

Focal Stack Photography: High-Performance Photography with a Conventional Camera

Kiriakos N. Kutulakos
 Dept. of Computer Science
 University of Toronto
 kyros@cs.toronto.edu

Samuel W. Hasinoff
 MIT CSAIL
 hasinoff@csail.mit.edu

Abstract

We look at how a seemingly small change in the photographic process—capturing a focal stack at the press of a button, instead of a single photo—can boost significantly the optical performance of a conventional camera. By generalizing the familiar photographic concepts of “depth of field” and “exposure time” to the case of focal stacks, we show that focal stack photography has two performance advantages: (1) it allows us to capture a given depth of field much faster than one-shot photography, and (2) it leads to higher signal-to-noise ratios when capturing wide depths of field with a restricted exposure time. We consider these advantages in detail and discuss their implications for photography.

1 Introduction

Despite major advances in digital photography in recent years, today’s cameras are identical to their film-based predecessors in terms of basic optics: they rely on the same three controls (aperture, focus and exposure time) and their optical performance is limited by the same fundamental constraints between aperture, depth of field, exposure time, and exposure level that govern traditional photography.

As a first step in pushing the performance limits of modern cameras further, this paper looks at how a seemingly small change in the photographic process—capturing a *focal stack* at the press of a button, instead of a single photo—can boost significantly the optical performance of a conventional camera. We show that cameras with this functionality have two performance advantages (Figures 1 and 2):

- they can capture a given depth of field much faster than one-shot photography allows, and
- they achieve higher signal-to-noise ratios when capturing wide depths of field at high speed.

A focal stack is a sequence of photos taken with distinct focus settings and possibly different apertures. Just like an individual photo has a well-defined exposure time and a well-defined depth of field (*i.e.*, a range of distances where subjects appear in focus), a focal stack can be thought of as having an exposure time and a depth of field too: its exposure time is simply the total time it takes to capture all photos in the stack, and its depth of field (DOF) is the union of DOFs of all these photos.

A few digital cameras already offer a rudimentary type of “focal stack photography” via a focus bracketing mode [1]. In that mode, lens focus is adjusted by a fixed increment after each shot in a rapid, multi-shot sequence. While this mode may sometimes confer an advantage, here we explore performance gains from a much more general ability: we envision a camera control system where the photographer sets the desired DOF and exposure level (or exposure time), presses the shutter release, and the camera captures the *optimal* focal stack for those settings. Depending on context, optimality can be expressed in terms of speed (*i.e.*, shortest exposure time for a given exposure level), image quality (*i.e.*, highest signal-to-noise ratio for a given exposure time), or both.

Focal stacks are certainly not a new concept. They have an especially long history in microscopy and macro photography, where lenses have very narrow DOFs [25, 20, 6, 18]. In these applications, producing an extended-DOF photo of a sample often requires capturing and merging a large focal stack [23, 21, 6]. Focal stacks are also a rich source of 3D shape information and have been used extensively for shape-from-focus and shape-from-defocus computations in computer vision [15, 7, 9, 11]. Nevertheless, we are not aware of work that has studied focal stack photography as a high-performance alternative to capturing just one photo.

Although the difference between focal stack and one-shot photography may appear superficial, focal stacks differ from traditional photos in two important ways.

First, the relations between aperture, depth of field, exposure time and signal-to-noise ratio that apply to individual photos *do not* apply to focal stacks. For example, we show in Section 3.1 that focal stacks with a given DOF and a given exposure level can be captured much faster than a single photo with the same specifications (Figure 1). In essence, focal stacks allow us to “break” some of the basic barriers imposed by lens optics that constrain one-shot photography.

Second, while one-shot photography gives us a readily-viewable photo with the desired specifications (DOF, exposure level, etc.), focal stacks require further processing: since no single photo in a stack spans a user-specified DOF completely, its photos must be merged (and perhaps even restored) to produce a one-shot equivalent, all-in-focus image. Fortunately, these merging and restoration problems have been well studied [26, 6, 18, 17] and have a very useful side-effect: they enable 3D reconstruction of the subject being photographed [6, 18].

