Check for updates

Optics Letters

Structure-aware parametric representations for time-resolved light transport

DIEGO ROYO,^{1,†} DZESHENG HUANG,^{2,†} YUN LIANG,^{2,*} BOYAN SONG,² Adolfo Muñoz,¹ DIEGO GUTIERREZ,¹ AND JULIO MARCO¹

¹Universidad de Zaragoza—I3A, Zaragoza, Spain

²South China Agricultural University, Guangzhou, China

*Corresponding author: yliang@scau.edu.cn

[†]These authors contributed equally to this paper.

Received 31 May 2022; revised 3 September 2022; accepted 7 September 2022; posted 8 September 2022; published 30 September 2022

Time-resolved illumination provides rich spatiotemporal information for applications such as accurate depth sensing or hidden geometry reconstruction, becoming a useful asset for prototyping and as input for data-driven approaches. However, time-resolved illumination measurements are high-dimensional and have a low signal-to-noise ratio, hampering their applicability in real scenarios. We propose a novel method to compactly represent time-resolved illumination using mixtures of exponentially modified Gaussians that are robust to noise and preserve structural information. Our method yields representations two orders of magnitude smaller than discretized data, providing consistent results in such applications as hidden-scene reconstruction and depth estimation, and quantitative improvements over previous approaches. © 2022 Optica Publishing Group

https://doi.org/10.1364/OL.465316

Introduction. Transient imaging methods analyze timeresolved light transport at very high temporal resolutions, with applications such as reconstruction of hidden geometry [1,3,4], object detection through scattering media [5], or material classification [6], with promising advances over recent years [7]. Current capture methods combine ultra-fast lasers with sensors such as single-photon avalanche diode (SPAD) arrays [8,9], providing dense spatial scanning and picosecond temporal resolution that yield rich spatiotemporal information of the captured scene.

These capture setups, however, introduce several limitations. First, analyzing the spatiotemporal structure of indirect illumination is fundamental in many transient imaging applications, but multiply scattered light can become too attenuated when reaching the sensor. Consequently, imaging applications may be subject to measurements with a low signal-to-noise ratio (SNR), which can decrease their performance. Second, dense temporal and spatial resolutions of the measurement space result in high memory and bandwidth requirements (with datasets of tens or hundreds of gigabytes) [10], which can become a bottleneck in light transport analysis and application design.

In this work, we provide a method for lightweight representations of time-resolved illumination with a threefold benefit: the representation space is up to two orders of magnitude smaller than the source data and two times smaller than previous compression approaches, it is robust to noise, and it preserves structural information fundamental in transient imaging applications. We rely on mixtures of exponentially modified Gaussian (EMG) distributions and design an optimization procedure that accounts for spatial gradients to preserve structural information.

Several previous works propose alternative representations of time-resolved light pulses with different goals. Note that none of these methods is structure-aware, as they only use temporal information from single points in the scene. Peters et al. [11] use Pisarenko and maximum entropy spectral estimates to reconstruct transient pulses and improve data quality by removing multipath interference in range imaging. For these same purposes, Kadambi et al. [12] recover per-pixel sparse time profiles expressed as a sequence of impulses. Other works are based on linear inverse problems [13], solved by numerical optimization, and frequency-domain reconstructions [14], which introduce Fourier analysis to reduce systematic errors. Closer to our representation space, Wu et al. [15] analyze direct and indirect illumination components by representing light pulses as a combination of one Gaussian and one EMG. Heide et al. [5] recover time-resolved illumination in turbid media from correlation-based sensors by sparsely encoding light with EMG distributions. Directly related to our goal, Liang et al. [2] recently introduced feature-based compact representations of time-resolved illumination using deep encoder-decoder neural architectures. However, their approach is biased toward line-ofsight training data and fails to preserve structural information that is fundamental in modern transient imaging applications. As shown in Fig. 1, our method yields representations of previously captured transient illumination histograms with higher quality that preserve structural information, providing significantly better performance on hidden geometry reconstructions than feature-based methods, even with a higher compression ratio.

