

Mirror, Mirror, on the Wall, Who's Got the Clearest Image of Them All? – A Tailored Approach to Single Image Reflection Removal

Daniel Heydecker*, Georg Maierhofer*, Angelica I. Aviles-Rivero*,
Qingnan Fan, Dongdong Chen, Carola-Bibiane Schönlieb and Sabine Süsstrunk

Abstract—Removing reflection artefacts from a single image is a problem of both theoretical and practical interest, which still presents challenges because of the massively ill-posed nature of the problem. In this work, we propose a technique based on a novel optimisation problem. Firstly, we introduce a *simple* user interaction scheme, which helps minimise information loss in reflection-free regions. Secondly, we introduce an H^2 fidelity term, which preserves fine detail while enforcing global colour similarity. We show that this combination allows us to mitigate some major drawbacks of the existing methods for reflection removal. We demonstrate, through numerical and visual experiments, that our method is able to outperform the state-of-the-art methods and compete with recent deep-learning approaches.

Index Terms—Reflection Suppression, Image Enhancement, Optical Reflection.

I. INTRODUCTION

This paper addresses the problem of single image reflection removal. Reflection artefacts are ubiquitous in many classes of images; in real-world scenes, the conditions are often far from optimal, and photographs have to be taken in which target objects are covered by reflections and artefacts appear in undesired places. This does not only affect amateur photography; such artefacts may also arise in documentation in museums and aquariums, or black-box cameras in cars (see Fig. 1). It is therefore unsurprising that the problem of how to remove reflection artefacts is of great interest, from both practical and theoretical points of view.

Although it is possible to reduce reflection artefacts by the use of specialised hardware such as polarisation filters [1], [2], [3], this option has several downsides. Firstly, even though the use of hardware can have a significant effect on removing the reflection, it only works when certain capture conditions are fulfilled, such as Brewster's angle [4]. In practise, it is difficult to achieve optimal capture conditions, which results in residual reflections [5], [6]. As a result, post-processing techniques are often needed for further improvement of the image. Moreover, for the purposes of amateur photography, the use of specialised hardware is expensive, and consequently less appealing.

*These three authors contributed equally and hold joint first authorship.

Daniel Heydecker, Georg Maierhofer, Angelica I. Aviles-Rivero and Carola-Bibiane Schönlieb are with the DAMTP and DPMMS, University of Cambridge. {dh489,gam37,ai323,cbs31}@cam.ac.uk

Qingnan Fan is with the Computer Science and Technology School, Shandong University. fqncchina@gmail.com; Dongdong Chen is with the University of Science and Technology of China. cddlyf@gmail.com

Sabine Süsstrunk is with the School of Computer and Communication Sciences, École Polytechnique Fédérale de Lausanne. sabine.sussstrunk@epfl.ch

As an alternative to the use of specialised hardware, a body of research has established a variety of computational techniques. These can be divided in those that use *multiple images*, and those that use a *single image*. The former techniques employ images from various view points (e.g. [7], [8], [9], [10]), with the aim of exploiting temporal information to separate the reflection artefacts from the observed target, while for the latter, carefully selected image priors are used to obtain a good approximation of the target object, for example [11], [12], [13], [14].

Although the use of multiple images somewhat mitigates the massively ill-posed problem created by the reflection removal formulation, the success of these techniques requires multiple images from several viewpoints and their performance is strongly conditional on the quality of the acquired temporal information. Moreover, in practice, acquisition conditions are non-optimal, which often results in image degradation, causing occlusions and blurring in the images. Therefore, either many images or post-processing are needed, which strongly restricts the applicability and feasibility of these methods to a typical end-user. These constraints make single-image methods a focus of great attention to the scientific community, since it is appropriate for most users, and this is the approach which we will take in this paper.

Mathematically, an image \mathbf{Y} containing reflection artefacts can be represented as a linear superposition [15] as:

$$\mathbf{Y} = \mathbf{T} + \mathbf{R}, \quad (1)$$

where \mathbf{T}, \mathbf{R} are $n \times m$ matrices representing the transmission layer and reflection layer, respectively. Therefore, the goal of a reflection suppression technique is to approximate \mathbf{T} from the acquired image \mathbf{Y} .

Although the body of literature for single-image reflection removal has proven promising results, this remains an open problem, and there is still potential for further enhancements. We consider the problem of how to get a better approximation of \mathbf{T} .

In this work, we propose a new approach, closely related to [14], and inspired by the observation that even *low-level* user input may contain a lot of information. Our technique relies on additional information, which gives the rough location of reflections. In our experiments, this is given by user-input; in principle, this could be done by an algorithmic or machine-learning technique. We recast the reflection removal

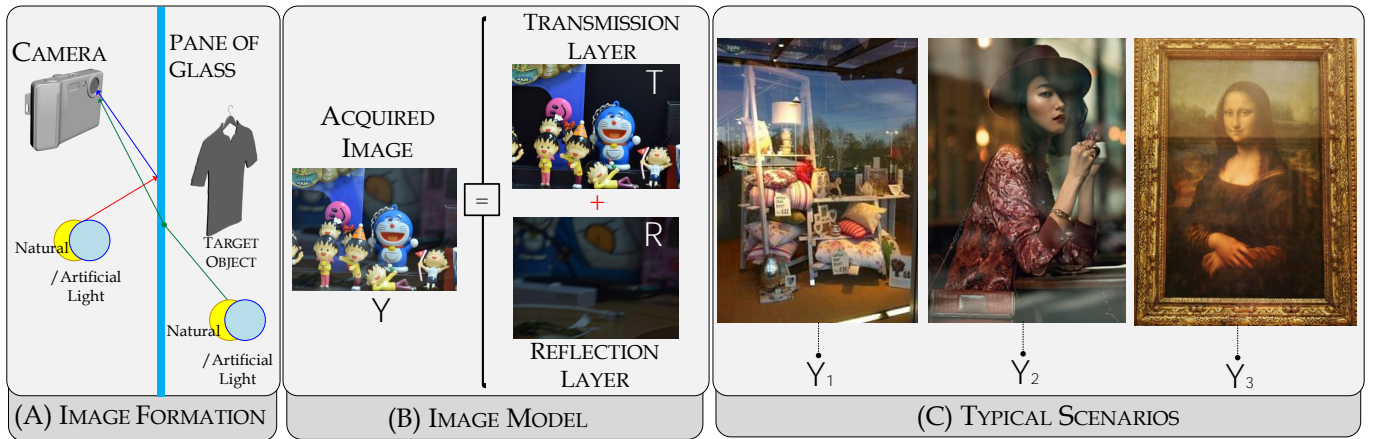


Fig. 1: (A) An illustration of the image formation in which a target object captured through a pane of glass will have reflection artefacts. (B) Based on the image model, an acquired image (Y) can be decomposed into two layers: Transmission (T) and Reflection (R). (C) images ($Y_{1,2,3}$) show a set of typical situations where there is no option but to take the picture through a pane of glass such as store display or in museums.

problem as an optimisation problem which is solved iteratively, by breaking it up into two more computationally tractable problems. Compared to existing solutions from the literature, we achieve a better approximation of T from a well-chosen optimisation problem, while simultaneously preserving image details and eliminating global colour shifts. Our contributions are as follows:

- We propose a computationally tractable mathematical model for single-image reflection removal, in which we highlight:
 - A *simple and tractable* user interaction method to select reflection-heavy regions, which is implemented at the level of the optimisation problem as a spatially aware prior term. We show that this improves the retention of detail in reflection-free areas.
 - A combined H^2 *fidelity term*, which combines L^2 and Laplacian terms. We show that this combination yields significant improvements in the quality of the colour and structure preservation.

We establish that the resulting optimisation problem can be solved efficiently by half-quadratic splitting.

- We validate the theory with a range of numerical and visual results, in different scenes and under varying capture conditions.
- We demonstrate that the combination of our fidelity term and prior term leads to a better approximation of T than state-of-the-art model based techniques, and can compete with the most recent deep-learning (DL) techniques.

II. RELATED WORK

The problem of image reflection removal has been extensively investigated in the computer vision community, in which solutions rely on using multiple images and single image data, alone or in combination with specialised hardware. In this section, we review the existing techniques in turn.

A number of techniques have been developed which use information from multiple images to detect and remove reflec-

tions. These include the use of different polarisation angles [5], [16], [6], [17], [3], adjustment of focal and flash settings [18], [1], [2], and the uses of relative motion and coherence [7], [8], [19], [20], [21], [22], [23], [24], [25]. A recent technique [26] seeks to improve on these methods by seeking to match the *transmitted* layer, while other techniques may erroneously match the reflected layer. Each of these techniques requires particular modelling hypotheses to be met, and advantageous capture conditions which may not be feasible in practice.

We now review the related works in single image techniques, as they are most applicable to everyday capture. A commonality of these techniques is the choice of a sparse gradient prior, which imposes a preference for output transmission layers, T , with few strong edges.

A user-intervention method was proposed in [11], which labels gradients as belonging to either transmission or reflection layer. They then propose to solve a constrained optimisation problem, with prior distribution given by the superposition of two Laplace distributions. A similar optimisation problem is used by [13], which replaces user-intervention labelling by a depth-of-field based inference scheme, while [27] relies on ghosting artefacts.

Our work is most closely related to the optimisation-based models and techniques of [12], [14]. The authors of [12] propose a smooth gradient prior on the reflection layer, and a sparse gradient prior on the transmission layer. This approach was adapted by Arvanitopoulos et al. in [14], who proposed a Laplacian-based fidelity term with a novel sparse gradient prior. This preserves (Gestalt) continuity of structure, while also reducing loss of high-frequency detail in the transmission layer. The algorithm they propose is both more effective, and more computationally efficient, than the other techniques discussed above.

The application of deep learning to reflection removal was pioneered by Fan et al. in [28]. In this work, the authors propose a deep neural network structure, which firstly predicts the edge map and then separates the layers. This technique

outperforms the algorithmic approach of [12]. Further work in this direction was made by Zhang et al. [29], who use a fully convolutional neural network with three loss terms, which help to ensure preservation of features and pixel-wise separation of the layers. Wan et al. [30] seek to use a loss function inspired by human perception to estimate the gradient of the transmission layer, and use this to concurrently estimate the two layers using convolutional neural networks, and Jin et al. [31] proposes a convolutional neural network with a resampling strategy, to capture features of global priors, and avoid the ambiguity of the average colour. Most recently, Yang et al [32] propose a bidirectional deep learning-scheme based on a cascade neutral network. This method first estimates the background layer \mathbf{T} , then uses this to estimate the reflected layer \mathbf{R} . Finally, the estimate on \mathbf{R} is used to improve the estimate of \mathbf{T} .

The philosophy of our approach is similar to that of [11]. Motivated by the principle that *humans are good at distinguishing reflections*, both our work and [11] seek to exploit further user input to assist an algorithmic technique. However, we emphasise that we are the first to propose a *simple and tractable* user interaction scheme: in evaluating our user interaction scheme in Section IV/E3, we will see that our user interaction scheme requires very little effort from the user, and that our algorithm performs well with even very crude selection. By contrast, the algorithm of [11] requires much more effort, and a much more detailed input.

III. PROPOSED METHOD

This section contains the three key parts of the proposed mathematical model: (i) the combined Laplacian and L^2 fidelity term, (ii) a *spatially aware* prior term, given by user input, and (iii) the implementation using quadratic splitting for computational tractability.

