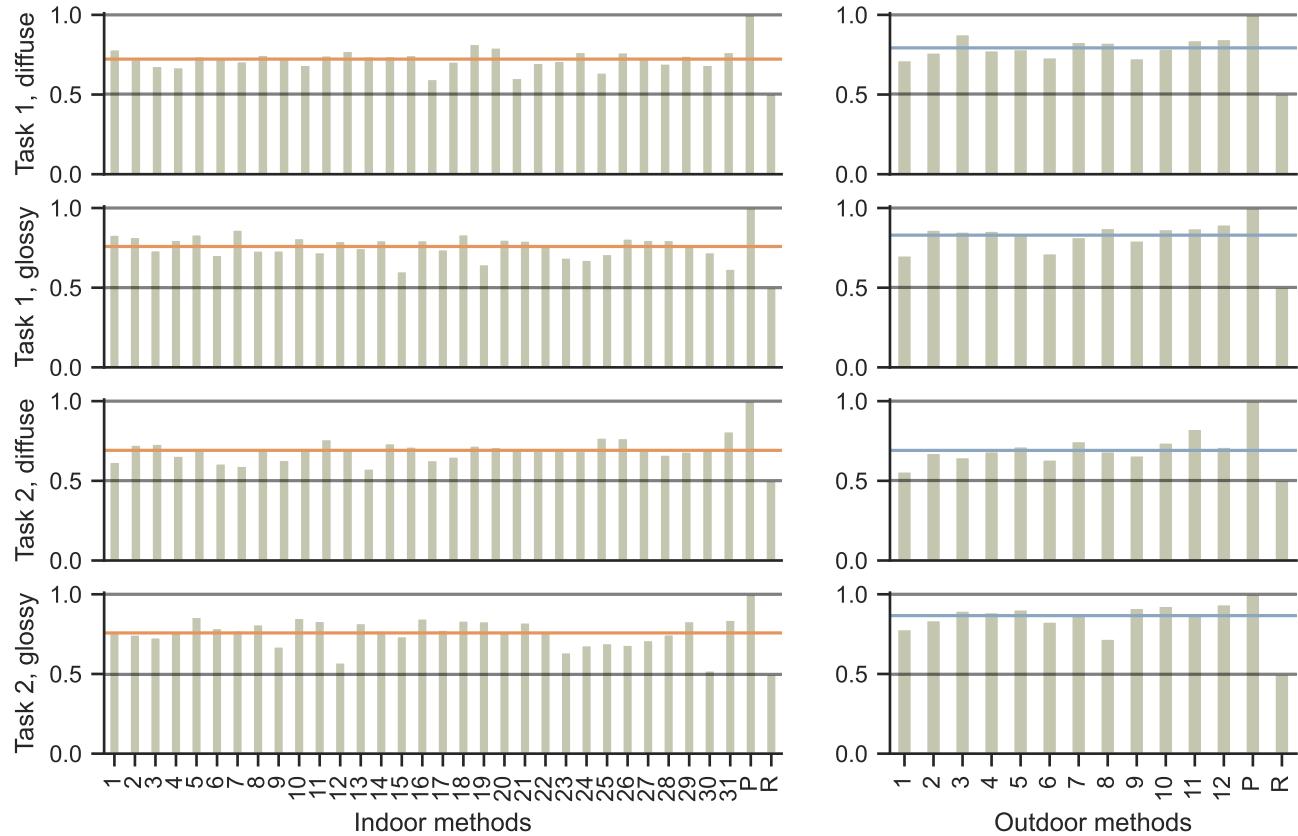


Figure 14. Agreement between the expected observer scores and the individual observer scores (columns) for all the lighting estimation methods (indoor: left bars; outdoor: right bars), for the different types of experiments (rows). The lower horizontal grey bar is set at chance level (labelled as “R”; ~ 0.5) and the higher one corresponds to the perfect observer (labelled as “P”; set at 1.0). The orange (indoor) and blue (outdoor) lines corresponds to the expected observer agreement score.