Inspired by previous works [5,15], we propose to use EMG distributions to compactly represent time-resolved illumination. Consider a transient camera, which adds a third temporal dimension *T* for an image *I* with size $W \times H$, i.e., $I[i,j,t] \in \mathbb{R}^{W \times H \times T}$. A time-resolved pixel $I_{\mathbf{p}}[t] \in \mathbb{R}^{1 \times 1 \times T}$ at $\mathbf{p} = \{i, j\}$ represents the



Fig. 1. Application of our method to non-line-of-sight (NLOS) data for a real scene. Top: original captured signal histogram. Bottom: adding noise to the histogram to simulate an exposure time 200 times shorter. The four letters in the scene are then reconstructed following the work of Liu *et al.* [1]. Right: results for the original captured and compressed signals obtained using the autoencoder network proposed by Liang *et al.* [2] and our method.

accumulation of light paths with a timestamp t traveling from the light sources to the sensor pixel **p** after being scattered through the scene elements. Time-resolved illumination typically has the temporal shape of aggregated radiance pulses with exponential decay. This behavior stems from multi-bounce convolutions of the source illumination pulse with the scene geometry. EMG distributions model a Gaussian with a parameterized exponential decay, and therefore arise as a convenient function to model the physical behavior of different illumination bounces in time [5,15]. We therefore propose to represent the response of a pulsed source measured at an ultra-fast sensor pixel by aggregating EMG distributions in a mixture model. We analyze the benefits of our method on real data captured on non-line-ofsight (NLOS) configurations [1] and simulated datasets [16,17], which have proved to faithfully represent data captured with real hardware. We demonstrate the consistency of our representation in such applications as NLOS reconstruction [1] and line-ofsight (LOS) depth estimation based on amplitude-modulated continuous-wave (AMCW) time-of-flight (ToF) sensors [16].

EMGs. An EMG distribution is defined as

$$\operatorname{EMG}(t;h,\mu,\sigma,\tau) = \frac{h\sigma}{\tau} \sqrt{\frac{\pi}{2}} \exp\left(\frac{1}{2} \left(\frac{\sigma}{\tau}\right)^2 - \frac{t-\mu}{\tau}\right) \\ * \operatorname{erfc}\left(\frac{1}{\sqrt{2}} \left(\frac{\sigma}{\tau} - \frac{t-\mu}{\sigma}\right)\right),$$
(1)

where *h* controls the pulse amplitude, τ controls the exponential decay rate, and the mean μ and standard deviation σ control the peak position and width of the Gaussian distribution. The complementary error function, erfc(·), is defined as

$$\operatorname{erfc}(x) = \frac{2}{\sqrt{\pi}} \int_{x}^{\infty} e^{-t^{2}} dt.$$
 (2)

In our method, we model the aggregation of multi-bounce light paths on a time-resolved single pixel $I_p[t]$ as an EMG mixture model (EMGMM) with K EMGs, denoted $I'_p[t]$:

1

$$I_{\mathbf{p}}[t] \approx I'_{\mathbf{p}}[t] = \text{EMGMM}(t; \boldsymbol{h}_{p}, \boldsymbol{\mu}_{p}, \boldsymbol{\sigma}_{p}, \boldsymbol{\tau}_{p}, K)$$
$$= \sum_{k=1}^{K} \text{EMG}(t; \boldsymbol{h}_{k}, \boldsymbol{\mu}_{k}, \boldsymbol{\sigma}_{k}, \boldsymbol{\tau}_{k}),$$
(3)

where each a_p term represents a vector of K EMG parameters a_k with k = 1, ..., K for the estimation of pixel **p**.

We formulate the estimation of a transient pixel $I_p[t] \approx I'_p[t]$ as an optimization problem, which attempts to compute the best



Fig. 2. Left: EMGMM with *K* EMG modules. Right: EMG substructure with four parameters and one input connected to the probability mass function (PMF) node corresponding to Eq. (1). We use exponential functions to ensure $h, \sigma, \tau > 0$. As $t \in [0, 1]$ is normalized, a sigmoid function is applied to μ .