Although the model for an image with reflection artefacts described in (1) is widely-used, our solution adopts the observation of [1], [12], [14] that the reflection layer is less in focus and often blurred, which we formalise as follows:

Observation 1. *In many cases, the reflected image will be blurred, and out of focus. This may be the case, for instance, if the reflected image is at a different focal distance from the transmitted layer. Moreover, reflections are often less intense than the transmitted layer.*

Based on this observation, the image model [1], [12] which we adapt is

$$\mathbf{Y} = w\mathbf{T} + (1 - w)(\mathbf{k} \star \mathbf{R}), \quad (2)$$

where \star denotes convolution, w is a weight $w \in [0, 1]$ that controls the relative strength of reflections, and \mathbf{k} is a blurring kernel.

A. Fidelity and Prior Terms.

We begin by discussing the prior term. Loss of some detail, in reflection heavy regions, is to be expected, and is a result of the ill-posed nature of reflection suppression. We seek to use low-level user input to reduce the loss of detail in *reflection-free regions*, motivated by the following observation:

Observation 2. *In many instances, the reflections are only present in a region of the image, and it is easy for an end user to label these areas. In regions where reflections are not present, all edges in \mathbf{Y} arise from \mathbf{T} , and so should not be penalised in a sparsity prior. Moreover, in certain instances, it may be particularly important to preserve fine detail in certain regions.*

For instance, for photographs containing a window, the reflections will only occur in the window, and not elsewhere in the image. To this end, we propose to incorporate a *region selection function* ϕ , taking values in $[0, 1]$, into a *spatially aware prior*:

$$P(\phi, \mathbf{T}) = \sum_{i,j} \phi_{ij} 1[\nabla_x T_{ij} \neq 0 \text{ or } \nabla_y T_{ij} \neq 0]. \quad (3)$$

Here, $1[\dots]$ denotes the indicator function for the set of indexes (i, j) where one of the gradients $\nabla_x \mathbf{T}, \nabla_y \mathbf{T}$ is nonzero. We assume that the region selection function ϕ is given by the user, along with the input. Although this is philosophically similar to the user intervention method of [11], our approach is drastically less effort-intensive: rather than labelling many edges, it is sufficient to (crudely) indicate which regions contain reflections. The practicalities of our technique will be discussed in Subsection C below. We will show that, by choosing $\phi_{ij} \approx 1$ on reflection-heavy regions and $\phi_{ij} \approx 0$ elsewhere, we can minimise the loss of detail in reflection-free areas, and avoid the ‘flattening’ effect described above. We also note that a naïve attempt to apply the approach of [14] to a region of the image produces noticeable colour shifts at the boundary of the selected region, which our spatially aware prior term avoids.

We now consider the fidelity term, seeking to build on the Laplacian fidelity term proposed by [14]; this choice of fidelity term penalises over-smoothing, and enforces consistency in fine details. Although this improves on the L^2 fidelity term of Xu et al. [33], one can still observe significant ‘flattening’ effects. These arise when there is significant colour variation over a large area: individual gradients are weak, and are neglected by the technique of [14]. This results in the whole region being given the same value in the output \mathbf{T} , producing unrealistic and visually unappealing results. Moreover, we also note that for any constant matrix \mathbf{C} the Laplacian is invariant under the transformation $\mathbf{T} \mapsto \mathbf{T} + \mathbf{C}$. As a result, the algorithm proposed by [14] risks producing global colour shifts; at the level of the optimisation problem, this reflects the non-uniqueness of minimisers. To eliminate this possibility, we propose a combined fidelity term:

$$d_\gamma(\mathbf{T}, \mathbf{Y}) = \|\Delta \mathbf{T} - \Delta \mathbf{Y}\|_2^2 + \gamma \|\mathbf{T} - \mathbf{Y}\|_2^2, \quad (4)$$

where $\Delta \mathbf{T}$ is the discrete Laplacian defined as $\Delta \mathbf{T} = \nabla_{xx} \mathbf{T} + \nabla_{yy} \mathbf{T}$, and γ is a positive parameter controlling the relative importance of the two terms. We will see, in numerical experiments, that this leads to results with more natural, saturated colours, and which are consequently more visually pleasing. We remark that other kernel filters are possible which would play the same role of measuring structure, such as the

discrete gradient ∇ , or more complicated elliptic second-order operators; we use the Laplacian for the following reasons. Firstly, the Laplacian penalises loss of high-frequency detail more strongly than first order operators such as ∇ , as can be seen by moving to Fourier space, and so our choice will preserve high-frequency details well. Secondly, the Laplacian is a simple measure of structure, and which is invariant under the (natural) symmetry of rotation.

Combining the prior and fidelity terms, as defined in (3) and (4), our optimisation problem is therefore

$$\mathbf{T}^* = \operatorname{argmin}_{\mathbf{T}} \left\{ \|\Delta \mathbf{T} - \Delta \mathbf{Y}\|_2^2 + \gamma \|\mathbf{T} - \mathbf{Y}\|_2^2 + \lambda P(\phi, \mathbf{T}) \right\}. \quad (5)$$

Here, λ is a regularisation parameter to be chosen later. The reader is invited to compare this optimisation problem to the similar problem of (localised) L^0 image smoothing, but to note the important difference of having a fidelity term including the image Laplacian. In the next section, we will detail how the proposed optimisation problem can be solved in a tractable computational manner by using quadratic splitting.

B. Solving the Optimisation Problem.

We solve the optimisation problem introduced in (5) by half-quadratic splitting. We introduce auxiliary variables $\mathbf{D}^x, \mathbf{D}^y$ as proxies for, respectively, $\nabla_x \mathbf{T}$ and $\nabla_y \mathbf{T}$. For ease of notation, we write \mathbf{D} for the pair $[\mathbf{D}^x, \mathbf{D}^y]$, and similarly $\nabla \mathbf{T}$ for the pair $[\nabla_x \mathbf{T}, \nabla_y \mathbf{T}]$. This leads to the auxiliary problem:

$$\mathbf{T}^*, \mathbf{D}^* = \operatorname{argmin}_{\mathbf{T}, \mathbf{D}} \left\{ \|\Delta \mathbf{T} - \Delta \mathbf{Y}\|_2^2 + \gamma \|\mathbf{T} - \mathbf{Y}\|_2^2 + \lambda P(\phi, \mathbf{D}) + \beta \|\mathbf{D} - \nabla \mathbf{T}\|_2^2 \right\} \quad (6)$$

where $\beta \in \mathbb{R}_{>0}$ is a penalty parameter yet to be chosen, and we use the shorthand

$$P(\phi, \mathbf{D}) = \sum_{i,j} \phi_{ij} 1[D_{ij}^x \neq 0 \text{ or } D_{ij}^y \neq 0]. \quad (7)$$

Notice that in the limit $\beta \rightarrow \infty$ the axillary penalty term ensures that we recover the solution to the original optimisation problem (5). Hence, we may solve the optimisation problem (6) by splitting into two more computational tractable problems. We alternate between optimising over \mathbf{T} and \mathbf{D} , while keeping the other fixed; at the same time, we increment β so that, after a large number of steps, \mathbf{D} is a good approximation of $\nabla \mathbf{T}$. We give details on the solution of each sub-problem below, and the full solution is presented in Algorithm 1.

►**Sub-problem 1: Optimisation over \mathbf{T} .** For a fixed \mathbf{D} , we wish to optimise:

$$\mathbf{T}^* = \operatorname{argmin}_{\mathbf{T}} \left\{ \|\Delta \mathbf{T} - \Delta \mathbf{Y}\|_2^2 + \gamma \|\mathbf{T} - \mathbf{Y}\|_2^2 + \beta \|\mathbf{D} - \nabla \mathbf{T}\|_2^2 \right\}. \quad (8)$$

The objective function is now quadratic in \mathbf{T} . We note that the discrete gradient ∇ and the discrete Laplacian Δ are both linear maps which take an $m \times n$ image matrix to an array of size $2 \times m \times n$ and $m \times n$ respectively. We may thus view these

Algorithm 1 Our Proposed Method

- 1: Start from $\mathbf{T} \leftarrow \mathbf{Y}$ and $\beta = \beta_{\min}$
 - 2: **while** $\beta \leq \beta_{\max}$ **do**
 - 3: Optimise over \mathbf{D} , for the current value of \mathbf{T} :
 Set $(D_{ij}^x, D_{ij}^y) = \begin{cases} (0, 0) & \text{if } |(\nabla_x T_{ij}, \nabla_y T_{ij})|_2^2 \leq \frac{\lambda_{ij}}{\beta}; \\ (\nabla_x T_{ij}, \nabla_y T_{ij}) & \text{o.w.} \end{cases}$
 - 4: Using ADAM [34] and (12), find the minimum \mathbf{T}^* of (8), and replace $\mathbf{T} \leftarrow \mathbf{T}^*$.
 - 5: Increment $\beta \leftarrow \kappa \beta$
 - 6: **end while**
 - 7: **return** \mathbf{T} .
-

linear maps as tensors, and use index notation to describe their action on an image (T_{ij}) as follows:

$$(\nabla_\mu \mathbf{T})_{ij} = \sum_{k,l} \nabla_{ijkl}^\mu T_{kl}; \quad 1 \leq i \leq m, \quad 1 \leq j \leq n, \quad \mu \in \{x, y\} \quad (9)$$

and similarly:

$$(\Delta \mathbf{T})_{ij} = \sum_{k,l} \Delta_{ijkl} T_{kl}; \quad 1 \leq i \leq m, \quad 1 \leq j \leq n. \quad (10)$$

With this notation, we can write the objective function as:

$$F_1(\mathbf{T}, \mathbf{D}) = \beta \sum_{\substack{1 \leq i \leq m \\ 1 \leq j \leq n \\ \mu \in \{x, y\}}} \left(D_{ij}^\mu - \sum_{1 \leq k \leq m, 1 \leq l \leq n} \nabla_{ijkl}^\mu T_{kl} \right)^2 + \sum_{1 \leq i \leq m, 1 \leq j \leq n} \left(\left(\sum_{1 \leq k \leq m, 1 \leq l \leq n} \Delta_{ijkl} (T_{kl} - Y_{kl}) \right)^2 + \gamma (T_{ij} - Y_{ij})^2 \right) \quad (11)$$

We observe that this is quadratic, and in particular smooth, in the components T_{ij} . Using the summation convention, we compute the gradient:

$$\frac{\partial}{\partial T_{ij}} F_1(\mathbf{T}, \mathbf{D}) = 2\Delta_{abij}\Delta_{abkl}(T_{kl} - Y_{kl}) + 2\gamma(T_{ij} - Y_{ij}) + 2\beta\nabla_{abij}^\mu(\nabla_{abkl}^\mu T_{kl} - D_{ab}^\mu). \quad (12)$$

We use this computation, together with ADAM [34], a first-order gradient descent method in stochastic optimisation, to efficiently optimise over \mathbf{T} .