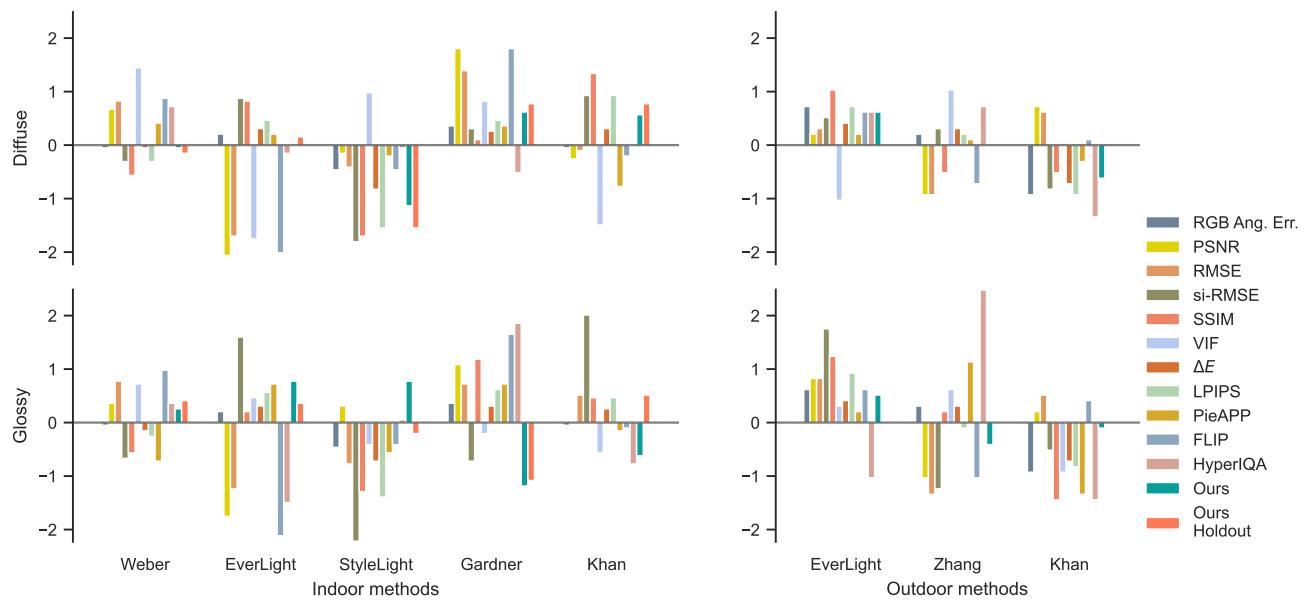


Figure 15. Thurstone Case V Law of Comparative Judgement scores for the different metrics as a function of the different lighting estimation methods (columns), for the different types of sphere materials used in the experiments (rows).

Type	Scene	Material	Method	Mean accuracy	Std accuracy
Task 1	Diffuse	Lasso		76.00	3.60
		ElasticNet		75.80	3.50
		Ridge		76.20	3.20
		LogisticRegression		76.80	3.40
		Perceptron		66.00	8.50
		ARDRegression		76.30	4.30
		BayesianRidge		76.30	3.30
		BernoulliNB		76.60	4.10
		HuberRegressor		76.20	3.50
		TheilSenRegressor		76.50	3.50
	Glossy	BayesianGaussianMixture		52.20	1.20
		GaussianMixture		52.20	1.20
		SVR		77.30	3.70
		Lasso		54.30	5.90
		ElasticNet		53.80	6.00
		Ridge		53.80	6.30
		LogisticRegression		56.40	4.80
		Perceptron		51.00	7.50
Task 2	Diffuse	ARDRegression		54.00	5.50
		BayesianRidge		54.60	8.10
		BernoulliNB		59.20	9.80
		HuberRegressor		54.80	5.80
		TheilSenRegressor		58.00	7.10
		BayesianGaussianMixture		50.60	0.60
		GaussianMixture		50.60	0.60
		SVR		56.90	5.90
		Lasso		57.20	7.20
		ElasticNet		56.80	6.40
	Glossy	Ridge		56.50	6.40
		LogisticRegression		55.60	6.60
		Perceptron		51.20	5.00
		ARDRegression		57.00	6.50
		BayesianRidge		57.00	6.90
		BernoulliNB		61.20	2.90
		HuberRegressor		57.00	6.30
		TheilSenRegressor		56.50	3.80
		BayesianGaussianMixture		53.50	1.30
		GaussianMixture		53.50	1.30
		SVR		58.60	6.90
		Lasso		59.70	9.50
		ElasticNet		59.70	9.60
		Ridge		60.20	10.00
		LogisticRegression		64.20	6.50
		Perceptron		54.30	6.80
		ARDRegression		58.60	10.50
		BayesianRidge		60.90	7.70
		BernoulliNB		65.10	5.10
		HuberRegressor		63.20	8.00
		TheilSenRegressor		59.80	7.30
		BayesianGaussianMixture		50.60	0.70
		GaussianMixture		50.60	0.70
		SVR		63.50	7.30

Table 1. Results of the different classical models trained with each of the psychophysical experiments data. The mean and standard deviation accuracy results on the k -folds ($k = 10$) experiments are given of all the models for each experiments.