EMGMM parameterization { h_p , μ_p , σ_p , τ_p } in order to minimize the pixel loss function \mathcal{L}_p :

$$\underset{h_{p},\mu_{p},\sigma_{p},\tau_{p}}{\arg\min} \mathcal{L}_{p}(I_{p},I'_{p}), \qquad (4)$$

where \mathcal{L}_p is lower when both time-resolved illumination distributions are more similar. Since we are working with EMG distributions, we define \mathcal{L}_p based on the Kullback–Leibler divergence (KLD) [18], a statistical distance metric used to measure the difference of probability distributions. For a pixel I_p and its reconstruction I'_p , the KLD is defined as

$$D_{\mathrm{KL}}(I'_{\mathbf{p}} \parallel I_{\mathbf{p}}) = \sum_{t \in T} I'_{\mathbf{p}}[t] \log \left(\frac{I'_{\mathbf{p}}[t]}{I_{\mathbf{p}}[t]}\right).$$
(5)

Note that KLD is asymmetric, and also produces low errors where $I'_{p} \approx 0$ and high errors where $I_{p} \approx 0$. To avoid these extreme cases, we use its asymmetric property to finally construct our pixel loss function \mathcal{L}_{p} :

$$\mathcal{L}_{p}(I_{\mathbf{p}}, I'_{\mathbf{p}}) = D_{\mathrm{KL}}(I_{\mathbf{p}} \parallel I'_{\mathbf{p}}) + D_{\mathrm{KL}}(I'_{\mathbf{p}} \parallel I_{\mathbf{p}})$$

= $\sum_{t \in T} (I_{\mathbf{p}}[t] - I'_{\mathbf{p}}[t])(\log (I_{\mathbf{p}}[t]) - \log (I'_{\mathbf{p}}[t])).$ (6)

Given that \mathcal{L}_p [Eq. (6)] and EMGs [Eqs. (1) and (2)] are differentiable, we use stochastic gradient descent to optimize for the best parameters in an EMG-based differentiable pipeline, as shown in Fig. 2. For an input time *t*, the pipeline should output the pixel intensity at that time $I'_p[t]$.

Other loss functions, such as the mean square error (MSE) [2], tend to be biased toward large-valued regions in the temporal domain, which may lead the optimization to fall into

Table 1. KLD Loss \mathcal{L}_I	$=\sum_{\mathbf{p}} \mathcal{L}_{\mathbf{p}}(\mathbf{I}_{\mathbf{p}},\mathbf{I}_{\mathbf{p}})$	in a	Region	of the
Staircase NLOS Scene	[17] ^a			

	KLD Loss			MSE Loss
	\mathcal{L}_p	$\mathcal{L}_{p} + \mathcal{L}_{g}(S)$	$\mathcal{L}_{p} + \mathcal{L}_{g}(R)$	$\mathcal{L}_p + \mathcal{L}_g (R)$
\mathcal{L}_{I}	0.012	0.010	0.008	0.010

^{*a*}Using K = 64 for pixel-independent (\mathcal{L}_p) , and *Sliding* (*S*) and *Random* (*R*) structure-aware fitting $(\mathcal{L}_p + \mathcal{L}_g)$. The *Random* fitting is also tested with an MSE loss.

local minima. Our KLD-based optimization reduces the error more uniformly across the entire temporal domain, as shown in Table 1, compared with the MSE-based approach.

For practical reasons, we pre-process the optimization input as follows: we first clip the temporal domain of each pixel $I_{\mathbf{p}}[t]$ to $t \in [t_{\text{start}}, T]$, where $I_{\mathbf{p}}[t_{\text{start}}]$ is the first non-zero value. We then normalize the clipped temporal domain to $t \in [0, 1]$ by subtracting t_{start} and dividing by the resulting length $T' = T - t_{\text{start}}$. The EMGMM [Eq. (3)] is evaluated in a reparameterized time interval $t \in [0, 1]$, where the reference pixel $I_{\mathbf{p}}[t]$ does not have leading zeros. We store t_{start} , T' along with the 4K EMGMM parameters $\mathbf{h}_{p}, \mu_{p}, \sigma_{p}, \tau_{p}$, and revert the clipping and normalization to compute the loss [Eq. (6)], resulting in a compression factor of T/(4K + 2) and an optimization runtime within O(KT).