►**Sub-problem 2: Optimisation over \mathbf{D} .** For a fixed \mathbf{T} , the optimisation problem in \mathbf{D} is given by

$$\mathbf{D}^* = \operatorname{argmin}_{\mathbf{D}} \left\{ \beta \|\mathbf{D} - \nabla \mathbf{T}\|_2^2 + \lambda P(\phi, \mathbf{D}) \right\}. \quad (13)$$

Although the objective function, F_2 , is neither convex nor smooth, due to the L^0 prior term, we observe that it separates as

$$F_2(\mathbf{T}, \mathbf{D}) = \sum_{i,j} \left[\beta (|D_{ij}^x - \nabla_x T_{ij}|^2 + |D_{ij}^y - \nabla_y T_{ij}|^2) + \lambda \phi_{ij} 1((D_{ij}^x, D_{ij}^y) \neq 0) \right]. \quad (14)$$

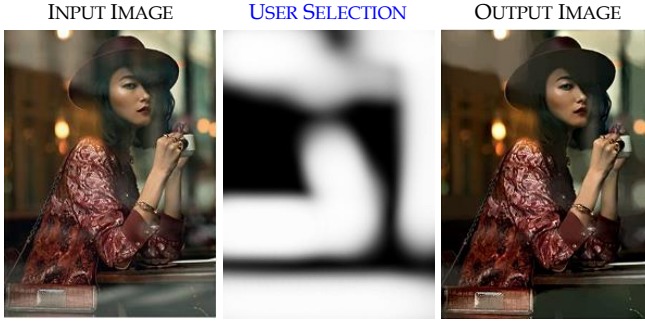


Fig. 2: From left to right. Input image, visualisation of the user interaction in practise and output image with our technique.

By explicitly solving the separated problems for each pair (D_{ij}^x, D_{ij}^y) , it is straightforward to see that a solution to (14) is given by

$$(D_{ij}^x, D_{ij}^y) = \begin{cases} (0, 0) & \text{if } |(\nabla_x T_{ij}, \nabla_y T_{ij})|_2^2 \leq \frac{\lambda \phi_{ij}}{\beta}; \\ (\nabla_x T_{ij}, \nabla_y T_{ij}) & \text{otherwise.} \end{cases} \quad (15)$$

Moreover, this minimiser is unique, provided that none of the edges are in the boundary case $|(\nabla_x T_{ij}, \nabla_y T_{ij})|_2^2 = \frac{\lambda \phi_{ij}}{\beta}$.

Hence, the optimisation (13) removes gradients below the *local* threshold $\frac{\lambda \phi_{ij}}{\beta}$. We will show, in numerical experiments, that this has the effect of smoothing *only* the selected regions, while keeping the strong edges which force continuity of structures, as was described in Section II.

The overall procedure of our method – in which previous individual steps are combined to solve the original optimisation problem (5) – is listed in Algorithm 1.

C. User Interaction Scheme.

We describe the user interaction scheme, and how the region selection function ϕ_{ij} may be obtained in practice. We recall that ϕ is responsible for passing information about the location of reflection into the algorithm, and that it takes values in the range $[0, 1]$ with

- ϕ_{ij} close to 1 if a reflection is present at pixel (i, j) and
- ϕ_{ij} close to 0 if no reflection is present at pixel (i, j) .

In practise a user, or an arbitrary instance that can recognise rough locations of reflections, is given an image, as in left-side of Fig. 2, and selects the regions in which reflections are present. A possible result can be seen in the middle part of Fig. 2, where the values of ϕ_{ij} are displayed as the grey-values in the image. This selection is then fed into our algorithm together with the input image to produce the reflection removed output as shown at right side of Fig.2.

In the absence of user interaction, we default to $\phi_{ij} \equiv 1$; that is, we assume reflections are present throughout the image.

It is noteworthy that the way this selection is performed is very simple and requires little effort. This makes it suitable for a range of applications, from an amateur human user, to algorithms that can recognise reflections, even in a very crude manner. For our experiments, the selection was performed by creating an overlay image in a raster graphics

editor, where white regions are marked with a rough brush on top of reflections. This process can be performed in a matter of seconds for each image. The results can, of course, improve with increasing selection quality, but even a rough selection produces significant improvements over no selection; see Section IV/E3 for experiments and discussion. Examples of region selection in practice are included in Section IV of the supplementary material.

D. Performance Reasoning of Parameters

Our procedure uses two parameters λ, γ , and an auxiliary parameter β in intermediary optimisation steps. We think of β as a coupling parameter, which determines the importance of the texture term in comparison to the coupling to the auxiliary variable \mathbf{D} . At later iterations, β is large and the coupling is strong, which justifies the use of \mathbf{D} as a proxy for $\nabla \mathbf{T}$.

The parameter λ determines the relative importance of preserving the structure versus preserving the texture. In terms of the model described above, it controls the importance of the penalty term $P(\phi, \mathbf{T})$ against the Laplacian $\|\Delta \mathbf{T} - \Delta \mathbf{Y}\|_2^2$. In regions where $\lambda \phi_{ij}$ is comparatively large, the sparsity of edges is much more important than the texture. Therefore, any edges which do not enforce structure will be washed out, and the region is smoothed during the optimisation over \mathbf{D} . On the other hand, in regions where $\lambda \phi_{ij}$ is comparatively small, the texture term dominates, and only very few edges are removed. In terms of the algorithm, this corresponds to controlling the edge threshold $\frac{\lambda \phi_{ij}}{\beta}$. This is illustrated in the supplementary material.

We also give an interpretation of why it is natural to increase β in this way. In the first stages of the iteration, β is very small, and so the threshold keeps only the largest magnitude edges, and sets most edges of reflection-heavy areas to 0. After each iteration, β increases and the threshold $\frac{\lambda \phi_{ij}}{\beta}$ decreases, and so the next iteration will preserve more edges. *Hence, in reflection-heavy areas, we include edges in decreasing order of magnitude*; this corresponds to looking at strongly-defined structures first, and then considering incrementally weaker structure. This is illustrated in the supplementary material.

We give a theoretical basis for excluding the limiting regimes of either $\gamma \ll 1$ or $\gamma \geq 1$. In the regime where $\gamma \ll 1$, we may consider a step of the gradient descent to be a step of ‘uncorrected’ gradient descent, with $\gamma = 0$, followed by a small correction $\gamma(\mathbf{Y} - \mathbf{T})$ to correct colour shift. For this reason, if $\gamma \ll 1$ is too small, our algorithm will not adequately correct for colour shifts. On the other hand, if $\gamma > 1$, then the L^2 term dominates the Laplacian term, and we expect blurring and loss of texture, as discussed in [14].

IV. EXPERIMENTAL RESULTS

In this section, we describe in detail the range of experiments that we conducted to validate our proposed method.

A. Data Description.

We evaluate the theory using the following three datasets. Firstly, we use real-world data from the SIR² benchmark



Fig. 3: (E1). Examples of the output, along with ground truth, of our approach compared against AR17 [14]. The examples with varying settings such as the focus in (A) and (B) and the glass thickness in (C) and (D). The three evaluation metrics of the reflection-free image are computed using the ground truth.

dataset [35]. The dataset is composed of 1500 images with size of 400×540 , and provides variety in scenes with different degrees of freedom in terms of aperture size and thickness of the glass. These variations allow us to test the respective algorithms in the presence of different effects, such as reflection shift. Moreover, it provides a ground truth that permits for quantitative evaluation. We also use the Berkeley dataset from [29], which contains 110 real image pairs (reflection and transmission layer) whose characteristics can be found in [29]. Finally, we also use a selection of ‘real-world’ images from [28], for which ground truths are not available. All measurements and reconstructions were taken from these datasets.

B. Evaluation Methodology.

We design a four-part evaluation scheme, where the evaluation protocol for each part is as follows.

(E1) The first part is a visual comparison of our method against AR17 [14]. We remark that in the case $\gamma = 0, \phi = 1$, our method reduces to that of AR17; this comparison therefore shows that the changes made to the objective function fulfil their intended purposes.

(E2) The main part of the evaluation is to compare our solution to the state-of-the-art methods. In (E2a) we compare to state-of-the-art algorithmic techniques LB14 [12], SH15 [27], AR17 [14], using FAN17 [28] as a benchmark. (E2b) is an evaluation against more recent advances in deep-learning FAN17 [28], WAN18 [30], ZHANG18 [29] and YANG18 [32]

on both real-world images and the Berkeley dataset. We present both numerical comparisons, averaged over the SIR² and Berkeley datasets in (E2a, E2b) respectively, and visual comparisons for a range of selected images from all three datasets.

(E3) We evaluate the impact of the user input, and show the results of our method with no region selection, with crude region selection and with more detailed region selection. This will justify our claim that crude region selection is sufficient to minimise loss of detail in reflection-free areas, but offers a substantial qualitative improvement on *no* region selection.

(E4) Finally, we demonstrate that, by comparison to the existing user interaction approach of Levin[11], we produce better results whilst requiring less effort from the end-user.

We address our scheme from both qualitative and quantitative points of view. The former is based on a visual inspection of the output \mathbf{T} , and the latter on the computation of three metrics: the structural similarity (SSIM) index [36], the Peak Signal-to-Noise Ratio (PSNR) and the inverted Localised Mean Squared Error (sLMSE). Explicit definition of the metrics can be found in Section VI of the Supplemental Material.

C. Parameter Selection.

For each of the approaches LB14 [12], SH15 [27] and AR17 [14], we set the parameters as described in the corresponding paper. Moreover, the comparison study was performed using the available codes from each corresponding

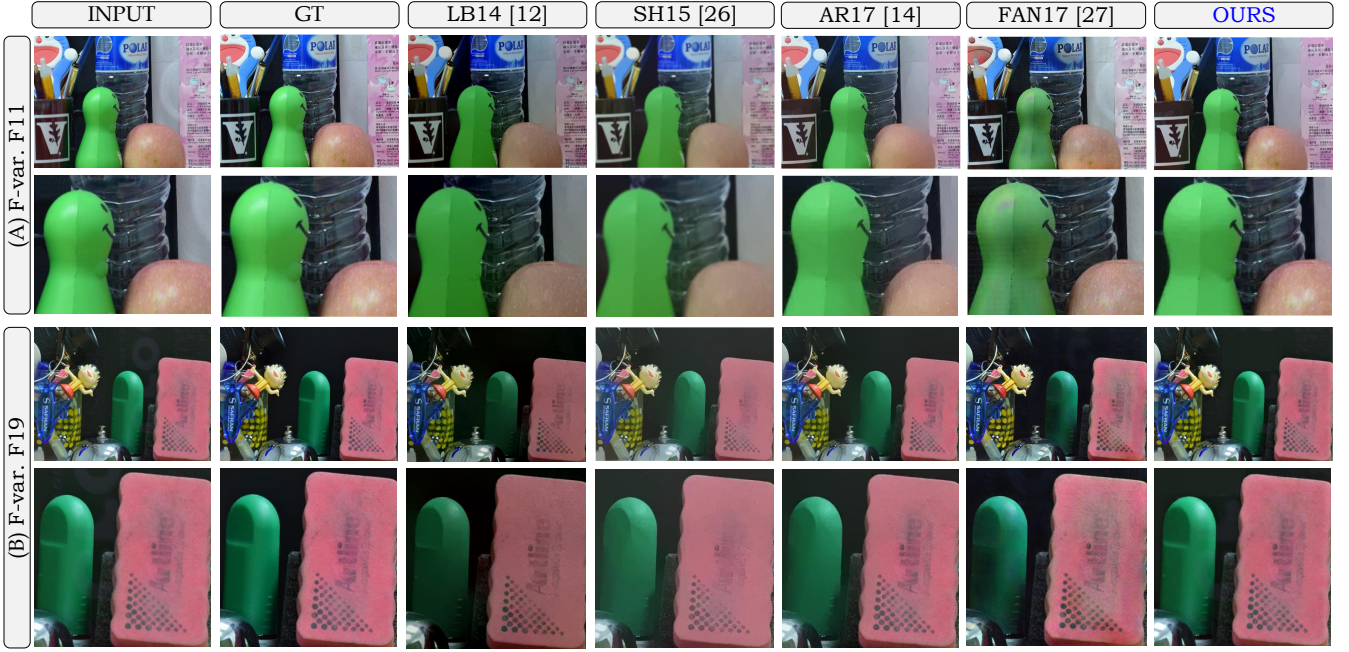


Fig. 4: (E2a). Visual comparison against the state-of-the-art of model-based approaches (including FAN17 [28] as baseline for comparison). The selected frames show variations in shape, colour and texture to appreciate the performance of the compared approaches. Overall, our approach gives a better approximation of \mathbf{T} by preserving colour and structure quality while keeping fine details. Details are better appreciated on screen.