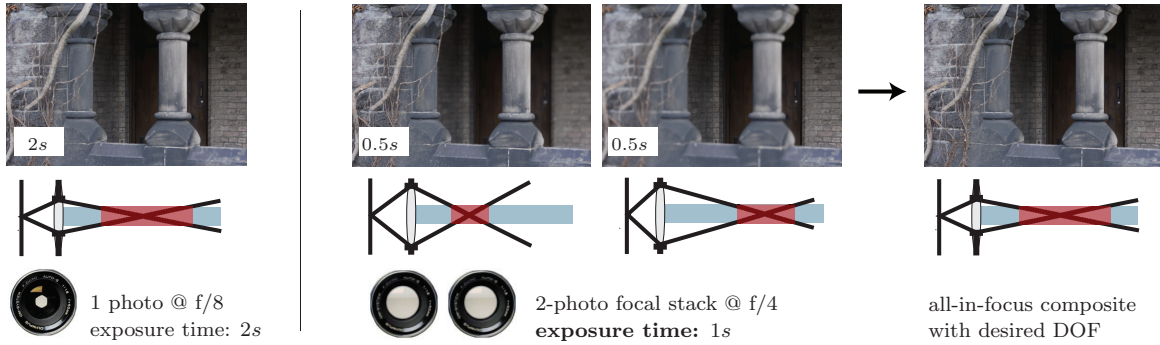


Figure 1: Speeding up photography with a desired depth of field. *Left:* Conventional one-shot photography. The desired DOF is shown in red. *Right:* Focal stack photography. Two wide-aperture photos span the same DOF as a one-shot narrow-aperture photo. Each wide-aperture photo requires 1/4 the time to reach the exposure level of the narrow-aperture photo, resulting in a 2× net speedup for the exposure time.

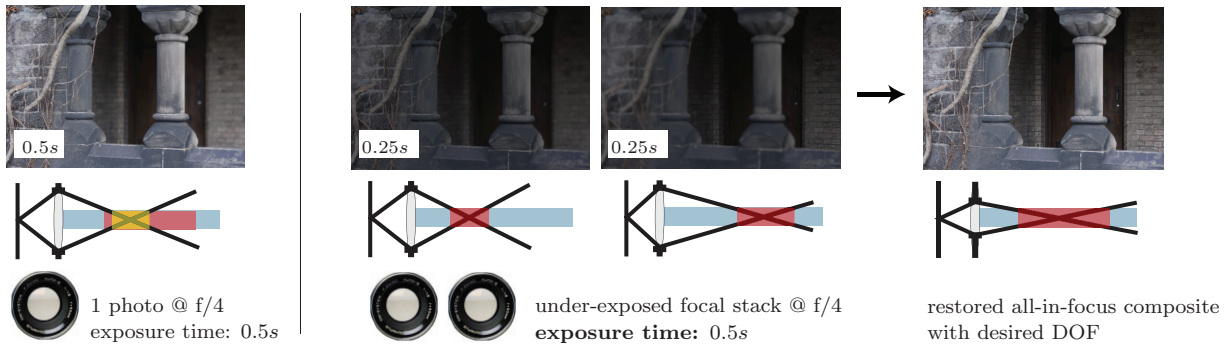


Figure 2: Increasing signal-to-noise ratio for restricted exposure time budgets. *Left:* Conventional one-shot photography. If exposure time is restricted to 0.5s, a well-exposed photo can span only a portion, shown in yellow, of the desired DOF. Subjects in the rest of the DOF will therefore be out of focus and, hence, will have reduced SNR. *Right:* Focal stack photography. Two wide-aperture, under-exposed photos span the desired DOF completely. Because under-exposure does not affect SNR as much as defocus blur in typical photography scenarios, the resulting composite will have higher SNR than a single well-exposed, wide-aperture photo.

2 Conventional Photography

Two of the most important choices when taking a photo are the photo’s exposure level and its depth of field. In conventional photography these choices are controlled indirectly, by choosing appropriate settings for the camera’s focus, aperture and exposure time. We review the basic relations between these quantities below and discuss how they constrain the photographic degrees of freedom.

Exposure level. The exposure level of a pixel is the total radiant energy integrated by the sensor element while the shutter is open (*i.e.*, number of photons). The exposure level can influence significantly the quality of a captured photo because when there is no saturation or thermal noise, a pixel’s signal-to-noise ratio (SNR) increases¹ at higher exposure levels [13]. For this reason, most modern cameras incorporate the notion of an “ideal” exposure level, *i.e.*, a level that provides good balance between SNR and likelihood of pixel saturation over the image. A typical choice is to capture photos with an average pixel intensity that is 13%

¹Thermal effects, such as dark-current noise, become significant only for exposure times longer than a few seconds [13].

of the maximum attainable value [16] (*i.e.*, 0.13×255 for 8-bit sensors).

