



**HAL**  
open science

# Deep learning depth-intensity reconstruction from compressive TCSPC SPAD-based imaging

Valentin Poisson, William Guicquero, Gilles Sicard

► **To cite this version:**

Valentin Poisson, William Guicquero, Gilles Sicard. Deep learning depth-intensity reconstruction from compressive TCSPC SPAD-based imaging. 2023. hal-04192408

**HAL Id: hal-04192408**

**<https://hal.science/hal-04192408v1>**

Preprint submitted on 31 Aug 2023

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Deep learning depth-intensity reconstruction from compressive TCSPC SPAD-based imaging

Valentin Poisson, William Guicquero, and Gilles Sicard

**Abstract**—Direct time-of-flight (D-ToF) image sensors based on Time-Correlated Single Photon Counting (TCSPC) systems face two main challenges, *i.e.*, the limitation of spatial resolution by the amount of data to store and the accuracy of the reconstruction under high background noise. The key contribution of this paper to overcome these issues consists in the introduction of a custom pixel-wise Compressive Sensing (CS) hardware implementation in combination with a deep learning reconstruction algorithm. Besides the reduction of the pixel pitch, the CS approach limits the probability of counters overflows, enabling larger photon counts operating mode, easing the depth/intensity reconstruction process in practice. Compared to prior work on deep learning models for TCSPC data, our proposed approach achieves a similar depth-intensity reconstruction accuracy in the typical low-photon flux mode of operation. However, when combined with the proposed CS hardware implementation compatible with high photon counting, our solution outperforms the most advanced SPAD sensing strategies as well as the best-in-class remote processing, both in terms of intensity-depth reconstruction performance and pixel pitch reduction.

**Index Terms**—Time-Correlated Single Photon Counting, Time-of-Flight imaging, SPAD image sensors, LiDAR, Compressive sensing, Cellular Automaton, Deep Generative Model.

## I. INTRODUCTION

Direct time-of-flight (D-ToF) sensors are now key devices for a wide range of imaging applications [1]–[3]. However, D-ToF measurements are subject to high level background noise illumination which poses a great challenge for state-of-the-art computational reconstruction algorithms for 3D imaging. To avoid noisy image reconstructions, traditional imaging algorithmic methods based on pixel-wise Maximum Likelihood (ML) estimation are proposed [4]. In addition, several works have also added spatial constraints, such as introducing a Total Variation (TV) regularization [5] or block-matching and 3-D filtering [6]. [7] even proposed an additional pixel-wise adaptive gating ML estimation method to discard photon detection times out of the gating interval. Furthermore, deep learning approaches [8], [9] that now outperform by far prior works [5]–[7], recently succeed to achieve a very high reconstruction fidelity under extremely low photon counts and a very low Signal-to-Background noise ratio (SBR).

D-ToF measurements noise is not the only major concern in the context of depth imaging. Indeed, although the depth sensors based on single-photon avalanche diodes (SPAD) can benefit from the most advanced 3D-stacking IC technologies

[10], some hardware-related limitations still remain. Those are mainly due to the Time-Correlated Single Photon Counting (TCSPC) data format itself. This data format consists in gathering ToF measurements (*i.e.*, round trip time of a light pulse whose laser source is synchronized with the sensor [11]) into histograms to further enables post-processing tasks. However, it intrinsically limits pixel pitch shrinkage, timing accuracy (*i.e.*, depth precision) and spatial/temporal resolution (*i.e.*, number of pixels and frame rate). In order to reduce the amount of TCSPC data, [12] proposed a method called Folded inter-frame Histogram (FiFH) which consists in building two smaller histograms, one representing the most significant bits (*i.e.*, a coarse temporal resolution) while the other is for the less significant bits. In addition, an extension to the FiFH method was proposed by [13] adding a control electronics to filter the second histogram using the estimated result from the first which increases the SNR and thus the accuracy of depth estimation. Similarly, [14] implemented an in-pixel zoom histogramming TDC architecture inspired by dichotomy partitioning where the bin equivalent duration is thus sequentially shrunk by half, requiring only two counters. Otherwise, the use of a signed Up/Down counter (UDC) by [15] further reduced the memory requirements. Unfortunately, although alleviating some hardware constraints, these methods imply a per-pixel complex SPAD scheduling while discarding most of the information except the histogram peak position.