F-var.	sLMSE			SSIM			PNSR		
	F11	F19	F32	F11	F19	F32	F11	F19	F32
LB14 [12]	0.835	0.832	0.833	0.784	0.804	0.791	21.659	21.869	21.678
SH15 [27]	0.901	0.852	0.874	0.779	0.813	0.765	21.642	22.046	21.620
AR17 [14]	0.983	0.984	0.984	0.820	0.825	0.824	22.748	22.705	22.851
FAN17 [28]	0.981	0.982	0.982	0.854	0.859	0.851	23.262	23.853	23.432
OURS	0.984	0.986	0.984	0.852	0.866	0.854	23.254	23.907	23.649

TG-var.	sLMSE			SSIM			PNSR		
	TG3	TG5	TG10	TG3	TG5	TG10	TG3	TG5	TG10
LB14 [12]	0.834	0.833	0.834	0.718	0.811	0.805	21.605	21.981	21.850
SH15 [27]	0.915	0.889	0.917	0.779	0.820	0.765	21.682	22.546	21.620
AR17 [14]	0.983	0.984	0.982	0.820	0.825	0.824	22.748	22.705	22.851
FAN17 [28]	0.981	0.981	0.981	0.850	0.852	0.852	23.415	23.403	23.470
OURS	0.984	0.984	0.984	0.846	0.851	0.861	23.374	23.421	23.507

TABLE I: (E2a). Measures averaged over all images in the solid-object dataset [35].

author. For FAN17 [28], we assumed a given trained network and with parameters set as described in that paper.

For our approach, we set the values of the ADAM method as suggested in [34]. Moreover, we set $\lambda = 2e - 3$, $\beta_{\max} = 1e5$ and $\kappa = 2$ and $\gamma = 0.012$. The choices of $\lambda, \beta_{\max}, \kappa$ follow [14] for analogous parameters, which is consistent with the reasoning in Subsection III-D. γ was chosen based on experimental results for a range of images disjoint from the test dataset, with a range of test values following the discussion in Subsection III-D

D. Results and Discussion.

We evaluate our proposed method following the scheme described in Section IV-B.

(E1). We begin by evaluating our method against AR17 [14]. We ran both approaches on the complete solid objects category of the dataset. In Fig. 3, we show four output examples with different settings (Aperture value $F=\{11, 32\}$ and thickness of glass $TG=\{3, 10\}$). Visual assessment agrees with the theory of our approach, in which we highlight the elimination of colour shifts and the preservation of the image details. Most notably, we see that our approach enforces global colour similarity and avoids blurring effects produced by the outputs of AR17 [14]; see, for example, outputs (A), (C) and (D). The detail in Fig. 3 highlights these effects, in particular in (A) the blur and colour loss effects in the *Winnie the Pooh* toy, in (C) the loss of edge details in the shirt collar (left toy) and the neck (white toy), and in (D) a blurring effect in the toy's

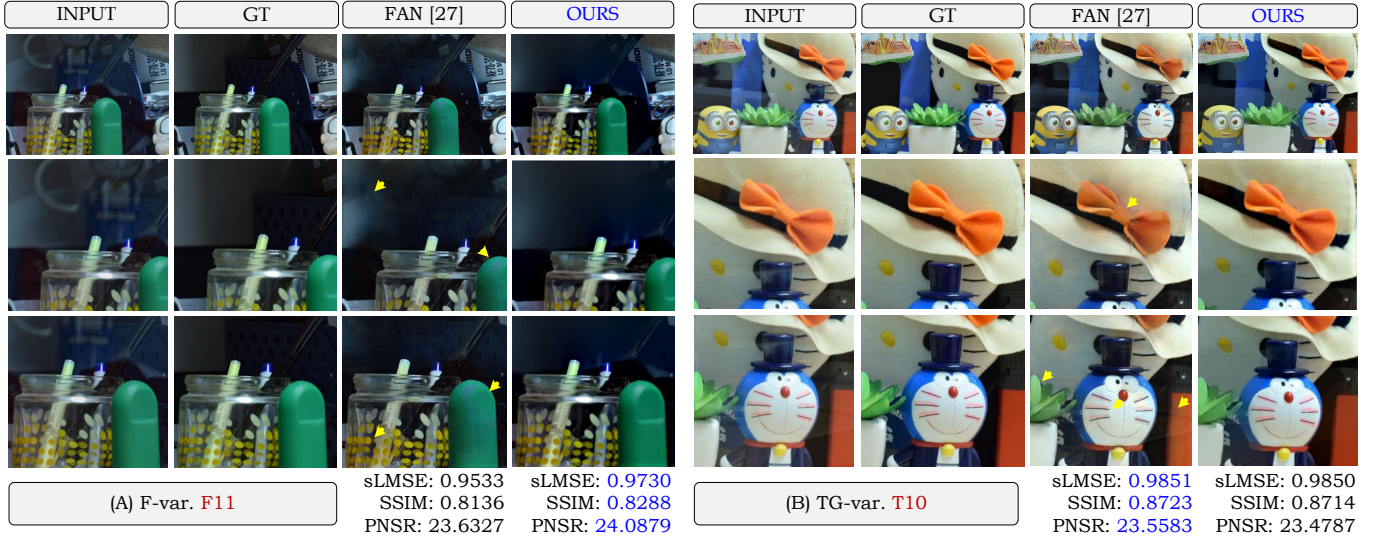


Fig. 5: (E2b). Two interesting cases in which we visually and numerically compare our approach against the work of Fan et al. [28]. We emphasise that even in cases when the metrics are higher for FAN17 [28], the output from our algorithm appears visually more appealing and natural. We highlight the false colour effects (see bow in (B)), loss of fine details (see green object in (A)) and reflection artefacts (see yellow markers in both) in the output of FAN17. Details are better appreciated on screen.

legs. In the detail of output (B), it can be seen that AR17 [14] fails to preserve the shadows and the colours of the flowers. This is further reflected in the numerical results, where our method reported higher values for the three evaluation metrics.

Overall, we noticed that often AR17 [14] fails to penalise colour shifts, due to the translation invariance of the Laplacian fidelity term. It also tends to produce blurring effects in reflection-free parts of the image, which our approach is able to prevent through our spatially aware technique.

(E2a). We now evaluate our approach against the model-based state-of-the-art methods (LB14 [12], SH15 [27], AR17 [14], and include FAN17 [28] as a baseline of comparison) using the full solid objects category of the SIR² dataset. As discussed above, we may view the results of AR17 as those of our algorithm in the special case $\gamma = 0$, and without user interaction ($\phi \equiv 1$) to evaluate the effect of these changes. We emphasise that results for our algorithm were generated with user interaction, *as this is a key part of our technique*.

We show the output of the selected methods and our proposed one for four chosen images along with the ground truth in Fig. 4. By visual inspection, we observe that outputs generated with LB14 [12] are darker than the desired output; see, for instance, the detail of (A). Moreover, LB14 fails to preserve texture and global colour similarity, as is apparent in (A) on the surface of the apple, and (B) on the pink block. By contrast, our approach was able to keep the details on both cases. Moreover, we observed that both SH15 [27] and AR17 [14] tend to have a noticeable colour shift and a significant loss of structure; as is visible on (B) the green pole. In particular, we highlight the green pole in (B), in which only our approach was clearly able to maintain the fine details.

We observe that the deep learning based solution FAN17 [28] shows good edge preservation, but often fails to correctly reproduce colour and texture, and produces notice-

able artefacts. This will be discussed further in (E2b). Overall, out of the evaluated model-based single-image reflection removal techniques, our approach consistently yields the most visually pleasing results. These observations are confirmed by further examples in Section II of the Supplemental Material.

For a more detailed quantitative analysis, we report the global results in Table I. The displayed numbers are the average of the image metrics across the whole body of ‘solid-object’ files in the dataset, in order to understand the general behaviour and performance of the algorithms.

We observe that both AR17 [14] and our approach outperform the remaining algorithms with respect to sLMSE. With respect to SSIM and PNSR, we also achieve significant improvements over most state-of-the-art techniques, most notably over the similar technique AR17 [14]. The only other approach evaluated here which performs similarly well is the deep learning approach FAN17 [28]. As was discussed above, a closer look at single images shows occasional difficulties of this approach, and the more reliable performance of our model-based method.

(E2b). Having extensively compared our new method to model-based approaches in (E2a), we now present a detailed comparison against recent advances in single-image reflection removal based on deep-learning. We compare against FAN17 [28], WAN18 [30], ZHANG18 [29] and YANG18 [32] on both the Berkeley dataset and real-world images.

Having used FAN17 [28] as a benchmark for comparison in (E2a), we first present a further comparison of this method against our technique. Indeed, from Table I, it may appear that FAN17 produces output of a similar quality to our technique. However, we notice that the outputs displayed in Fig. 4 suggest that our method produces visually nicer results; to validate this, we present further experiments in Fig. 5. The images displayed are two cases from the SIR² dataset, in which we