Conventional cameras provide only two ways to control exposure level—the diameter of their aperture and the exposure time. Assuming that all light that passes through the aperture reaches the sensor plane, the exposure level L is equal to

$$L = \tau D^2, \quad (1)$$

where τ is exposure time, D is the effective aperture diameter, and the units of L are chosen appropriately.

Depth of field (DOF). We assume that focus and defocus obey the standard thin lens model [22, 24] (Figure 3). This model relates three positive quantities: the focus setting v , defined as the distance from the sensor plane to the lens; the distance d from the lens to the in-focus scene plane; and the focal length f , representing the “focusing power” of the lens (Eq. (A) in Table 1).

Apart from the idealized pinhole, all apertures induce spatially-varying amounts of defocus for points in the scene. If the lens focus setting is v , all points at distance d from the lens will be in focus. A scene point at distance $d' \neq d$, however, will be defocused: its image will be a circle on the sensor plane whose diameter

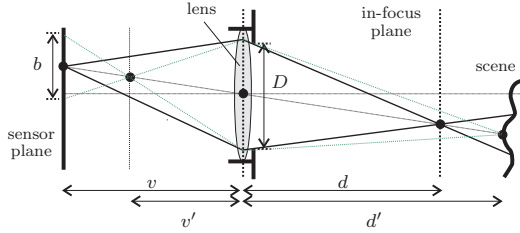


Figure 3: Blur geometry for a thin lens.

b is called the *blur diameter*. For any given distance d , the thin-lens model tells us exactly what focus setting we should use to bring the plane at distance d into focus, and what the blur diameter will be for points away from this plane (Eqs. (B) and (C), respectively).

For a given aperture and focus setting, the *depth of field* is the interval of distances in the scene whose blur diameter is below a maximum acceptable size. Since every distance in the scene corresponds to a unique focus setting, every DOF can also be expressed as an interval $[\alpha, \beta]$ in the space of focus settings. This alternate DOF representation gives us especially simple relations for the aperture and focus setting that produce a given DOF (Eqs. (D)-(E)) and, conversely, for the DOF produced by a given aperture and focus setting (Eq. (F)).

Note that a key property of the depth of field is that it shrinks when the aperture diameter increases: from Eq. (C) it follows that for a given out-of-focus distance, larger apertures always produce larger blur diameters.

Capturing an ideally-exposed photo with a given DOF. Now suppose that we want to capture a single photo with a specific exposure level L^* and a specific depth of field $[\alpha, \beta]$. How can we capture this photo? The basic DOF geometry, along with Eq. (1), tell us that we have little choice over camera settings: there is only one aperture diameter that can span the given depth of field exactly (Eq. (D)), and only one exposure time that can achieve a given exposure level with that diameter. This exposure time is given by

$$\tau^{one} = L^* \cdot \left(\frac{\beta - \alpha}{c(\beta + \alpha)} \right)^2. \quad (2)$$

Informally, Eq. (2) tells us that the larger the desired DOF, the longer it will take to capture an ideally-exposed photo (Figure 1, left). Unfortunately, this has an important practical side-effect: large exposure times can lead to motion blur when we photograph moving scenes or when the camera is not stabilized [27]. This limits the range of scenes that can be photographed with ideal exposure levels, with the range depending on scene radiance; the physical limits of the camera (*i.e.*, possible apertures and shutter speeds); and subjective factors such as the acceptable levels of motion blur and defocus blur.

Capturing a photo with a given DOF and restricted exposure time. In rapidly changing environments, it is often not possible to expose photos long enough to achieve an ideal exposure. This means that we must compromise something—we must either

(A) Thin lens law	(B) Focus for distance d	(C) Blur diameter for distance d'
$\frac{1}{v} + \frac{1}{d} = \frac{1}{f}$	$v = \frac{fd}{d-f}$	$b = D \frac{f d' - d }{d'(d-f)}$
(D) Aper. diam. for DOF $[\alpha, \beta]$	(E) Focus for DOF $[\alpha, \beta]$	(F) DOF for aper. diam. D , focus v
$D = c \frac{\beta + \alpha}{\beta - \alpha}$	$v = \frac{2\alpha\beta}{\alpha + \beta}$	$\alpha, \beta = \frac{Dv}{D \pm c}$

Table 1: Eqs. (A)–(F): Basic equations governing focus and DOFs for the thin-lens model. The maximum blur diameter within the DOF is assumed to be a user-specified tolerance value c .

reduce the exposure time (thereby inducing under-exposure) or reduce the DOF of the captured photo (thereby inducing defocus blur), or both. In other words, a restricted time budget imposes limits on a photo’s signal-to-noise ratio for subjects in the desired DOF (Figure 2, left).