On the other hand, in addition to being highly relevant in terms of depth reconstructions accuracy, deep learning data processing methods will most likely play an important role in overcoming SPAD hardware limitations (*i.e.*, overall data throughput, data storage, photon detection efficiency). For example, [16]–[19] developed Deep Learning algorithms –namely image-guided depth up-sampling– which consists in reconstructing a high resolution depth map from a full-scale RGB frame combined with its associated low resolution depth map, thereby indirectly relaxing hardware constraints. Besides, the Compressive Sensing (CS) strategy [20] has been investigated in the context of SPAD imaging in combination with DL data processing methods, by the use of a spatial light modulator (*e.g.*, a Digital Micro-mirror Device, DMD) in front of the sensor [21]–[24]. Unfortunately, optical-CS seems to be impractical for consumer electronics as long as they rely on the use of bulky optical systems which makes the system sensitive to process and temperature variations, while requiring complex and unstable calibrations.

Instead, this paper depicts a custom deep learning ap-

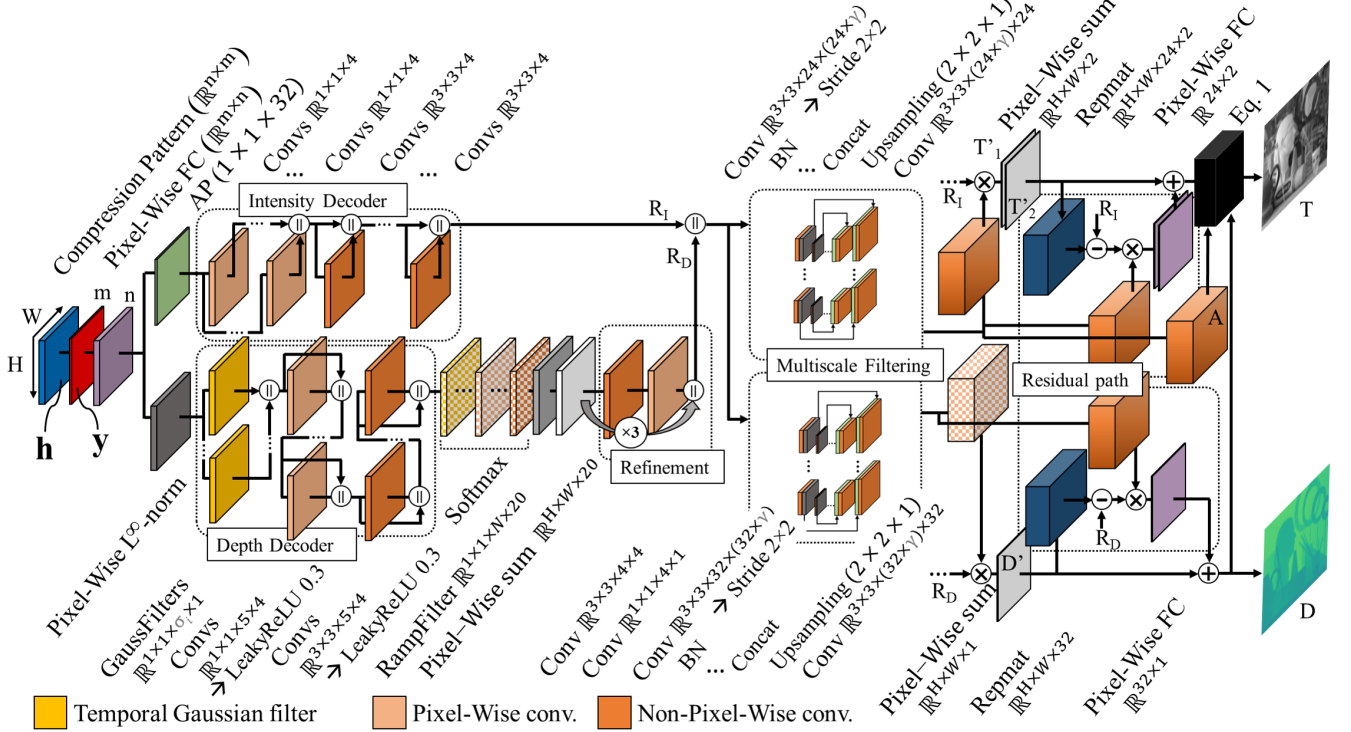


Fig. 1: Deep Neural Network topology for depth-intensity reconstruction from compressed histograms.

proach using a compression pattern aiming at converting high-dimensional raw data to low-dimensional data, replacing the canonical in-pixel TCSPC sensing scheme. The first section introduces a photon-efficient depth reconstruction that integrates a learned-compression pattern layer. This was first designed without taking into consideration hardware constraints, *i.e.*, only with the purpose of demonstrating the possible dimensionality reduction of the TCSPC histograms without much intensity-depth information loss. Then, the second section presents a pixel-wise histogram CS scheme from its mathematical formulation to its possible implementation using an in-pixel Cellular Automaton (CA), supplanting the combination of the counter-based TDC (*i.e.*, TCSPC) and the learned pattern linear projection. Finally, the two last sections provide hardware synthesis and experimental results to highlight the advantages and competitiveness of the proposed CS method in terms of pixel area optimization and reconstruction accuracy.