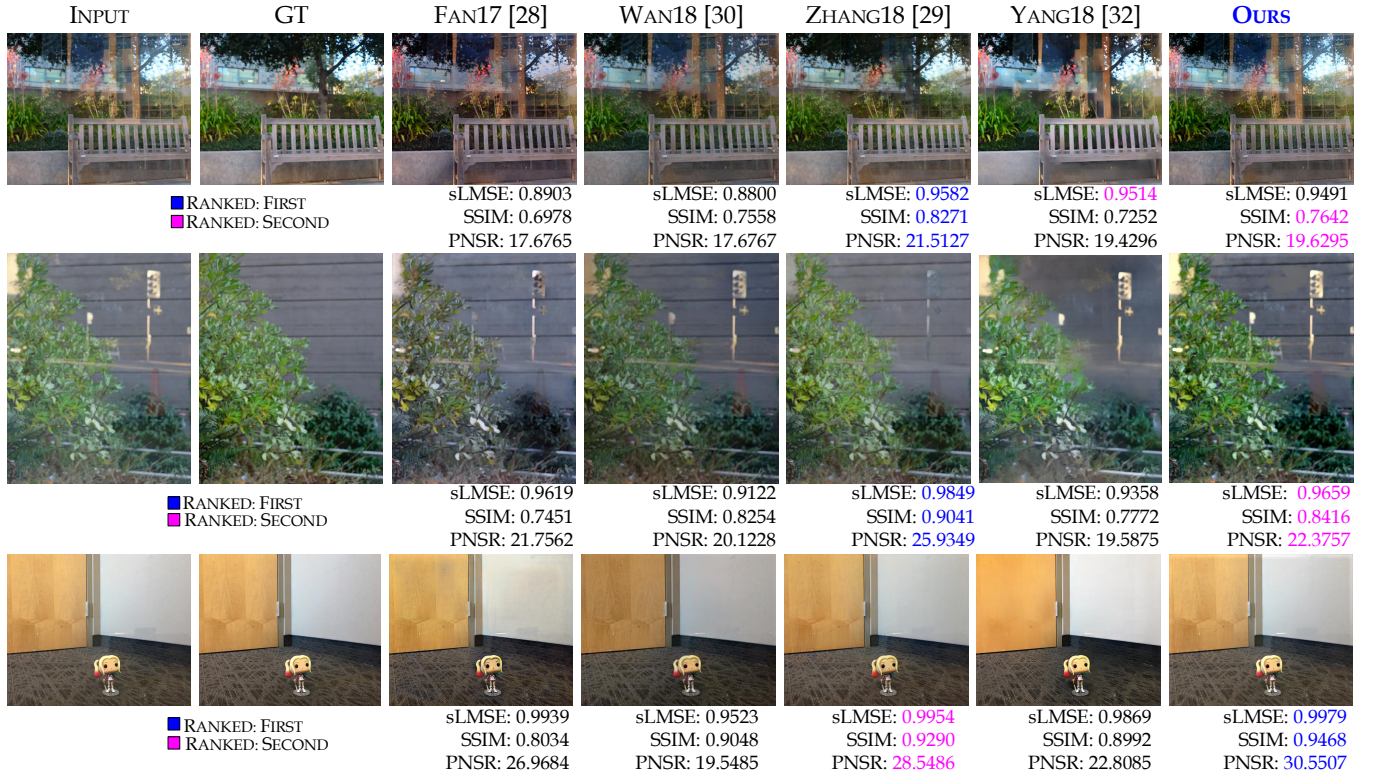


Fig. 6: (E2b). Visual and numerical comparison of our technique vs. Deep-learning techniques on a selection of images from the Berkeley dataset. Details are better appreciated on screen.

THE BERKLEY DATASET					
	FAN17 [28]	WAN18 [30]	ZHANG18 [29]	YANG18 [32]	OURS
sLMSE	0.8407	0.8090	0.8638	0.8398	0.8647
SSIM	0.7022	0.6982	0.7923	0.6911	0.7315
PNSR	18.2989	18.300	21.6203	17.8673	18.7833
	RANKED FIRST		RANKED SECOND		

TABLE II: (E2b). Numerical comparison of our technique vs. Deep-learning techniques for the entire Berkeley dataset. The numerical values are computed as the averages of the similarity metrics over all images.

observe difficulties similar to those in Fig. 4. In Figs. 4A, 5A, FAN17 has wrongly identified a specular reflection in the transmitted layer as belonging to the reflected layer, producing unpleasant artefacts. We also highlight incomplete reflection removal in the examples in Fig. 5, false-colour effects in Figs. 4 and 3B, and unwanted colour flattening in Fig. 5A.

Next, we present a visual comparison of a selection of images from the Berkeley dataset in Fig. 6. The images include the values of the similarity metrics compared to the ground truth in each case. We observe that FAN17 [28], WAN18 [30] suffer from poor colour retention in these test images, while YANG18 [32] induces a significant amount of blurring (see the door in the bottom picture, and the edges of the plant in the middle one). ZHANG18 [29] performs very well both visually and numerically, although the quality of its performance somewhat decreases when compared on a different dataset as we do in Fig. 7. Our method readily competes with ZHANG18 [29] in terms of similarity metrics, but also is able to preserve structure and color much better than

the remaining approaches, while still removing a comparable amount of the reflections.

In Table II we present the similarity measures which are computed as the average over all images in the Berkeley dataset. With respect to sLMSE, our method outperforms all other techniques, in particular FAN17 [28], WAN18 [30] and YANG18 [32] by a significant margin. With respect to SSIM and PNSR, our method performs similarly well, and places second behind ZHANG18 [29].

Finally, we test all of the DL methods on a selection of real-world images in Fig. 7. We observe that most of the competing methods suffer from poor colour preservation, which is especially visible in ZHANG18 [29] with respect to the skin colour in middle and upper image, and incomplete removal of the reflections. In FAN17 [28] especially we notice the introduction of artefacts on the arms in the top picture and nearby the head in the bottom one. Our method, while not completely removing the reflections, still ensures good preservation of colour and important structure, and hence

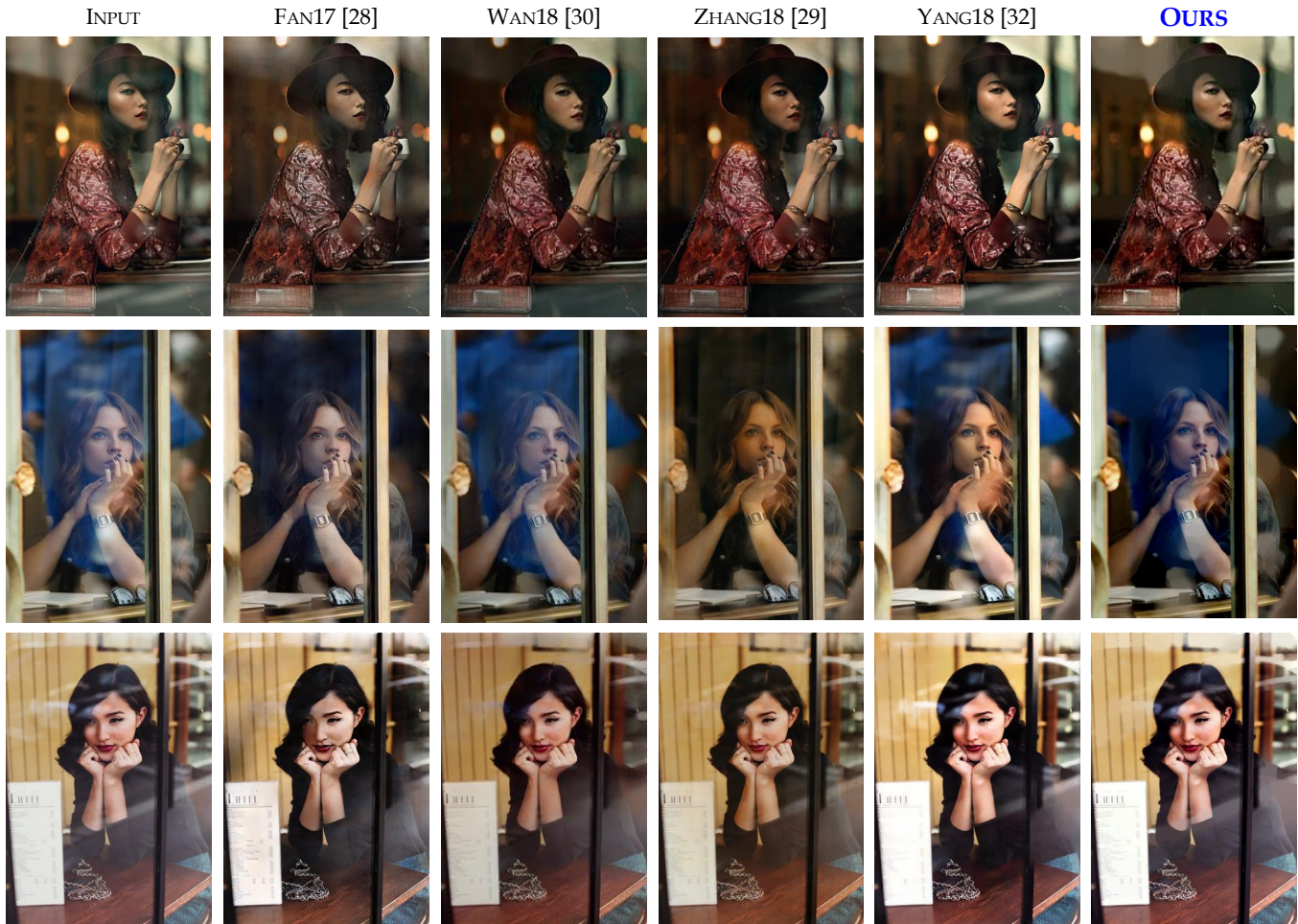


Fig. 7: (E2b). Comparison of our technique vs Deep-Learning techniques on real-world images. We note that our technique is able to suppress the reflections while avoiding flattening in the skin tone and avoiding false-colour effects. This is an example of our motivation in Observation 2: colour flattening on the skin is much more noticeable than the same effect on the props. Images are from the real-world dataset [28] and no ground truths are available.

in terms of output quality readily competes with the deep-learning based methods. Additional experiments, which further validate this conclusion, may be found in Section III of the Supplemental Material.

The above comparison against the most recent deep learning approaches for this problem, demonstrates that *at this point in time, our model-based method readily competes with deep learning in terms of output quality*. The authors note that traditionally, deep learning has achieved ground breaking success in tasks involving labelling or classification [37], [38]. The good visual results generated by deep network usually benefit from the statistical information covered in the large body of training samples. However, a plain fully convolutional neural network does not impose the same kind of rigid and intuitive constraints as model-based approaches; for example, piecewise smoothness is not enforced. Such a limitation in the deep network results in inconsistent reflection removal within a single image, as seen in Figs 5, 6, 7. While in this paper the deep-learning based techniques provide an important benchmark, their classification as ‘single-image’ techniques

raises definitional issues that might be interesting for the community to discuss. This discussion can be found in Section V of the Supplemental Material.

(E3). In Fig. 8, we analyse the impact of the user-interaction, again including FAN17 [28] as a baseline for comparison. In the first subfigure, we present the results of our approach without region selection, and with both crude and detailed region selection. Without region selection, there is noticeable blurring and flattening: see, for example, the green object in the first example and the apple in the second. However, even with very crude region selection, our technique is able to mitigate these to produce a visually better result. In the second subfigure, we show the result of our technique with and without region selection on two examples from the real-world dataset where region selection makes a substantial visual difference to the output. In both cases, without region selection, the output has a lot of colour flattening on the skin of the model, leading to a very unnatural and unrealistic output. We therefore conclude that *even very crude selection of the reflection regions results in good reflection removal, and*

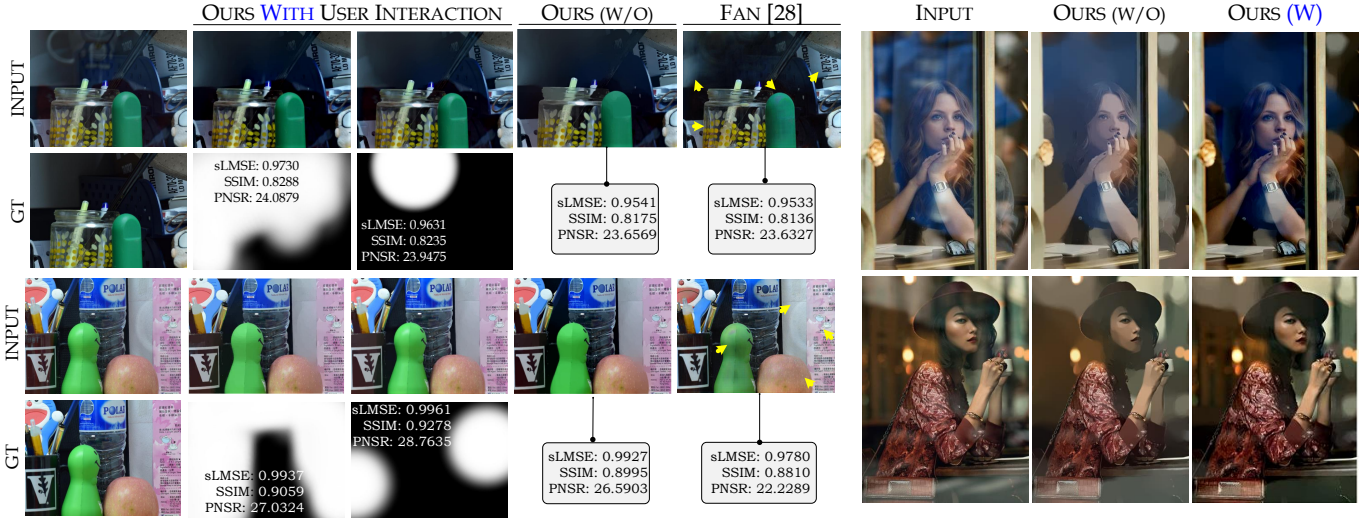


Fig. 8: (E3). From left to right: The impact of the user-interaction on the outputs computed by OUR approach (with and without user interaction), with FAN[28] as a benchmark. Examples of cases where region selection leads to noticeable qualitative improvements in avoiding flattening.