3 Focal Stack Photography

In the following we briefly outline two ways that focal stack photography can enhance the optical performance of a conventional camera.

3.1 Achieving reduced exposure times

How quickly can we capture an ideally-exposed focal stack with a given DOF? Unlike conventional one-shot photography which is bound by Eq. (2), focal stacks do not need to cover the entire DOF in one photo. This added flexibility can lead to significant reductions in exposure time.

The efficiency of focal stack photography comes from the different rates at which exposure time and DOF change: if we increase the aperture diameter and adjust exposure time to maintain a constant exposure level, the lens DOF shrinks (at a rate of about $1/D$), but the time needed to get a single ideally-exposed photo shrinks much faster (at a rate of $1/D^2$). This opens the possibility of “breaking” time barrier of Eq. (2) by capturing a sequence of photos that span the DOF in less total time than τ^{one} .

Figure 1 shows a simple illustration of this idea: by splitting the desired DOF into two parts and capturing an ideally-exposed photo for each part with a matching DOF, we halve the necessary exposure time.

As a general rule, partitioning a given DOF into smaller segments and capturing one photo per segment almost always confers an exposure time advantage. The only exceptions to this rule occur when (1) the lens does not have an aperture large enough to allow further subdivision, (2) the camera has a significant per-shot overhead (*e.g.*, due to electronic or mechanical delays), or (3) discretization effects become significant, *i.e.*, the camera’s discrete aperture settings require capturing a focal stack with a DOF significantly larger than the desired one.

A complete analysis of how to achieve the shortest-possible exposure time for a given DOF and camera can be found in [12]. In particular, we show that we can compute the optimal focal stack by solving the