## II. DEEP NEURAL NETWORK COMBINED WITH A PIXEL-WISE TCSPC LEARNED PROJECTION

### A. Network architecture

To improve the performance of TCSPC LiDAR systems, especially under noisy measurements, this paper proposes a Deep Neural Network model that aims at jointly inferring depth,  $D$ , and intensity maps,  $T$ , from raw TCSPC SPAD data  $\mathbf{h}$ . This topology is divided into modules; the compression stage, the depth and intensity decoders, both providing multiple depth and intensity initial reconstructions, and the multiscale filtering (MF) that adaptively tune the weights of each reconstruction sub-channel. Last layers of this topology perform pixel-

wise reconstructions selections combined with residual skip connections among the depth and intensity decoders outputs filtered by the MF module (*i.e.*, respectively lower part and upper part of Fig. 1) which stack several reconstructions from various filters, gathering both local, region-based information and pixel-wise temporal information. Luminance/Intensity  $T$  reconstructions additionally rely on a physical model with assumptions on frame-based normalization.

Regarding the linear compression module (red block in Fig. 1), this layer converts high-dimensional TCSPC data to low-dimensional data in the shape of a pixel-wise vector aiming at preserving semantic information. In other words, this layer corresponds to a Conv2D1x1 in the sense of CNN, *i.e.*, a pixel-wise Fully Connected (FC) layer. Note that, this will be the layer considered as the CS layer in the next sections.

On the depth decoder side, to remove histogram scaling variability from one pixel to another, this module first embeds a pixel-wise  $L^\infty$  normalization (dark gray block in Fig. 1). Then, “Gaussian filters” in the temporal domain with various  $\sigma_i$  radius, ranging from 0 to 14 (N.B. pulse width is about 16 bins) is performed (yellow blocks in Fig. 1). A second filter, which is a per-pixel cascaded temporal filter, is then applied (orange blocks in Fig. 1). In order to take advantage of local temporal and spatial information simultaneously, 3D filters (Conv3D) are also cascaded afterwards (deep orange blocks in Fig. 1). The LeakyReLU activation is used rather than the canonical ReLU activation function in order to alleviate the “dying ReLU” problem. Every filter then provides a depth estimation output thanks to the combination of a Softmax (gridded block in Fig. 1) and a ramp layer (light gray in