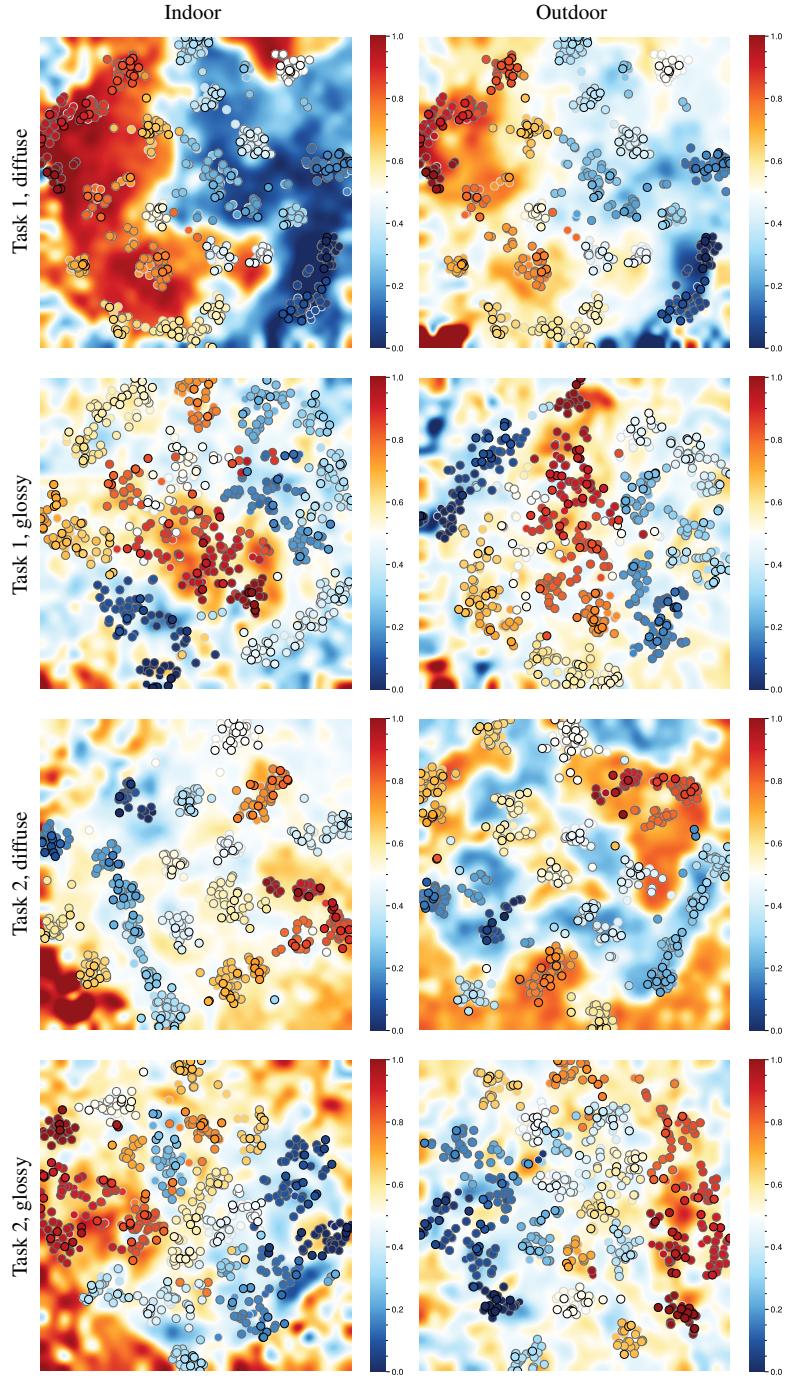


Figure 16. Visualisation of the learnt SVR metric for the four different experiments (rows), applied to an indoor/outdoor (columns) validation dataset. The projection of the ten classical metrics to \mathbb{R}^2 is done using UMAP [15]. The point corresponds to the pair comparisons of the stimuli (given as input to the network) and the background corresponds to the learnt function by the SVR. The white line corresponds to the boundary between the left and right choice of image. The colour of the data points and the background correspond to the choice of picking the left or right image given as input to the metric. The data points with the black contours indicate data points used in validation, and the grey contours correspond to the support vectors.