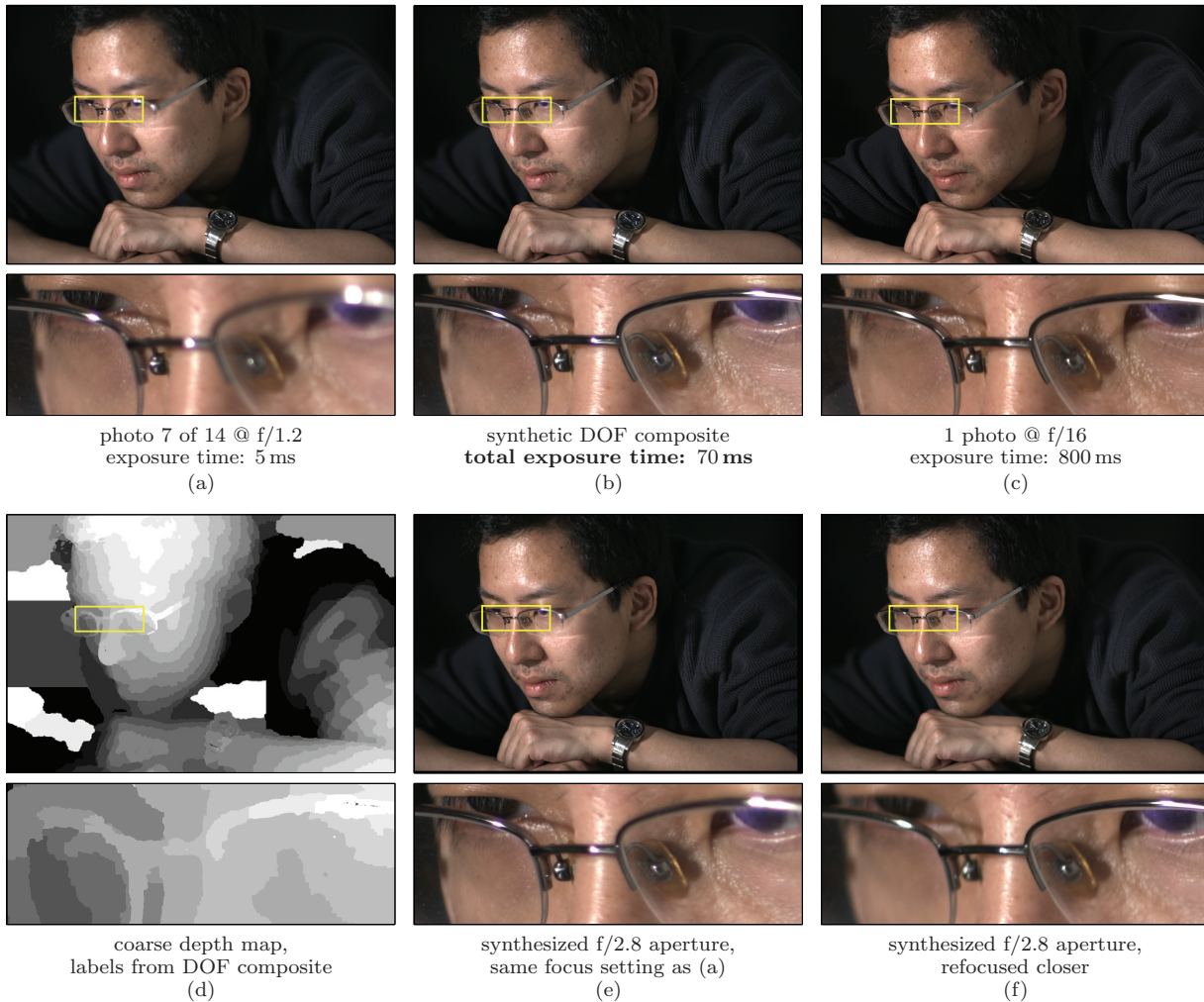


Figure 4: Rapid photography with a given DOF at the ideal exposure level. We captured real photos using a Canon EOS-1Ds Mark II camera with an EF85mm 1.2L lens. The desired DOF was $[100\text{cm}, 124\text{cm}]$, corresponding to a single f/16 photo. (a) A sample wide-aperture photo from the focal stack. (b) DOF composite synthesized using [6]. (c) A single narrow-aperture photo spanning the desired DOF; it requires a much longer exposure time. (d) Coarse depth map, computed from the labeling used to compute (b). Even though depth artifacts may occur in low-texture regions, this does not affect the quality of the all-in-focus image. (e) Synthetically changing aperture size, focused at the same setting as (a). (f) Synthetically changing focus, for the same synthetic aperture as (e).

following integer linear programming problem:

$$\text{minimize } \sum_{i=1}^m n_i \left[\frac{L^*}{D_i^2} + \tau^{\text{over}} \right] \quad (3)$$

$$\text{subject to } \sum_{i=1}^m n_i \log \frac{D_i - c}{D_i + c} < \log \frac{\alpha}{\beta} \quad (4)$$

$$n_i \geq 0 \text{ and integer,} \quad (5)$$

where D_i is the i -th aperture used in the optimal focal stack; n_i is the number of photos in the optimal stack taken with aperture D_i ; τ^{over} is the camera’s per-shot overhead; $[\alpha, \beta]$ is the desired DOF; c is the maximum acceptable blur diameter within the DOF; and L^* is the ideal exposure level.

Figure 5 shows the optimal subdivisions computed by solving the integer program for a specific camera and DOF, and for a range of camera overheads. The figure shows that focal stack photography results in significant speedups even in the presence of non-negligible per-shot camera overheads. In practice, we can pre-compute the most efficient focal stack for each photography scenario by solving the linear program for a

whole range of DOFs, and storing the results onboard the camera.

Results from a portrait photography experiment with a high-end digital SLR are shown in Figure 4. We applied Eqs. (3)-(5) to determine the optimal focal stack (in this case, it contained 14 f/1.2 photos) under an assumption of zero camera overhead. To merge the captured focal stack into a single photo, we used an existing depth-from-focus and compositing technique [6] from the computer graphics literature. In addition to an all-in-focus photo and a depth map, the captured focal stack allows us to “reshape” the camera’s DOF synthetically [12], producing photos with novel camera settings (Figure 4(e)-(f)).