Modality	SBR	Shin [5]	Rapp [7]	Peng [8]	Yao [9]	ours					
						m=256	m=64	m=32	m=16	m=12	m=8
Depth	10 : 2	0.0570	0.0479	0.0198	0.0182	0.0159	0.0207	0.0204	0.0176	0.0186	0.0202
	2 : 10	0.1906	0.0612	0.0494	0.0528	0.0596	0.0836	0.1016	0.0770	0.0893	0.1381
	5 : 50	0.2502	0.0592	0.0264	0.0241	0.0311	0.0349	0.0454	0.0347	0.0371	0.0473
	2 : 100	0.3188	0.0939	0.0308	0.0293	0.0357	0.0373	0.0482	0.0386	0.0394	0.0480
	Avg.	0.1849	0.0703	0.0335	0.0327	0.0382	0.0483	0.0586	0.0467	0.0509	0.0686
Intensity	10 : 2	12.68	16.91	<i>N/A</i>	<i>N/A</i>	19.96	21.54	20.52	20.28	20.38	20.70
	2 : 10	8.136	15.10	<i>N/A</i>	<i>N/A</i>	24.63	27.31	28.52	10.59	14.10	17.87
	5 : 50	8.708	15.95	<i>N/A</i>	<i>N/A</i>	25.41	28.00	27.76	13.45	18.98	21.52
	2 : 100	8.301	11.61	<i>N/A</i>	<i>N/A</i>	21.27	22.64	20.82	20.74	20.08	21.05
	Avg.	9.322	14.90	<i>N/A</i>	<i>N/A</i>	22.47	24.61	24.29	16.24	18.20	20.05

TABLE I: Quantitative comparisons of several Intensity-Depth reconstruction methods under various SBR. Note that depth results are reported as an average RMSE (m) and intensity results are reported as an average PSNR (in dB).

Fig. 1) that act as a trainable peak detector (as a softargmax [25] would do). Finally, three layers of Conv2D3x3 are used for an additional refinement stage.

For the intensity decoder module, since the original temporal dimension is considered irrelevant for intensity reconstructions with respect to the induced complexity on the topology, an average pooling (third-axis AP) layer (green block in Fig. 1) first reduces the scale of the input tensor in the temporal axis. Then in the same way as the depth decoder module, several convolutional layers are cascaded, each one providing intensity reconstructions that are concatenated in order to provide a set of intensity reconstructions. Note that for a proper behaviour of the convolutions, a custom 3D/2D padding has been implemented with a reflect padding instead of a zero padding for all modules of our model.

MF modules both take as input the intensity and depth maps provided by the decoders. It leverages local collaborations between pixels reconstructions through downscaling and then upscaling as a U-net structure [26] with an expanding coefficient value  $\gamma$  that relates to the number of channels increase at each scale. These modules for depth and intensity reconstructions selection incorporate, Conv2D3x3 with a stride of  $2 \times 2$  for downscaling steps (orange blocks in Fig. 1), and Conv2D3x3 with an  $2 \times 2 \times 1$  upsampling (US) (green blocks in Fig. 1). The MF attention module for depth selection ends with two separate Conv2D1x1 layers, one with a Softmax activation and the other with a Tanh activation respectively represented by a grid pattern and an orange-black gradient (see Fig. 1). The output of the softmax enables to further control the pixel-wise channel selection among the variety of depth reconstructions  $R_D$ , through a pixel-wise multiplication layer followed by a summation performed by a pixel-wise FC, providing the selected reconstruction  $D'$ . Besides, the output of the tanh activation function adds the depth reconstruction residual error to the selected reconstruction  $D'$ . Instead, the MF attention module for intensity estimation is terminated by three  $1 \times 1$  convolutional layers with tanh activation functions, in order to estimate the number of photons per pixel from the laser source  $T'_1$  and from the background illumination source  $T'_2$ . In the same fashion as for the depth module, the second tanh activation function performs the signal and noise photon