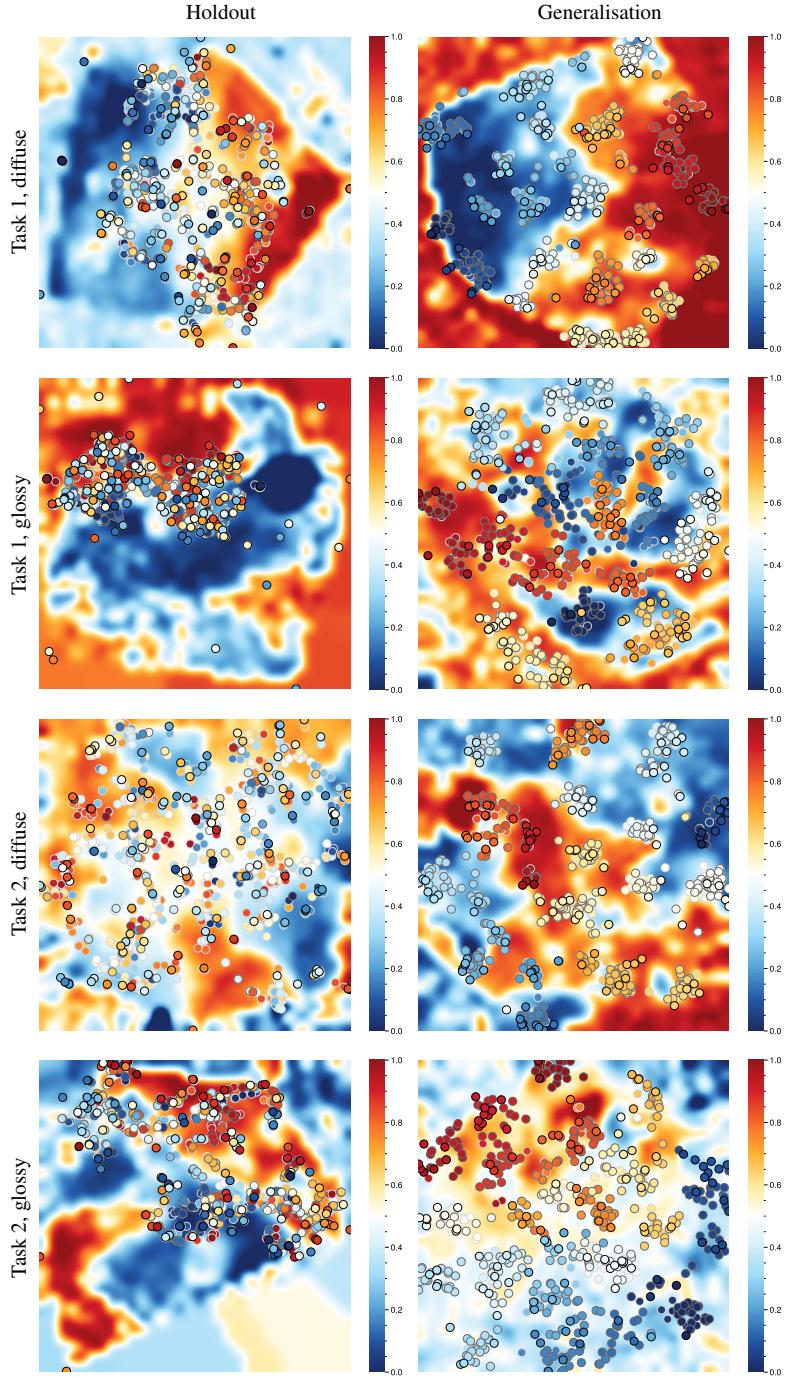


Figure 17. Visualisation of the learnt SVR metric for the four different experiments (rows), applied to the metric trained with the holdout approach (see sec. 5.1 of the main paper for details) in the left column and for the metric trained for the generalisation test (see sec. 5.2 of the main paper for details) in the right column. The projection of the ten classical metrics to \mathbb{R}^2 is done using UMAP [15]. The point correspond to the pair comparisons of the stimuli (given as input to the network) and the background corresponds to the learnt function by the SVR. The white line corresponds to the boundary between the left and right choice of image. The colour of the data points and the background correspond to choice of picking the left or right image given as input to the metric. The data points with the black contours indicate data points used in validation and the grey contours correspond to the support vectors.