3.2 Achieving high SNR with restricted exposure times

Given a fixed time budget, what is the best focal stack for capturing a given DOF? A constrained time budget prevents us from simultaneously spanning the

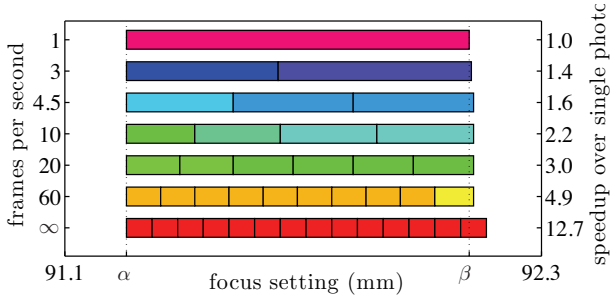


Figure 5: Optimal focal stacks computed for the Canon EF85mm f/1.2L lens, the DOF interval $[100\text{cm}, 124\text{cm}]$, and a one-shot exposure time $\tau^{\text{one}} = 1.5\text{s}$. Each row illustrates the optimal focal stack for a specific per-shot camera overhead, expressed in frames per second. To visualize each stack, we show how the DOF is partitioned by individual photos in the stack. Colors encode the aperture size of each photo (this lens has 23 distinct aperture settings). Note that as overhead increases, the optimal stack has fewer photos with larger DOFs (*i.e.*, smaller apertures).

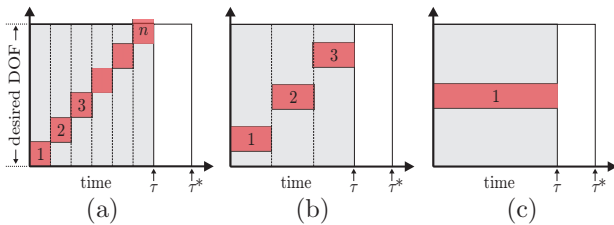


Figure 6: Focal stack photography with restricted exposure times. Red rectangles represent individual photos in a stack. Given a time budget of at least τ^* we can span the full DOF at the ideal exposure level (Sec. 3.1); restricting the budget to $\tau < \tau^*$ causes a reduction in SNR. (a) A simple policy to meet the restricted time budget is to reduce the exposure time of each photo proportionally, leading to increased noise. (b)-(c) Focal stacks with fewer photos yield brighter images, but at the expense of defocus since the DOF is spanned incompletely.

DOF and achieving the ideal exposure level. Nevertheless, it is possible to find the focal stack that yields the highest signal-to-noise ratio, *i.e.*, one that optimally balances defocus and noise (Figure 6).

In broad terms, dense focal stacks achieve the highest SNR because defocus from a conventional lens causes a more severe degradation than under-exposure. We confirmed this intuition through a detailed analysis of lens defocus, sensor noise, and the resulting restoration problem [10]: this analysis involved both a theoretical component (frequency-based restoration) and a series of simulations that covered a broad range of photography conditions and used data from high-end cameras, photographic lenses, and imaging sensors. Our analysis showed that capturing fewer photos is only beneficial for severely limited time budgets or cameras with high per-shot overhead.

From a theoretical standpoint, establishing the highest-SNR focal stack under a restricted time budget is more complex than the analysis in Section 3.1 for two reasons. First, since the optimal focal stack

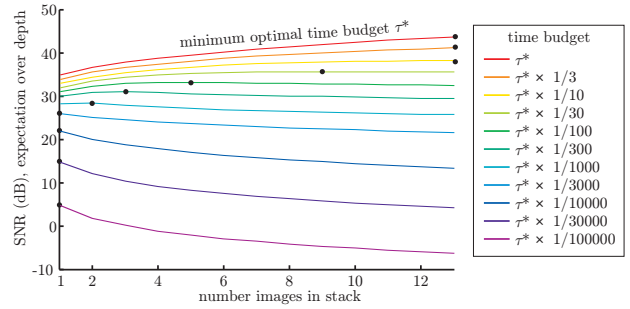


Figure 7: SNR for different focal stack sizes, for the same DOF and camera as in Figure 5, and all photos taken with an f/1.23 aperture. Each curve represents a different time budget, relative to the minimum time τ^* required to capture a well-exposed focal stack spanning the DOF. The black dot on each curve indicates the optimal stack size for that time budget.

may not span the DOF, we need a more detailed image formation model to quantify degradations due to defocus. In our approach, we represent defocus using the modulation transfer function (MTF) of the lens, which provides a frequency-based description of attenuation due to defocus blur. We represent the MTF with a classic diffraction model that takes diffraction into account [14], and use the blur diameter of Eq. (C) as its main parameter.