number estimation residual error to be added to the previous estimation  $T'_1$  and  $T'_2$ . An additional physically-driven layer is used for the final intensity estimation (illustrated by the black block in Fig. 1) taking as input the depth estimation  $D$ , the photons quantity estimated from laser pulse  $T'_1$  and from background illumination  $T'_2$  and the third convolutional layer outputs of the MF attention module. This black block embeds the knowledge of physical laws (cf. Eq. 1 where  $J$  is a all-ones matrix) that govern the given imaging measurement system dataset in the learning process (e.g., distance inverse-square law and frame-based intensity normalization). Finally, the last convolutional layer with tanh activation functions (denoted  $A$ ) performs a pixel-wise selection of the intensity estimator from the weighted sum of the background photon counts estimator or the laser pulse photon counts estimator. This pixel-wise selection enables an accurate intensity reconstruction even when there is few photon coming from the laser pulse due to the distance measurement inverse square law and conversely when there is few photon coming from background illumination.

$$T = A \odot \frac{T'_2}{T'_2} + (J - A) \odot \frac{T'_1}{T'_1} \odot \frac{D^2}{D^2}. \quad (1)$$

### B. Experimental results using learned-compression pattern

To highlight the interest of the proposed solution, a quantitative benchmark has been conducted on "pseudo-realistic" SPAD raw data generated from the simulation model presented in [8] with SPAD control asynchronicity. This simulation ToFs model encapsulates several physical parameters such as the laser source power, distance inverse-square law, scene point brightness, SPAD control asynchronicity, dark count rate [27] and background noise measurements [28]. The simulation model generates the SPAD TCSPC measurements train dataset from the NYU V2 dataset [29], and the test dataset from the Middlebury dataset [30]. To ensure a proper matching between the train dataset and the test dataset, and to avoid unnecessary training for long range depths, all images whose dynamic range is greater than the one of the test dataset have been removed. In addition, a data augmentation is performed, consisting in generating 12 samples of each NYU image with a dynamic range below 3m under 12 SBRs (i.e., an average of 1, 2, 3, 5, 10 signal photon counts and 2, 10, 50, 100 noise photon