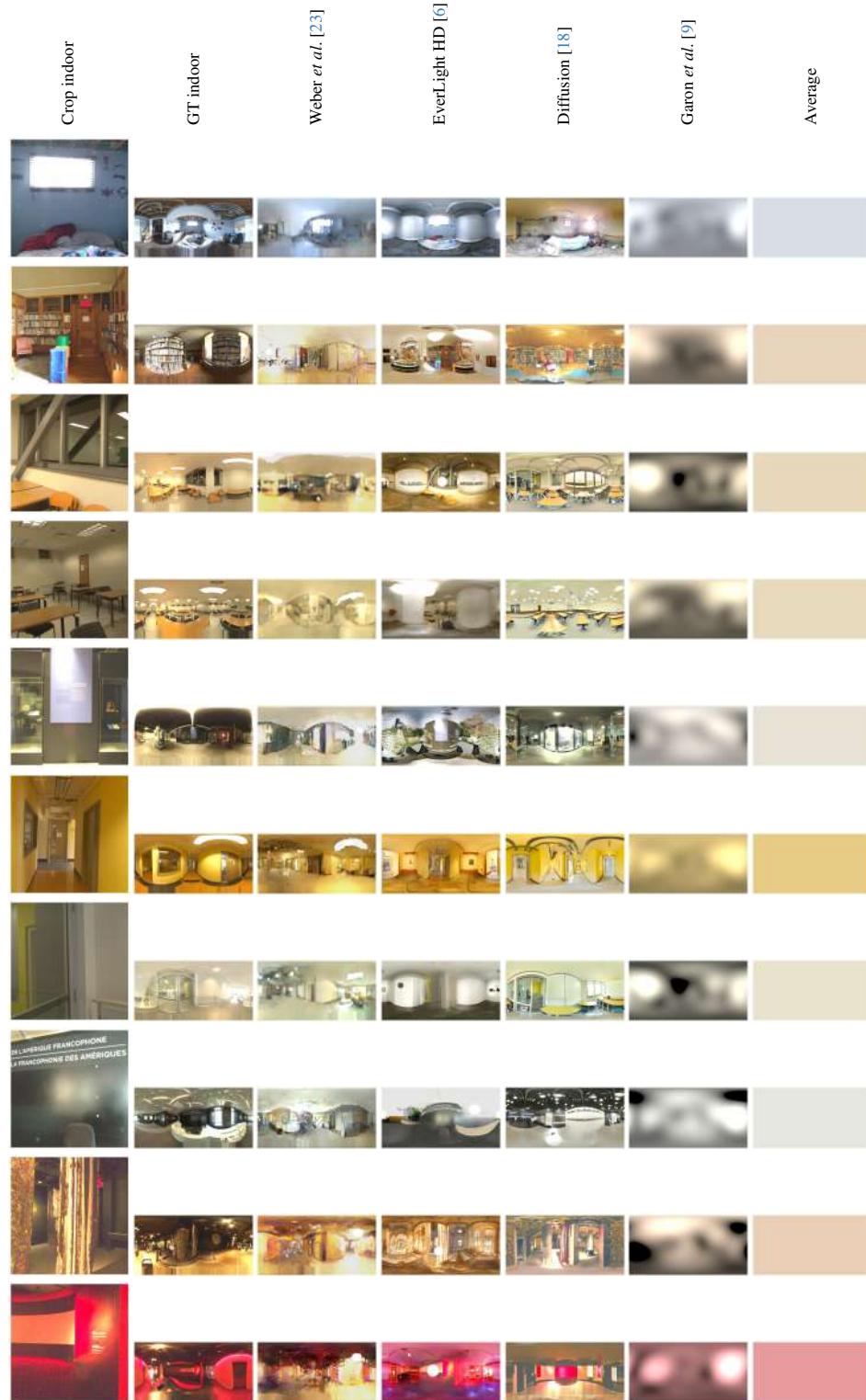


Figure 18. IBLs generated by the different generalisation lighting methods (columns) for each scene (rows). The first column corresponds to the region extracted from the indoor scene, corresponding to a 50° FoV, is taken from the centre of the full GT panorama (for most scenes), shown in the second column. The IBLs are reexposed and tonemapped with $\gamma = 2.2$ for display.