Reconstructing a transient image. Eq. (4) defines the optimization scheme to represent a single pixel $I_{\mathbf{p}}[t]$, with $\mathbf{p} = \{i, j\}$ using an EMGMM. In a full transient image $\equiv I[i, j, t]$, neighboring pixels usually present a significant spatiotemporal structure. We propose an optimization methodology to use spatiotemporal information to improve the pixel representation in a transient image with a twofold benefit. First, reduction of noise in low-SNR areas by leveraging noisy information from nearby pixels. Second, preserving the spatial structure of the signal by accounting for the spatial gradient. In particular, we use the $N \times N$ window around each pixel $\mathbf{p} / - \mathbf{p} |_{\infty} < N/2$. We introduce a spatial gradient loss term \mathcal{L}_g that fosters consistency between neighboring pixels, reducing the influence of per-pixel noise:

$$\mathcal{L}_{g}(I_{W}, I'_{W}) = \sum_{t \in T} \left(\|G_{up}(I_{W}, t) - G_{up}(I'_{W}, t)\|_{2} + \|G_{down}(I_{W}, t) - G_{down}(I'_{W}, t)\|_{2} + \|G_{left}(I_{W}, t) - G_{left}(I'_{W}, t)\|_{2} + \|G_{right}(I_{W}, t) - G_{right}(I'_{W}, t)\|_{2} \right).$$
(7)

Each gradient G is defined for a neighborhood W where each element is computed using its immediate neighbors as

$$\begin{cases} G_{up}(I_{W},t) = I[i,j-1,t] - I[i,j,t] \\ G_{down}(I_{W},t) = I[i,j+1,t] - I[i,j,t] \\ G_{left}(I_{W},t) = I[i-1,j,t] - I[i,j,t] \\ G_{right}(I_{W},t) = I[i+1,j,t] - I[i,j,t] \end{cases}, \mathbf{p} = \{i,j\} \in \mathcal{W}, (8)$$

adding zero-padding on the image to satisfy the calculation for edge pixels. The final optimization problem is

$$\arg\min_{\boldsymbol{h}_{W},\boldsymbol{\mu}_{W},\boldsymbol{\sigma}_{W},\boldsymbol{\tau}_{W}} \mathcal{L}_{g}(\boldsymbol{I}_{W},\boldsymbol{I}_{W}') + \sum_{\mathbf{p}\in W} \mathcal{L}_{p}(\boldsymbol{I}_{\mathbf{p}},\boldsymbol{I}_{\mathbf{p}}'),$$
(9)

reconstructing $N \times N$ pixels in the image simultaneously. The resulting optimization accounts for a number of EMGMMs, as in Fig. 2, for each pixel in the neighborhood, calculating the

gradient afterwards. Typical values are $N \in \{3, 5, 7\}$. From our experiments, N = 5 provides the best trade-off between performance and time, with little difference from other values. Also, the data need to be pre-processed, as explained previously. The value of t_{start} is obtained as the minimum for all pixels $\mathbf{p} \in W$ based on the first pixel that receives a photon. The compression ratio for the whole neighborhood is $(N^2 \cdot T)/(N^2 \cdot 4K + 2)$, with an execution time of the optimization within $O(N^2KT)$.

Initialization. Since our optimization is based on a stochastic gradient descent, a good initialization is crucial to avoid convergence to bad local minima. Light transport in real-world scenes decays exponentially: short, high-energy paths arrive early and have higher temporal frequency, while lower-energy, multiply scattered paths with longer optical lengths are smooth and arrive later in time. We use these observations to estimate convenient initial values for the peak μ , width σ , and decay rate τ of each of the *K* EMG modeled pulses. Note that the amplitude *h* is just a scale factor, so we arbitrarily initialize it as h = 1.