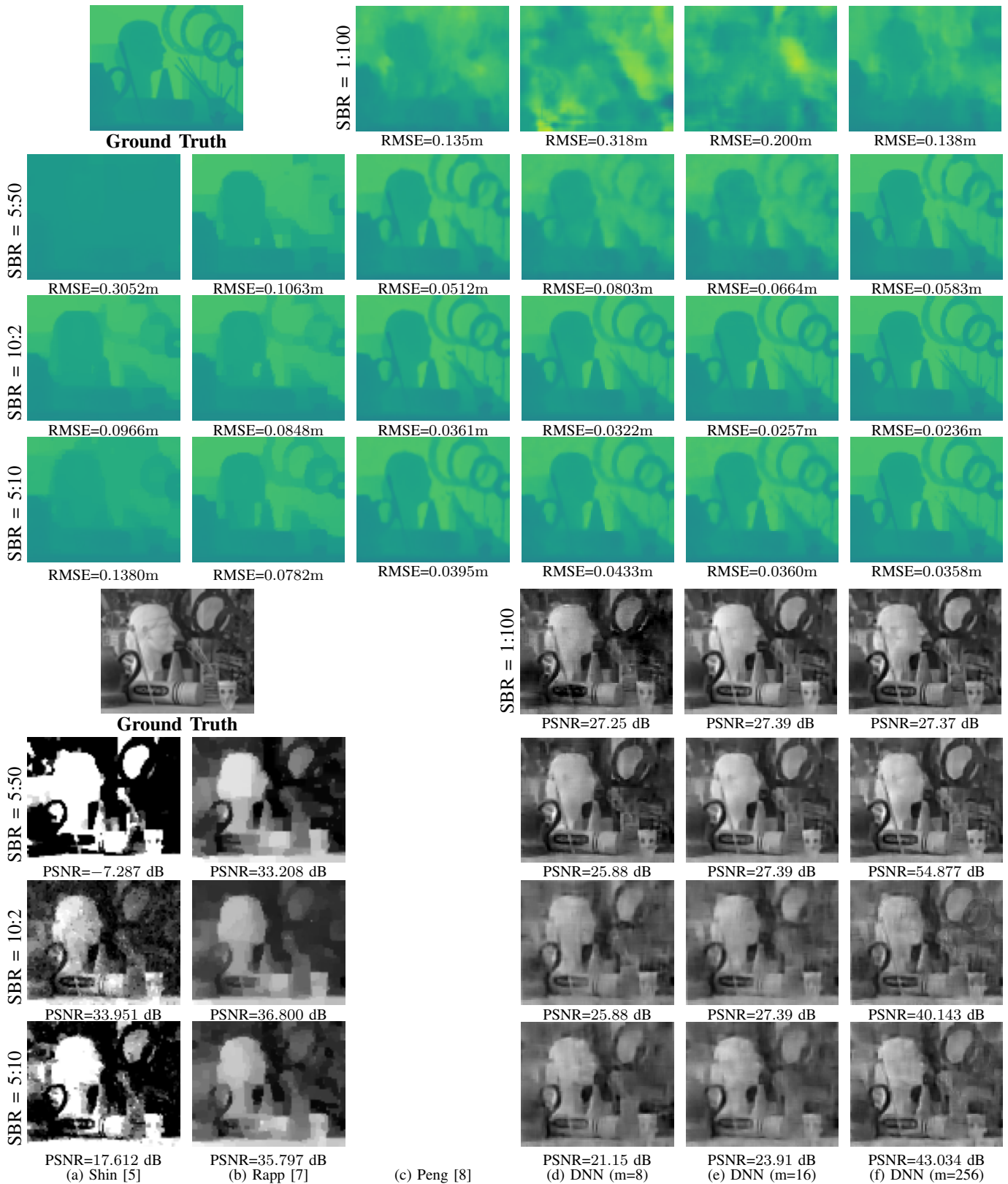


Fig. 2: Intensity-Depth reconstructions for various SBR, reported with RMSE (in meters) and PSNR (in dB) metrics. Note that, [8] does not provide any estimation of the intensity and under the 1:100 SBR [5] and [7] completely fail for both tasks.

counts). The simulation model finally provides a 3D output data volume with the 2D pixel array resolution and a third axis of 256 values (*i.e.*, a time bin resolution of 80ps and a dynamic range of 20ns (6m)).

This subsection therefore reports our proposed DNN performances in comparison to prior works [5], [7]–[9] without any further hardware considerations (*i.e.*, considering an optimized learned-compression pattern) and under low photon counts. For a fair comparison, our proposed DNN and that of [8], [9] were trained with the same training data *i.e.*,  $64 \times 64 \times 256$  tensors and similarly to the original work, *i.e.*, using PyTorch, with a batch size of 4, a random initialization, an Adam optimizer [31], and a learning rate of  $10^{-4}$  with a learning rate decay of 0.9 after each epoch. Instead, our DNN was trained using TensorFlow2, with a batch size of 8 and a number of epochs of 40 with a learning rate of  $10^{-3}$  and a decay of 0.95 after each epoch, starting from the  $10^{\text{th}}$  epoch. Note that for a proper convergence of our model, the training process of our neural network is performed in a 2-stage fashion so that the ‘‘Gaussian filters’’ and the ramps are made trainable only during the second stage (*i.e.*, 80 epochs in total). For the sake of simplicity, a Mean Squared Error (MSE) loss has been used with the Adam optimizer for all experiments reported in this paper. In order to limit the proposed DNN model size,  $\gamma$  is fixed to 1.5.

Even though the proposed DNN only achieves an average depth reconstructions root mean square error (RMSE) of 0.0382 m, *i.e.*, a similar accuracy to [8], [9] *cf.*, Tab. I. It provides decent depth reconstructions (*cf.* Fig.2) in comparison to [5] and [7] with respectively a 80% and 46% lower average depth RMSE for  $m = 256$  (uncompressed case). With a compression corresponding to  $m = 32$ , our method still improves the depth reconstruction performance by respectively 69% and 18%, compared to [5], [7]. Note that, our method clearly outperforms [5], [7] for any SBR, except in the 2:10 SBR case. In addition, an acceptable intensity reconstruction is obtained for  $m \geq 16$ , especially when compared to prior works. These results pave the way to possible data dimensionality reduction through the use of a learned compression pattern (*i.e.*, a linear projection represented by the red block in Fig. 1). Consequently, to reduce the pixel pitch in practice, while enabling the high photon flux operating mode, we propose to implement a data-agnostic linear projection in the name of a CS scheme, instead of a learned compression pattern (*cf.* the following section).