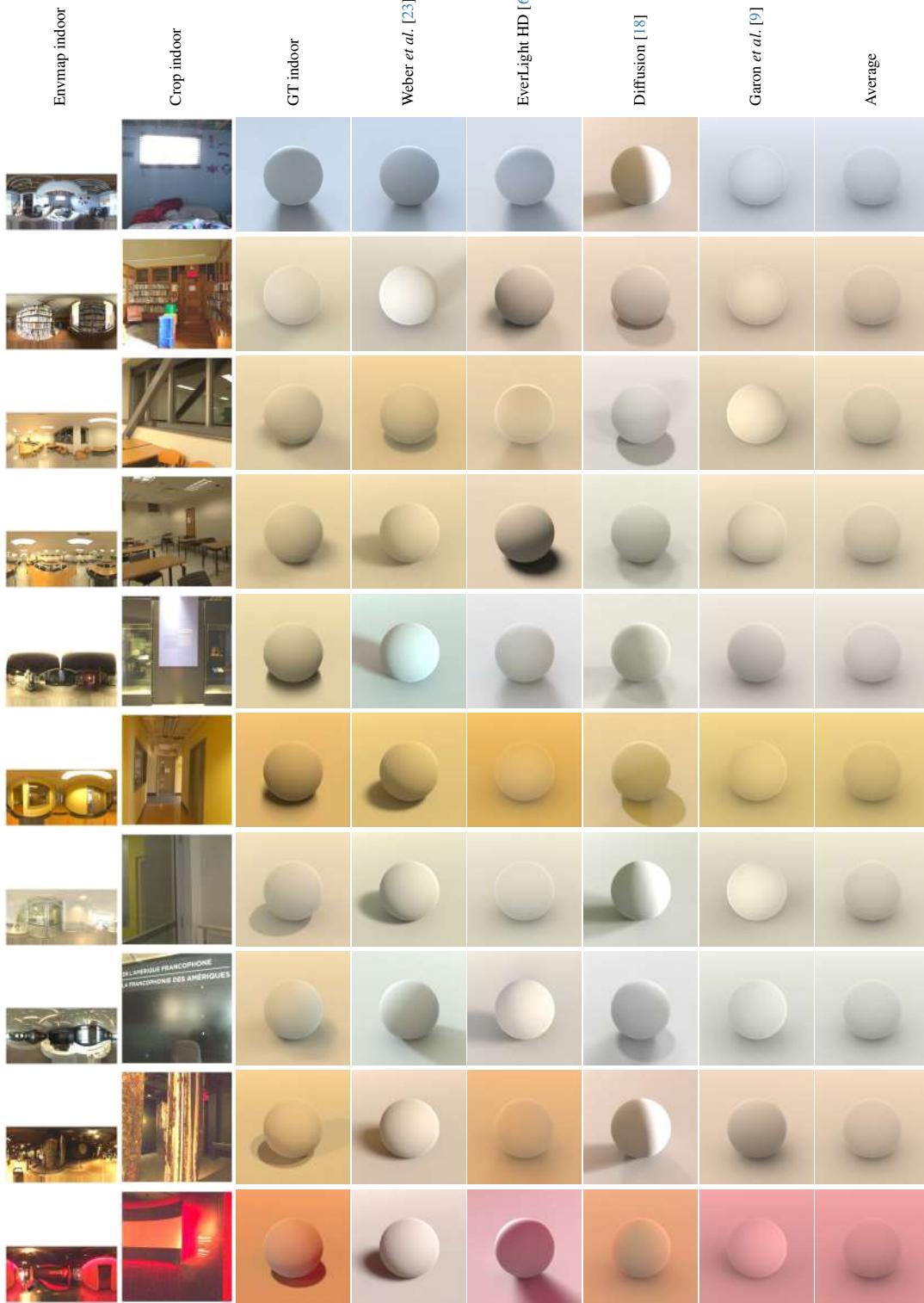


Figure 19. Stimuli used for the generalisation task 1 experiment with the diffuse sphere. The full HDR panorama (first column) is reexposed and tonemapped with $\gamma = 2.2$ for display. The region extracted from the scene (second column), corresponding to a 50° FoV, taken from the centre of the full panorama (for most scenes). The rendered stimuli using the ground truth IBLs (first column) are shown in the third column. The other columns are the rendered stimuli using the IBLs (shown in fig. 1) produced by the different lighting estimation methods.



Figure 20. Stimuli used for the generalisation task 1 experiment with the glossy sphere. The full HDR panorama (first column) is reexposed and tonemapped with $\gamma = 2.2$ for display. The region extracted from the scene (second column), corresponding to a 50° FoV, taken from the centre of the full panorama (for most scenes). The rendered stimuli using the ground truth IBLs (first column) are shown in the third column. The other columns are the rendered stimuli using the IBLs (shown in fig. 1) produced by the different lighting estimation methods.