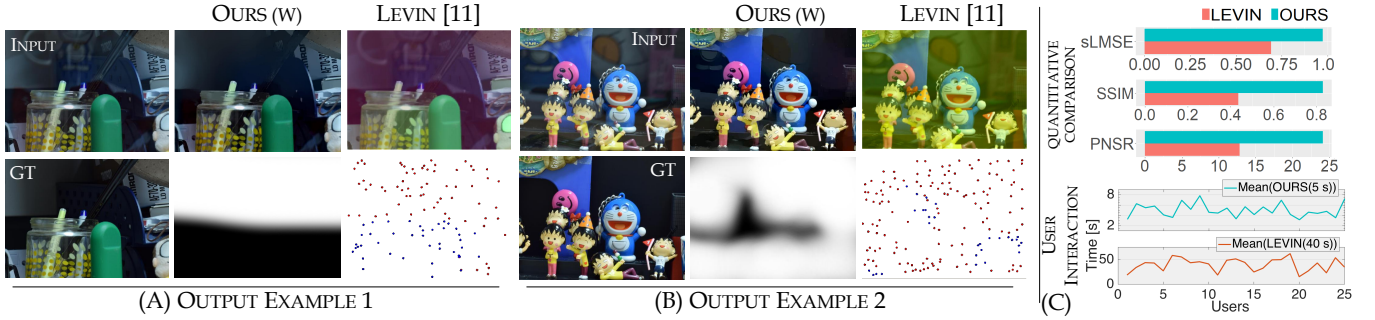


Fig. 9: (E4). (A-B): Visual comparison of the user-interaction schemes in LEVIN [11] and OURS based on a specific example (C): Quantitative comparison of the two schemes on the solid object corpus of the SIR² dataset, and user-interaction time based on a selection of images from this dataset.

that crude region selection noticeably improves on no region selection. This justifies our claim of a providing a simple and effective user-interaction scheme.

(E4). We also compare our method to the existent user-interaction by Levin [11]. We demonstrate that in comparison, our method produces qualitatively and quantitatively better results, while requiring significantly less effort from the end-user. This underlines one of the main messages of this paper, that *we provide a simple user-interaction method, which gives a significant improvement in the quality of the output.*

In Fig. 9 we compare the amount of user interaction required and the quality of the resulting output for both methods. Firstly, in the bottom half of (A-B), the user-interaction for both methods is shown. For our method, the user is asked to determine the location of reflections in the image by marking the rough location in white; several examples of this user-selection are provided in Section IV of the Supplemental Material. In Levin’s approach, the user is asked to select foreground gradients in red and background gradients in blue. We can also see the corresponding output of the algorithm, which can be visually observed to be significantly improved

using our method.

In Fig. 9 (C) we compare the specific effort of user-interaction between Levin [11] and our proposed method. For this we asked a group of 25 colleagues to perform the user-interaction on both schemes and try to achieve the best quality removal as quickly as possible. We observe that, on average, our approach took our colleagues around 5 seconds per image, while Levin’s method required around 40 seconds, an increase of around 700%. The corresponding quantitative results can be seen in the upper half of Fig. 9 (C). The numerical values are the metrics averaged over the entire output from 25 users working on the solid-object dataset. In particular each user was given 6 different settings (3 types of focus and 3 types of thickness) of reflections for each of the 20 images in the dataset, and was then asked to perform the user selection for both methods. We see that the similarity metrics are significantly improved using our new method. This shows that our method *requires significantly less effort from the end user than other existent approaches*, while at the same time significantly improving the quality of reflection removal.

V. CONCLUSIONS

This paper addresses the challenging problem of single image reflection removal. We propose a technique in which two novelties are introduced to provide reflection removal of higher quality. The first is an *spatially aware prior term, exploiting low-level user interaction*, which tailors reflection suppression to preserve detail in reflection-free areas. The second is an H^2 *fidelity term*, which combines advantages of both L^2 and Laplacian fidelity terms, and promotes better reconstruction of faithful and natural colours. Together, these result in better preservation of structure, detail and colour. We demonstrate the potential of our model through quantitative and qualitative analyses, in which it produces better results than all tested model-based approaches and readily competes with recent deep learning techniques. Future work might include the use of deep learning techniques to automatically select regions, which would avoid the need for user interaction, while preserving many of the advantages of our technique.

ACKNOWLEDGMENT

This work was supported by the UK Engineering and Physical Sciences Research Council (EPSRC) grant EP/L016516/1 for the University of Cambridge Centre for Doctoral Training, the Cambridge Centre for Analysis. Support from the CMIH University of Cambridge is greatly acknowledged.

REFERENCES

- [1] Y. Y. Schechner, N. Kiryati, and R. Basri, "Separation of transparent layers using focus," *International Journal of Computer Vision (IJCV)*, vol. 39, no. 1, pp. 25–39, 2000.
- [2] A. Agrawal, R. Raskar, S. K. Nayar, and Y. Li, "Removing photography artifacts using gradient projection and flash-exposure sampling," *ACM Transactions on Graphics (TOG)*, vol. 24, no. 3, pp. 828–835, 2005.
- [3] N. Kong, Y.-W. Tai, and J. S. Shin, "A physically-based approach to reflection separation: from physical modeling to constrained optimization," *IEEE Transactions on Pattern Analysis and Machine Intelligence (PAMI)*, vol. 36, no. 2, pp. 209–221, 2014.
- [4] A. Lakhtakia, "General schema for the brewster conditions," *Optik*, vol. 90, no. 4, pp. 184–186, 1992.
- [5] H. Farid and E. H. Adelson, "Separating reflections and lighting using independent components analysis," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, vol. 1, 1999, pp. 262–267.
- [6] N. Kong, Y.-W. Tai, and S. Y. Shin, "High-quality reflection separation using polarized images," *IEEE Transactions on Image Processing*, vol. 20, no. 12, pp. 3393–3405, 2011.
- [7] R. Szeliski, S. Avidan, and P. Anandan, "Layer extraction from multiple images containing reflections and transparency," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2000, pp. 246–253.
- [8] B. Sarel and M. Irni, "Separating transparent layers through layer information exchange," *European Conference on Computer Vision (ECCV)*, pp. 328–341, 2004.
- [9] S. N. Sinha, J. Kopf, M. Goesele, D. Scharstein, and R. Szeliski, "Image-based rendering for scenes with reflections," *ACM Transactions on Graphics (TOG)*, vol. 31, no. 4, pp. 100–1, 2012.
- [10] X. Guo, X. Cao, and Y. Ma, "Robust separation of reflection from multiple images," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2014, pp. 2187–2194.
- [11] A. Levin and Y. Weiss, "User assisted separation of reflections from a single image using a sparsity prior," *IEEE Transactions on Pattern Analysis and Machine Intelligence (PAMI)*, vol. 29, no. 9, 2007.
- [12] Y. Li and M. S. Brown, "Single image layer separation using relative smoothness," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2014, pp. 2752–2759.
- [13] R. Wan, B. Shi, T. A. Hwee, and A. C. Kot, "Depth of field guided reflection removal," in *IEEE International Conference on Image Processing (ICIP)*, 2016, pp. 21–25.
- [14] N. Arvanitopoulos Darginis, R. Achanta, and S. Süsstrunk, "Single image reflection suppression," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.
- [15] H. Barrow and J. Tenenbaum, "Recovering intrinsic scene characteristics," *Comput. Vis. Syst.*, vol. 2, 1978.
- [16] Y. Y. Schechner, J. Shamir, and N. Kiryati, "Polarization and statistical analysis of scenes containing a semireflector," *Journal of the Optical Society of America (JOSA)*, vol. 17, no. 2, pp. 276–284, 2000.
- [17] N. Kong, Y.-W. Tai, and S. Y. Shin, "A physically-based approach to reflection separation," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2012, pp. 9–16.
- [18] Y. Y. Schechner, N. Kiryati, and J. Shamir, "Blind recovery of transparent and semireflected scenes," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, vol. 1, 2000, pp. 38–43.
- [19] Y. Li and M. S. Brown, "Exploiting reflection change for automatic reflection removal," in *IEEE International Conference on Computer Vision (ICCV)*, 2013, pp. 2432–2439.
- [20] C. Sun, S. Liu, T. Yang, B. Zeng, Z. Wang, and G. Liu, "Automatic reflection removal using gradient intensity and motion cues," in *ACM Multimedia Conference*, 2016, pp. 466–470.
- [21] T. Xue, M. Rubinstein, C. Liu, and W. T. Freeman, "A computational approach for obstruction-free photography," *ACM Transactions on Graphics (TOG)*, vol. 34, no. 4, p. 79, 2015.
- [22] A. Nandoriya, M. Elgharib, C. Kim, M. Hefeeda, and W. Matusik, "Video reflection removal through spatio-temporal optimization," in *IEEE International Conference on Computer Vision (ICCV)*, 2017, pp. 2430–2438.
- [23] S. M. A. Shah, S. Marshall, and P. Murray, "Removal of specular reflections from image sequences using feature correspondences," *Machine Vision and Applications*, vol. 28, no. 3-4, pp. 409–420, 2017.
- [24] C. Simon and I. K. Park, "Reflection removal for in-vehicle black box videos," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015, pp. 4231–4239.
- [25] J. Y. Cheong, C. Simon, C.-S. Kim, and I. K. Park, "Reflection removal under fast forward camera motion," *IEEE Transactions on Image Processing*, vol. 26, no. 12, pp. 6061–6073, 2017.
- [26] B.-J. Han and J.-Y. Sim, "Glass reflection removal using co-saliency based image alignment and low-rank matrix completion in gradient domain," *IEEE Transactions on Image Processing*, 2018.
- [27] Y. Shih, D. Krishnan, F. Durand, and W. T. Freeman, "Reflection removal using ghosting cues," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015, pp. 3193–3201.
- [28] Q. Fan, J. Yang, G. Hua, B. Chen, and D. Wipf, "A generic deep architecture for single image reflection removal and image smoothing," *IEEE International Conference on Computer Vision (ICCV)*, 2017.
- [29] X. Zhang, R. Ng, and Q. Chen, "Single image reflection separation with perceptual losses," *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
- [30] R. Wan, B. Shi, L.-Y. Duan, A.-H. Tan, and A. C. Kot, "CRRN: Multi-scale guided concurrent reflection removal network," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 4777–4785.
- [31] M. Jin, S. Süsstrunk, and P. Favaro, "Learning to see through reflections," in *Computational Photography (ICCP), 2018 IEEE International Conference on*. IEEE, 2018, pp. 1–12.
- [32] J. Yang, D. Gong, L. Liu, and Q. Shi, "Seeing deeply and bidirectionally: a deep learning approach for single image reflection removal," in *European Conference on Computer Vision*. Springer, Cham, 2018, pp. 675–691.
- [33] L. Xu, C. Lu, Y. Xu, and J. Jia, "Image smoothing via l0 gradient minimization," in *ACM Transactions on Graphics (TOG)*, vol. 30, no. 6. ACM, 2011, p. 174.
- [34] D. Kingma and J. Ba, "Adam: A method for stochastic optimization," *International Conference for Learning Representations (ICLR)*, 2014.
- [35] R. Wan, B. Shi, L.-Y. Duan, A.-H. Tan, and A. C. Kot, "Benchmarking single-image reflection removal algorithms," in *IEEE International Conference on Computer Vision (ICCV)*, 2017, pp. 3942–3950.
- [36] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, "Image quality assessment: from error visibility to structural similarity," *IEEE Transactions on Image Processing*, vol. 13, no. 4, pp. 600–612, 2004.
- [37] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in neural information processing systems (NIPS)*, 2012, pp. 1097–1105.
- [38] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, vol. 521, no. 7553, p. 436, 2015.