For each parameter, we initialize values v_i in log space, which provides a 5–10% error improvement with respect to linearspace initialization for a similar number of epochs. We uniformly sample $\xi_i \in [\log(v_{\min}), \log(v_{\max})]$ and then convert to linear range as $v_i = \exp(\xi_i)$. First, we sample *K* values for μ and \sqrt{K} values for σ and τ , of which we generate all possible *K* combinations. The first bounces with smaller μ are narrower, so they are assigned the smaller values of σ and τ . This is done for each of the values of *K* of μ and (σ, τ) pairs. Under this initialization, illumination pulses centered at later timestamps μ will have a wider support σ and a slower decay τ .

To reconstruct an image $I[i, j, t] \in \mathbb{R}^{W \times H \times T}$, the order for which we optimize the pixels is important when using spatial consistency, as in Eq. (9). A first approach would be to slide through the $N \times N$ image pixels while optimizing each pixel neighborhood. This can produce square artifacts, so we use a *Random* sampling of the image and optimize a number of neighborhoods until every pixel has converged. Our quantitative analysis, given in Table 1, shows that, while the *Sliding* approach provides better results than independently fitting each pixel, the *Random* approach provides the best results overall.

Results. We evaluate the performance of signal compression for our *Single-pixel* [Eq. (4)] and *Structure-aware* [Eq. (9)] models, and compare it with that of the 3D convolutional autoencoder network recently proposed by Liang *et al.* [2], designed for the same purposes.

In our method, the number of EMG components K introduces a trade-off between quality and efficiency. Figure 3 shows a comparison with the *Single-pixel* model. K = 4 is enough to surpass deep-learning methods, increasing temporal similarity for higher values of K. Execution times range from 5 s per pixel with K = 4 EMGs, to 10 s with K = 64 EMGs (Intel Xeon Gold 6140 CPU, using 10 threads). Figure 4 uses K = 64 EMGs.



Fig. 3. Temporal slices and KLD loss metric $\mathcal{L}_{I} = \sum_{\mathbf{p}} \mathcal{L}_{p}(I_{\mathbf{p}}, I'_{\mathbf{p}})$ in the *Bathroom* scene [16], comparing each compressed image I' with its original I for different numbers K of EMGs, yielding much better performance than previous work [2].



Fig. 4. Signal representation of a 10×10 region of *Staircase* [17], as in Fig. 5. The last row shows a spatial slice $I_p[t]$ for the pixel **p** marked with a cross. The vertical dotted lines correspond to four temporal slices at different instants (A–D).



Fig. 5. Comparison for a LOS *Church* scene [16]. Our method using *Single-pixel* mode (green) outperforms previous works (orange) in both signal reconstruction (spatial and temporal slices) and AMCW ToF depth estimation (bottom right).

It showcases the importance of the 5×5 window with spatial gradient constraints, added in the Structure-aware model, compared with the Single-pixel model. Following this, Fig. 1 shows a real-world NLOS imaging application, where the capture setups produce noisier measurements, compared with simulated data. Our work provides a higher-quality representation of both the signal and hidden geometry reconstruction, while also improving on the robustness to noise. The compression ratio is also higher, $T/(4 \cdot K + 2) \approx 42$, reducing the resulting file size from 365 mebibytes to 8.7 mebibytes. Finally, Fig. 5 shows results for a LOS scene [16] with K = 16, resulting in approximately 62 times less parameters than the original signal with the Singlepixel model. This almost doubles the compression ratio of the autoencoder network of Liang et al. [2], while providing more accurate representations, as shown in the spatial (A-D) and temporal (bottom left) slices of I[i, j, t]. The bottom-right images of Fig. 5 show a comparison with a practical application of AMCW ToF depth estimation, where our output (right) yields more consistent depth results than the output of Liang et al. (center) [2].

In conclusion, we present a method to efficiently represent time-resolved light transport data based on EMGs, significantly reducing the number of coefficients required to represent timeresolved transport. Our optimization loss, based on first-order spatial differences, preserves the spatial structure and reduces noise, both fundamental aspects in transient imaging applications. We demonstrate the benefits of our method in both LOS and NLOS scenarios, overcoming previous approaches targeted to transient light transport compression. Relating the EMG components of our representation to higher-order illumination bounces may help to increase the quality of scene-understanding applications such as multipath interference correction in depth imaging or hidden-scene reconstruction. Additionally, analysis of the frequency components of our EMG-based representation could be exploited in wave-based NLOS imaging algorithms [1,3] to improve their computational efficiency.