Second, because the focal stacks we consider have varying exposure levels, we also must consider the noise properties of the sensor [13]. To model the essential characteristics of sensor noise, we rely on a two-component affine model [19, 8] that contains a multiplicative term, approximating Poisson-distributed shot noise, and a constant term, accounting for read noise and quantization:

$$\varepsilon(x, y)^2 \sim \mathcal{N}\left(0, \sigma_s^2 \left(\frac{L}{L^*}\right) \tilde{I}(x, y) + \sigma_c^2\right), \quad (6)$$

where \tilde{I} is the ideally-exposed noise-free photo. Note that because the multiplicative component scales with the relative exposure level $\frac{L}{L^*}$, it does not affect SNR when the photo is underexposed.

Quantitative results from our simulation analysis are shown in Figure 7 for a particular camera and DOF. The figure shows how the SNR of the restored, all-in-focus photo varies as a function of focal stack size and available time budget; also indicated is the optimal stack size in each case. These results confirm that under-exposed focal stacks spanning a DOF completely (Figure 6a) result in the highest SNR, unless the time budget is severely restricted. Indeed, only when the exposure time is 30 times less than required to capture an ideally-exposed focal stack does noise become high enough to tilt the balance toward incompletely spanning the DOF (Figure 6(b)). In the limit, for time budgets reduced by a factor of 3000 or more, photos are so under-exposed that little more than the DC component can be captured by the focal stack. In such cases, one-shot photography is the only viable option.

The results in Figure 7 can be thought of as generalizing the geometric treatment of Section 3.1, and essentially arrive at the same conclusion: whether the goal is to reduce exposure time or to increase SNR, dense wide-aperture focal stacks confer a significant advan-

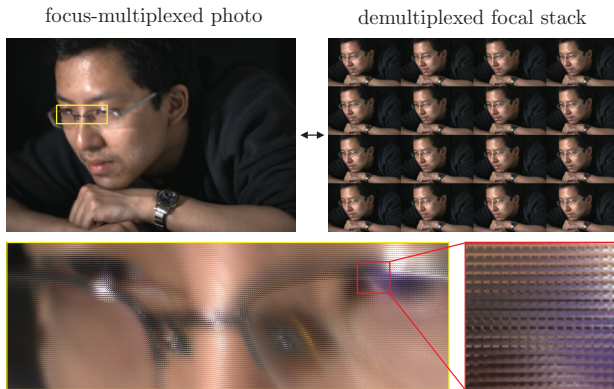


Figure 8: A focus-multiplexed photo representing the focal stack in the experiment of Figure 4. Each 4×4 pixel block corresponds to sixteen different focus settings stored in scanline order. By rearranging the pixels in the 16 Mpixel focus-multiplexed image, we obtain a focal stack with sixteen 1 Mpixel photos.

tage over one-shot photography. The specific parameters of these stacks (number of photos, apertures and focal settings) can be computed in advance for different DOF sizes, scene brightness levels, and time budgets, and stored onboard the camera.

4 A Focal Stack Camera?

How could one realize focal stack photography in practice? A basic approach would be to modify the firmware of existing digital cameras to capture images in burst mode and refocus programmatically [2]. This would enable capturing optimal focal stacks pre-computed for each photography setting.

To be most useful, focal stacks should be captured with low per-photo camera overhead. The current trend in digital cameras suggests that overheads will continue to decrease rapidly. Already, the line between digital photography and video has become blurred, with recent cameras capturing 2 Mpixel high-definition video with single-lens reflex lenses [3], and others achieving 60 fps at full 6 Mpixel resolution [4]. Moreover, lenses can focus and re-focus very quickly: typical lenses with ultrasonic motors need less than 3ms to refocus between images in a focal stack.

For very tight exposure time budgets (*e.g.*, less than 10ms), another important bottleneck is the data transfer between sensor and memory. While this rate continues to improve, transfer bottlenecks can be alleviated by trading off stack density and image resolution. For instance, instead of recording a full-resolution focal stack, the camera could store a “focus-multiplexed” image, where many lower-resolution photos from the stack are packed into a single full-resolution pixel array (Figure 8).