SUPPLEMENTAL MATERIAL FOR THE ARTICLE:

Mirror, Mirror, on the Wall, Who's Got the Clearest Image of Them All? A Tailored Approach to Single Image Reflection Removal

Daniel Heydecker*, Georg Maierhofer*, Angelica I. Aviles-Rivero*
Qingnan Fan, Dongdong Chen, Carola-Bibiane Schönlieb and Sabine Süsstrunk

VI. OUTLINE

This document extends the practicalities and visual results presented in the main paper in order to show further details of our approach and experiments. This is structured as follows.

- **Section II:** We offer further visual results of our and the state-of-the-art model-based approaches using the SIR² dataset. As in the main paper, we will show that our technique is able to perform noticeably better than other model-based approaches.
- **Section III:** We give further visual comparison of our technique against competing Deep-Learning techniques using the Berkeley dataset [29]. This further validates our claim of being able to compete with the DL approaches.
- **Section IV:** We show how our user-interaction technique may be implemented in practice, and demonstrate examples generated by the authors.
- **Section V:** Continuing a point raised in the main text, we discuss a definitional issue over whether competing DL techniques can be considered ‘single-image’.
- **Section VI:** In the interests of clarity and completeness, we give an explicit definition and motivation of the metrics used for the quantitative analyses.

VII. SUPPLEMENTARY VISUAL RESULTS WITH SIR² DATASET

In this section, we extend the comparison of visual results of Fig. 4 from the main paper. The comparison includes LB14 [12], SH15 [27], AR17 [14], using FAN17 [28] as a benchmark.

In Fig. 10 shows four further examples from the solid object part of the SIR² dataset and we note that amongst these methods, ours presents the most visually appealing results. In particular we note that LB14 [12] suffers from colour shift at the fan in the third image, and SH15 [27] and AR17 [14] suffer significant loss of structure in the third and fourth images, downsides that are not observed in our technique. Fan on the other hand removes a significantly smaller portion of the reflections in these images – and this incomplete removal can also be noticed in the further detailed comparison of FAN [28] and ours in Fig. 11.

VIII. FURTHER VISUAL RESULTS OF OUR AND DL-BASED APPROACHES

In addition to the experimental results displayed in the main paper, we present some further examples of our technique vs

competing DL techniques on elements of the Berkeley dataset [29], displayed in Fig. 12.

We see that many of the output images for the competing DL approaches suffer from the same problems described in the main text. The results of FAN17 [28] introduce very visible artefacts in images 1, 4 and 5, which often make the reflections *more* visible than in the input image, and have false-colour effects in images 3 and 7. YANG18 [32] has substantial blurring and loss of detail, which is visible on the carpet in images 4 and 7 and the sign in image 5, and WAN18 [30] introduces substantial blurring throughout. As in the main text, ZHANG18 [29] usually performs extremely well on this dataset, but the outputs of images 3 and 7 display noticeable and unpleasant false colour effects. By contrast, our output is able to mitigate the loss of detail and false-colour effects, while competing with the deep-learning approaches in suppressing the reflection layer.

IX. SUPPLEMENTARY VISUAL RESULTS OF THE USER INTERACTION

In Fig. 13 we display a number of examples of the user selection in practise. In particular we display a range of images from all our datasets and show the user-selection as performed in the experiments. In the graph the input images are shown together with the corresponding region selection: Here the user selects a region to be white, if a reflection is seen in that part of the image and black otherwise. This information is then translated (as the relative gray values) into the region selection function ϕ as described in section II.A of the main paper.

X. IS IT ‘SINGLE IMAGE’?

The need in Deep-Learning for a large set of training data raises the definitional issue of whether the technique could be considered as a ‘single-image’ technique. The definitions of what constitutes a single-image technique reads:

Definition 1. A *Single-Image technique* is one which uses the information from a single input \mathbf{Y} to extract the transmission layer \mathbf{T} .

In practical terms, this leads to disadvantages similar to those of multiple image techniques; as from a mathematical point of view, DL has the same goal than the multiple-image case: to reduce the strongly ill-posedness problem created in the single-image case.

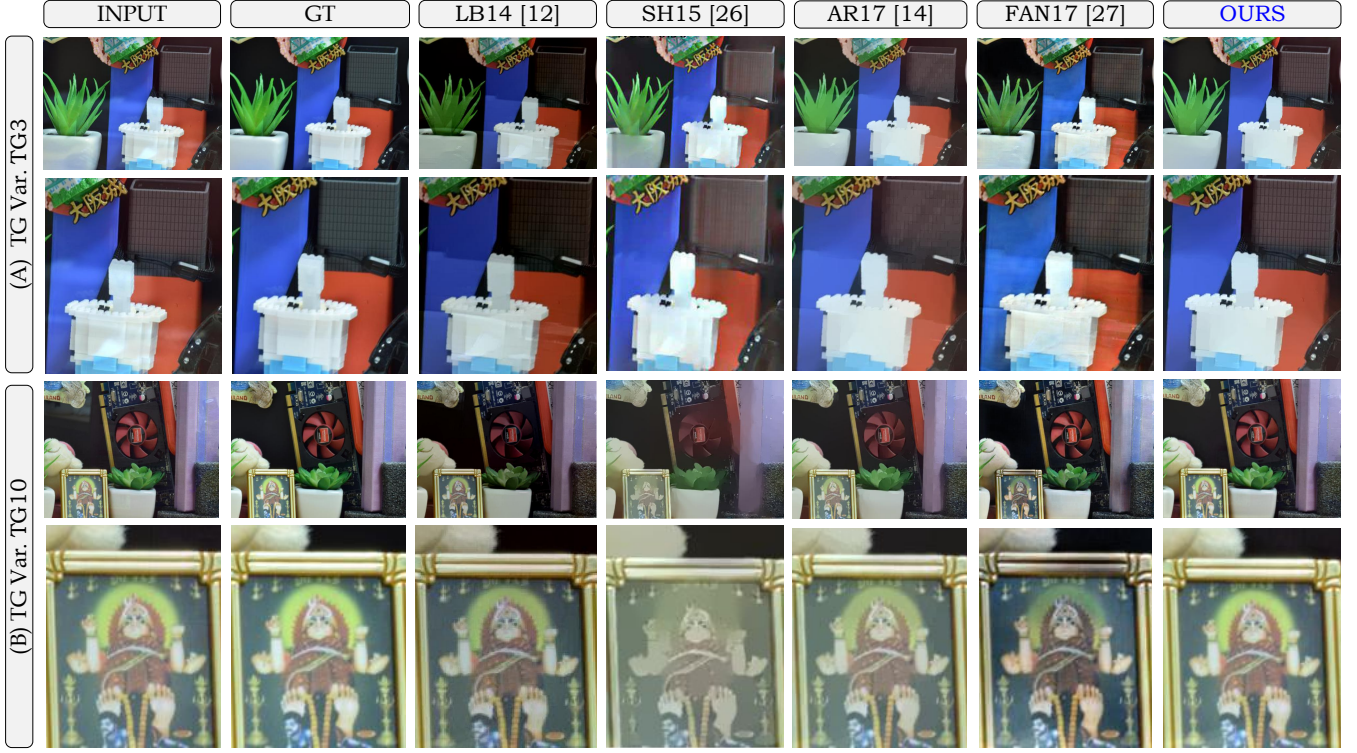


Fig. 10: Visual comparison against the state-of-the-art of model-based approaches (including FAN17 [28] as baseline for comparison). The selected frames show variations in shape, colour and texture to appreciate the performance of the compared approaches. Overall, our approach gives a better approximation of \mathbf{T} by preserving colour and structure quality while keeping fine details. Details best appreciated on screen.

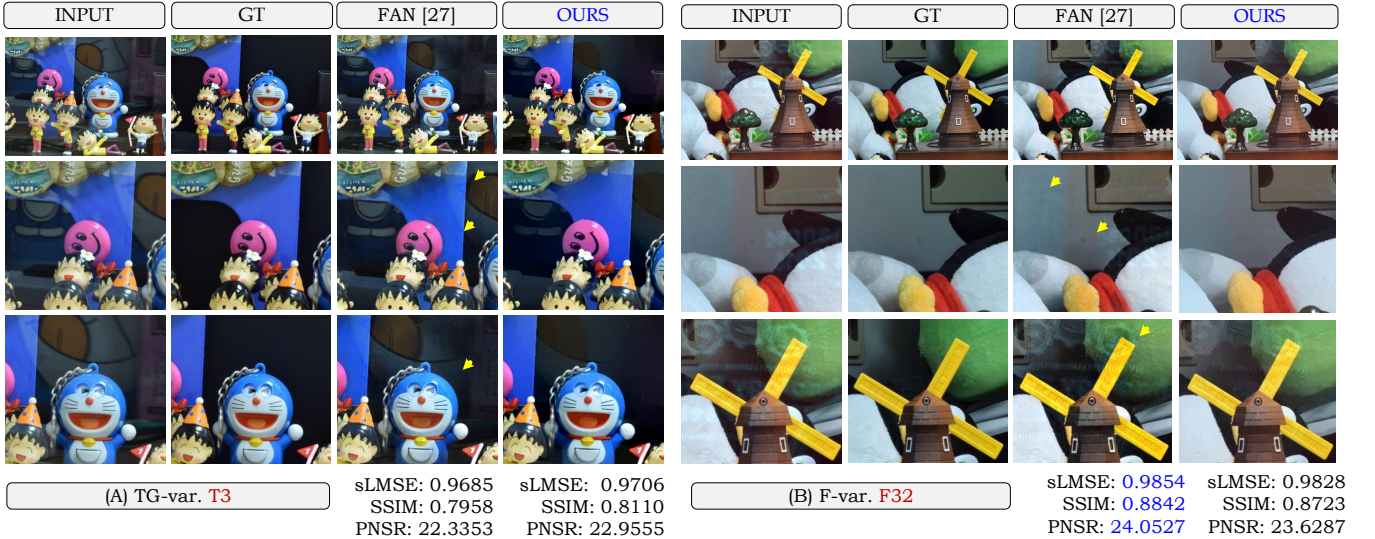


Fig. 11: Two interesting cases in which we visually and numerically compare our approach against the work of Fan et al. [28]. We emphasise that even in cases when the metrics are higher for FAN17 [28], the output from our algorithm appears visually more appealing and natural. Details best appreciated on screen.