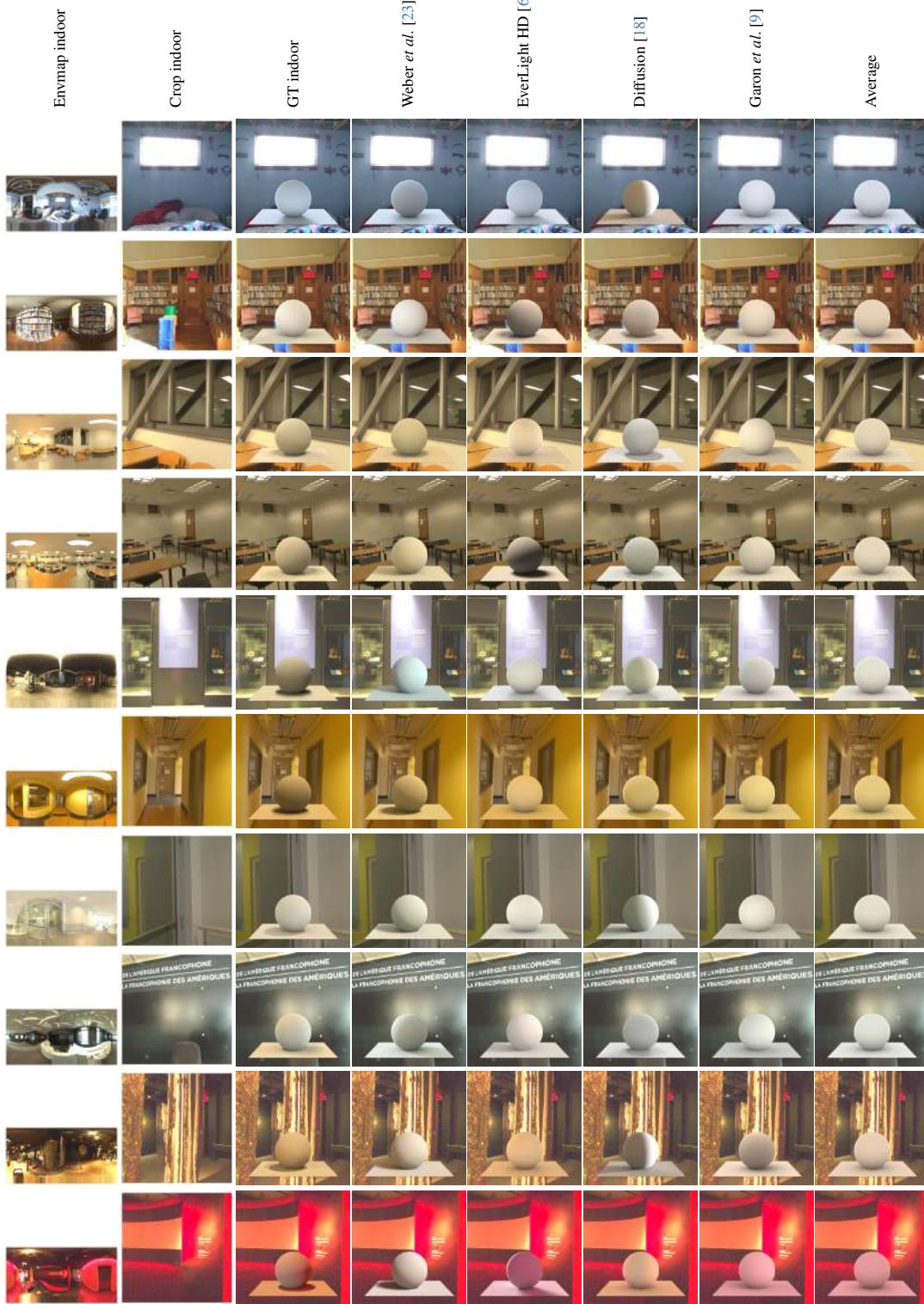


Figure 21. Stimuli used for the generalisation task 2 experiment with the diffuse sphere. The full HDR panorama (first column) is reexposed and tonemapped with $\gamma = 2.2$ for display. The region extracted from the scene (second column), corresponding to a 50° FoV, taken from the centre of the full panorama (for most scenes). The rendered stimuli using the ground truth IBLs (first column) are shown in the third column. The other columns are the rendered stimuli using the IBLs (shown in fig. 1) produced by the different lighting estimation methods.