In principle, focus-multiplexed images can be captured with a conventional camera using the strategies outlined in this paper. An open question, however, is whether it is possible to capture such images in one shot by altering the camera’s optics, in a spirit similar to the plenoptic camera [5, 18]. We believe that such a camera would be broadly useful in photography, and its design is a subject of our ongoing efforts.

Acknowledgements. This work was supported in part by the Natural Sciences and Engineering Research Council of Canada under the RGPIN, RTI and PDF programs, and by an Ontario Premier’s Research Excellence Award.

References

- [1] Canon PowerShot G10 Review, <http://www.dpreview.com/reviews/canong10/>.
- [2] CHDK, <http://chdk.wikia.com/>.
- [3] Canon EOS 5D Mark II Review, <http://www.dpreview.com/reviews/canoneos5dmarkii/>.
- [4] Casio Pro EX-F1, http://exilim.casio.com/products_exf1.shtml.
- [5] E. H. Adelson and J. Y. A. Wang. Single lens stereo with a plenoptic camera. *TPAMI*, 14(2):99–106, 1992.
- [6] A. Agarwala, M. Dontcheva, M. Agrawala, S. Drucker, A. Colburn, B. Curless, D. Salesin, and M. Cohen. Interactive digital photomontage. In *SIGGRAPH*, pp. 294–302, 2004.
- [7] N. Asada, H. Fujiwara, and T. Matsuyama. Edge and depth from focus. *IJCV*, 26(2):153–163, 1998.
- [8] R. N. Clark. Digital camera sensor performance summary, <http://clarkvision.com/>, 2009.
- [9] P. Favaro and S. Soatto. A geometric approach to shape from defocus. *TPAMI*, 27(3):406–417, 2005.
- [10] S. W. Hasinoff. *Variable-Aperture Photography*. PhD thesis, University of Toronto, Dept. of Computer Science, 2008.
- [11] S. W. Hasinoff and K. N. Kutulakos. Confocal stereo. In *ECCV*, vol. 1, pp. 620–634, 2006.
- [12] S. W. Hasinoff and K. N. Kutulakos. Light-efficient photography. In *ECCV*, vol. 4, pp. 45–59, 2008.
- [13] G. E. Healey and R. Kondepudy. Radiometric CCD camera calibration and noise estimation. *TPAMI*, 16(3):267–276, 1994.
- [14] H. H. Hopkins. The frequency response of a defocused optical system. *Proc. of the Royal Society of London, Series A*, 231(1184):91–103, 1955.
- [15] B. K. P. Horn. Focusing. Technical Report AIM-160, Massachusetts Institute of Technology, 1968.
- [16] ISO 2721:1982. Photography—Cameras—Automatic controls of exposure, 1982.
- [17] A. Levin, W. T. Freeman, and F. Durand. Understanding camera trade-offs through a Bayesian analysis of light field projections. In *ECCV*, vol. 4, pp. 88–101, 2008.
- [18] M. Levoy, R. Ng, A. Adams, M. Footer, and M. Horowitz. Light field microscopy. In *SIGGRAPH*, pp. 924–934, 2006.
- [19] C. Liu, R. Szeliski, S. Kang, C. Zitnick, and W. Freeman. Automatic estimation and removal of noise from a single image. *TPAMI*, 30(2):299–314, 2008.
- [20] S. K. Nayar and Y. Nakagawa. Shape from focus: an effective approach for rough surfaces. In *ICRA*, vol. 2, pp. 218–225, 1990.
- [21] J. Ogden, E. Adelson, J. R. Bergen, and P. Burt. Pyramid-based computer graphics. *RCA Engineer*, 30(5):4–15, 1985.
- [22] A. P. Pentland. A new sense for depth of field. *TPAMI*, 9(4):523–531, 1987.
- [23] R. J. Pieper and A. Korpel. Image processing for extended depth of field. *App. Optics*, 22(10):1449–1453, 1983.
- [24] W. J. Smith. *Modern Optical Engineering*. McGraw-Hill, New York, 3rd edition, 2000.
- [25] N. Streibl. Three-dimensional imaging by a microscope. *JOSA A*, 2(2):121–127, 1985.
- [26] N. Xu, K. Tan, H. Arora, and N. Ahuja. Generating omnifocus images using graph cuts and a new focus measure. In *ICPR*, vol. 4, pp. 697–700, 2004.
- [27] L. Yuan, J. Sun, L. Quan, and H.-Y. Shum. Image deblurring with blurred/noisy image pairs. In *SIGGRAPH*, 2007.