XI. DEFINITIONS OF THE METRICS

For clarification purposes, we explicitly define the specific form of the three metrics used in our comparison study. It is particularly interesting since the results can differ from the

ones reported in [35] due to the different forms of the metrics. It is to be noted that not all used metrics are explicitly defined in [35], leading to some ambiguity concerning the specific form of, for example, LMSE that is used. Our metrics are

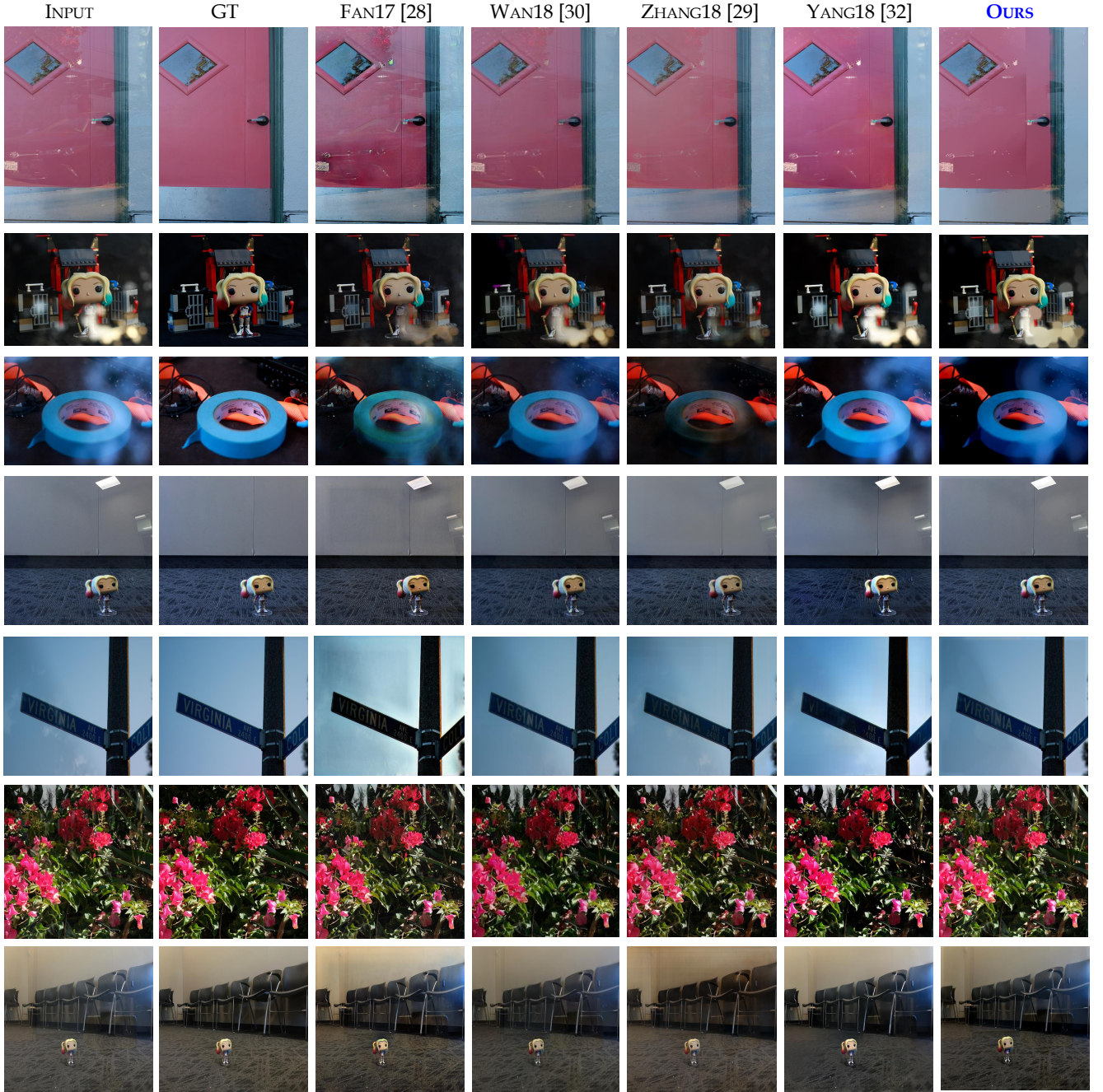


Fig. 12: Further visual comparison of our technique vs deep-learning techniques, on images from the Berkely dataset [29]. Details best appreciated on screen.

computed as follows:

- sLMSE (Inverted Localised Mean Squared Error) is computed as follows: Let S be an approximation of some ground truth \hat{S} . We compute the LMSE as the MSE over patches S_ω of size 20×20 , shifted by 10 each stage, such that

$$LMSE(S, \hat{S}) = \sum_{\omega} \|S_\omega - \hat{S}_\omega\|_2^2 \quad (16)$$

Then we normalise and produce an inverted measure such that the error measure is 1 if the approximation is good,

and zero otherwise.

$$sLMSE(S, \hat{S}) = \frac{LMSE(S, \hat{S})}{LMSE(S, 0)}. \quad (17)$$

- SSIM (Structural Similarity Index) is computed in the standard way. Let again S be an approximation of some ground truth \hat{S} , and let $\mu_S, \mu_{\hat{S}}, \sigma_S, \sigma_{\hat{S}}, \sigma_{S\hat{S}}$ be the averages, variances and covariance of S and \hat{S} respectively.

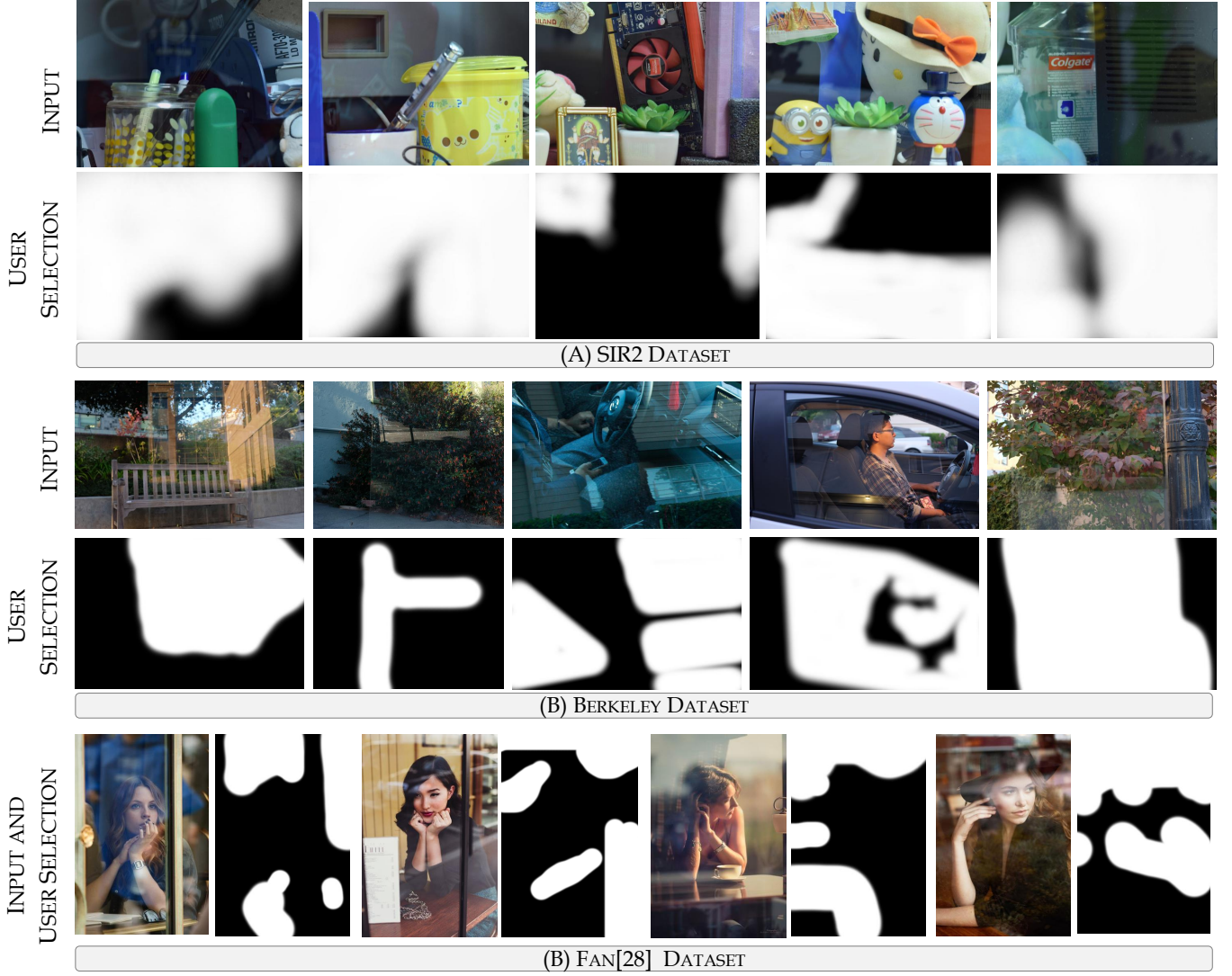


Fig. 13: Demonstration of the region selection performed by users on a number of images from the three relevant datasets. We show the input and corresponding user selection, these selections were also used in the main experiments of the paper.

Then the SSIM is calculated as

$$SSIM(S, \hat{S}) = \frac{(2\mu_S\mu_{\hat{S}} + c_1)(2\sigma_{S\hat{S}} + c_2)}{(\mu_S^2 + \mu_{\hat{S}}^2 + c_1)(\sigma_S^2 + \sigma_{\hat{S}}^2 + c_2)}. \quad (18)$$

Here, c_1, c_2 are variables to stabilise the division in case of weak denominator. In our implementation, these are chosen to be:

$$c_1 = (0.01 * L)^2 \quad (19)$$

$$c_2 = (0.03 * L)^2 \quad (20)$$

with L being a dynamic range variable that depends on the class of the image (e.g. $L = 1$ for type single images).

- PSNR (Peak Signal-to-Noise Ratio) is also computed in the standard way. Let again S be an approximation of some ground truth \hat{S} . Firstly, the full MSE is computed via

$$MSE(S, \hat{S}) = \frac{1}{N} \|S_\omega - \hat{S}_\omega\|_2^2 \quad (21)$$

where N is the number of total pixels in S . We then compute the PSNR as follows:

$$PSNR(S, \hat{S}) = 10 \log_{10} \left(\frac{MAX_I^2}{MSE(S, \hat{S})} \right), \quad (22)$$

where MAX_I is the maximal possible pixel value in the images S, \hat{S} (e.g. 255 for 8-bit images).