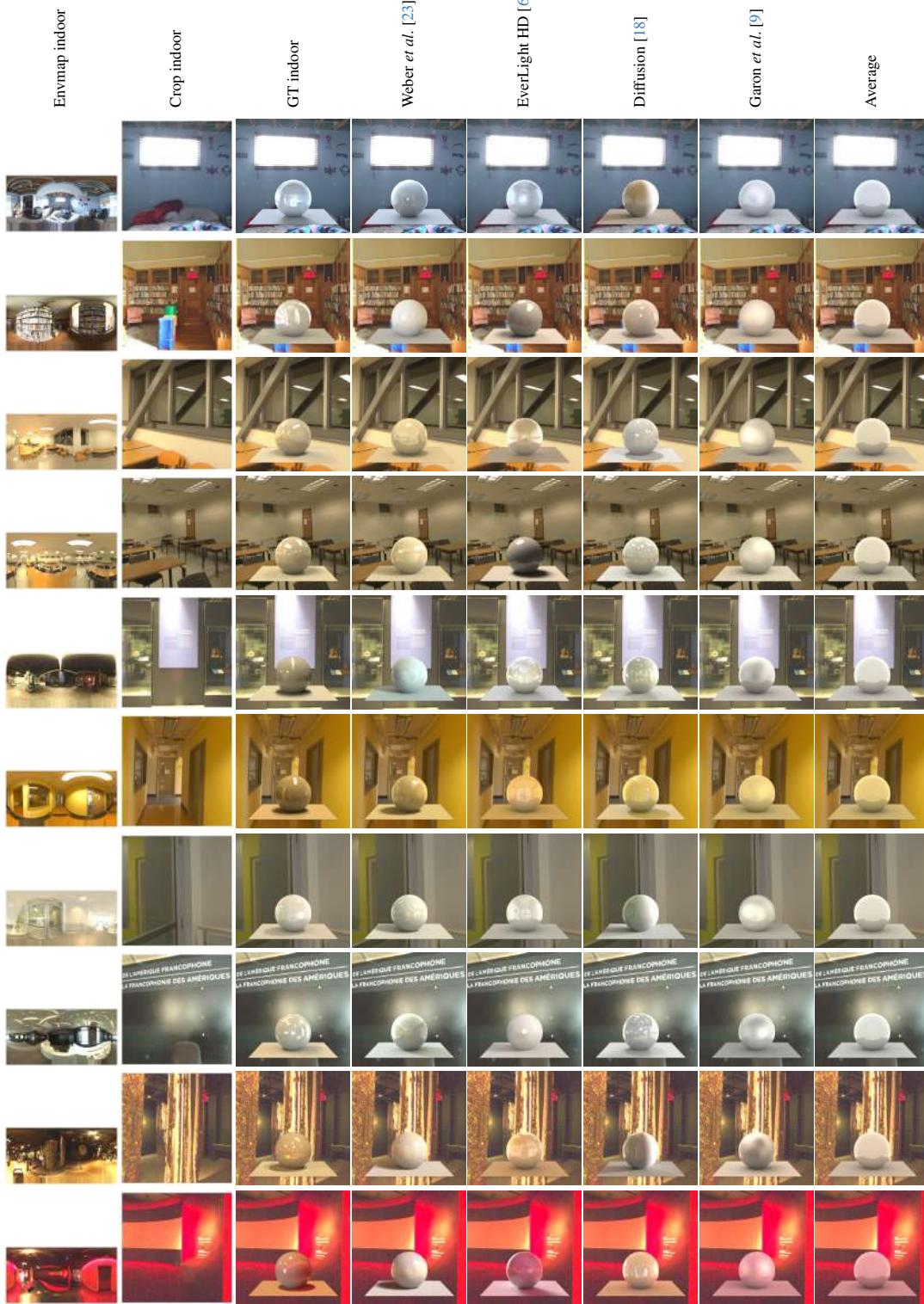


Figure 22. Stimuli used for the generalisation task 2 experiment with the glossy sphere. The full HDR panorama (first column) is reexposed and tonemapped with $\gamma = 2.2$ for display. The region extracted from the scene (second column), corresponding to a 50° FoV, taken from the centre of the full panorama (for most scenes). The rendered stimuli using the ground truth IBLs (first column) are shown in the third column. The other columns are the rendered stimuli using the IBLs (shown in fig. 1) produced by the different lighting estimation methods.

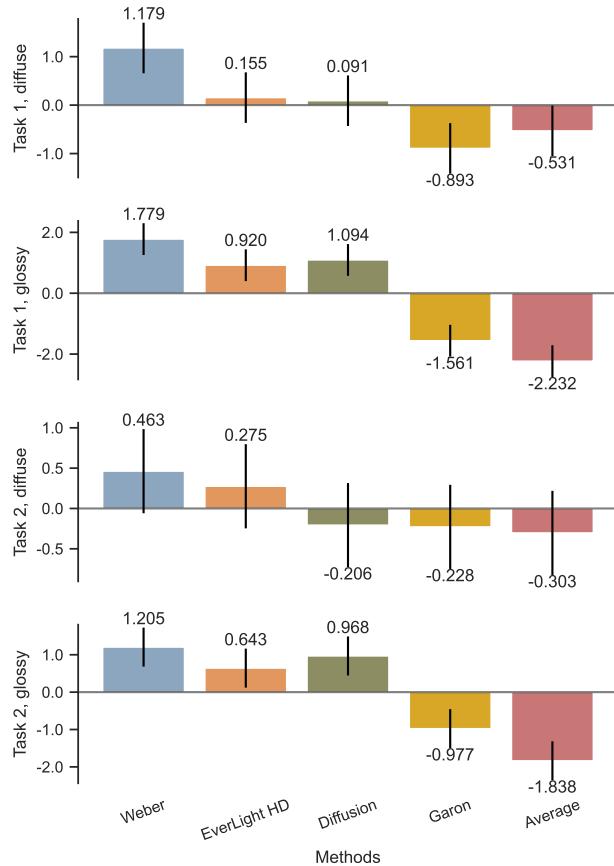


Figure 23. Thurstone Case V Law of Comparative Judgement scores from the observers as a function of the different generalisation lighting estimation methods (columns), for the different types of sphere materials used in the experiments (rows). Error bars correspond to 95 % confidence interval.