Funding. H2020 European Research Council (682080); Agencia Estatal de Investigación (PID2019-105004GB-I00); Gobierno de Aragón (MP30_21); Gobierno de Aragón; Key Technologies Research and Development Program of Guangzhou (202206010091); Science and Technology Planning Project of Guangdong Province (2019A050510034); Guangzhou Key Laboratory of Intelligent Agriculture (201902010081).

Disclosures. The authors declare no conflicts of interest.

Data availability. Data underlying the results presented in this paper come from the *DeepToF* [16] and *Z-NLOS* [17] datasets, as well as from Code 1, Ref. [19].

REFERENCES

- X. Liu, I. Guillén, M. La Manna, J. H. Nam, S. A. Reza, T. H. Le, A. Jarabo, D. Gutierrez, and A. Velten, Nature 572, 620 (2019).
- Y. Liang, M. Chen, Z. Huang, D. Gutierrez, A. Mu noz, and J. Marco, Opt. Lett. 45, 1986 (2020).
- D. B. Lindell, G. Wetzstein, and M. O'Toole, ACM Trans. Graph. 38, 1 (2019).
- S. Xin, S. Nousias, K. N. Kutulakos, A. C. Sankaranarayanan, S. G. Narasimhan, and I. Gkioulekas, in *IEEE Computer Vision and Pattern Recognition (CVPR)*, (2019), pp. 6800–6809.
- F. Heide, L. Xiao, A. Kolb, M. B. Hullin, and W. Heidrich, Opt. Express 22, 26338 (2014).
- S. Su, F. Heide, R. Swanson, J. Klein, C. Callenberg, M. Hullin, and W. Heidrich, in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, (2016), pp. 3503–3511.
- A. Jarabo, B. Masia, J. Marco, and D. Gutierrez, Visual Informatics 1, 65 (2017).
- M. Renna, J. H. Nam, M. Buttafava, F. Villa, A. Velten, and A. Tosi, Instruments 4, 14 (2020).
- J. H. Nam, E. Brandt, S. Bauer, X. Liu, M. Renna, A. Tosi, E. Sifakis, and A. Velten, Nat. Commun. 12, 6526 (2021).
- J. Marco, A. Jarabo, J. H. Nam, X. Liu, M. Ángel Cosculluela, A. Velten, and D. Gutierrez, in 2021 IEEE/CVF International Conference on Computer Vision (ICCV) (2021).
- C. Peters, J. Klein, M. B. Hullin, and R. Klein, ACM Trans. Graph. 34, 1 (2015).
- A. Kadambi, R. Whyte, A. Bhandari, L. Streeter, C. Barsi, A. Dorrington, and R. Raskar, ACM Trans. Graph. 32, 1 (2013).
- F. Heide, L. Xiao, W. Heidrich, and M. B. Hullin, in IEEE Computer Vision and Pattern Recognition (2014), p. 3222.
- J. Lin, Y. Liu, M. B. Hullin, and Q. Dai, in IEEE Computer Vision and Pattern Recognition (2014), p. 3222.
- D. Wu, A. Velten, M. O'Toole, B. Masia, A. Agrawal, Q. Dai, and R. Raskar, Int. J. Comput. Vis. **107**, 366 (2014).
- J. Marco, Q. Hernandez, A. Mu noz, Y. Dong, A. Jarabo, M. Kim, X. Tong, and D. Gutierrez, ACM Trans. Graph. 36, 1 (2017).
- M. Galindo, J. Marco, M. O'Toole, G. Wetzstein, D. Gutierrez, and A. Jarabo, in SIGGRAPH '19 (ACM, 2019).
- 18. S. Kullback and R. A. Leibler, Ann. Math. Statist. 22, 79 (1951).
- D. Royo, "Structure-aware parametric representations for timeresolved light transport," figshare (2022), https://doi.org/10.6084/m9.figshare.21253359.