### III. DATA-AGNOSTIC COMPRESSION PATTERN USING SHUFFLED CELLULAR AUTOMATA

Due to the limitation related to the one-hot type of encoding in standard TCSPC systems, the second goal of this paper is to replace the learned-compression pattern in the DNN topology in Sec. II by an on-the-fly data-agnostic CS hardware implementation. Let us recall that LiDAR systems based on the TCSPC principle [32] make repeated measures of the propagation time of a light pulse emitted

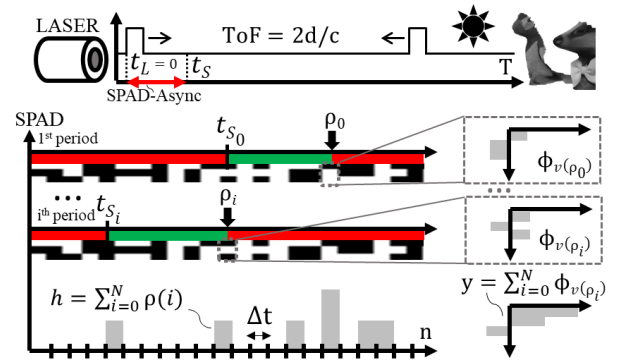


Fig. 3: SPAD Operation System Overview: direct time-of-flight measurement (D-ToF) of a light pulse reflected by a target using a TCSPC or a CS-TCSPC system.

by a transceiver (*i.e.*, pulsed laser) and then received by a receiver (*i.e.*, SPAD sensor). These propagation time of flights (ToF) are then stored in the shape of a histogram. Prior works [5], [7], [8], [16] commonly consider the low photon counts mode of operation allowing to neglect the pile-up effect but is very restrictive and constraining in terms of reconstruction performance. In fact, SPAD pixels can only detect the first incident photon during each laser cycle, after which it enters into a dead time (*i.e.*, any further photons can be detected). In high photon counts mode of operation, the measured histograms thus exhibit an exponentially decaying shape leading to complex reconstruction algorithms. However, this unwanted histogram shaping issue can be efficiently bypassed thanks to an asynchronous SPAD control [33]. This recent technique consists in temporally misaligning the SPAD measurement windows with the TDC and the pulse laser –still synchronous to each other– by a constant circular shift  $t_{S_i}$  (*cf.*, Fig. 3), in order to smooth out pile-up distortions [34]. Asynchronous control furthermore relaxes hardware constraints because of reducing the probability of counter overflows by spreading out the histogram; thus motivating its consideration for this work.

Thanks to this asynchronous SPAD control, the resulting histogram model exhibits pseudo-sparsity (*i.e.*, is said a compressible signal because only few coefficients have high magnitude values) under various background photon counts as illustrated in Fig. 4 (a). This condition is required to fulfill the CS theoretical background [20] and to properly take advantage of a ‘‘chaotic’’ encoding at sensor level as presented in [35] and [36]. Therefore, mathematically speaking, we proposed to replace the first learned-compression pattern layer (red color block in Fig. 1) by an untrained linear projection of a measured ToF histogram  $\mathbf{h} \in \mathbb{N}^n$ . This is performed using the matrix  $\Phi \in \{-1, +1\}^{m \times n}$  (Rademacher-like distribution considered for its universality property [37]), providing the measurement vector  $\mathbf{y} \in \mathbb{Z}^m$  (Eq. 2). The choice of a signed modulation also limits the probability of counters overflows since all the bins of the histogram will count more or less the same amount of noise-related photons, which will lead to a value of the counters centered around zero (*i.e.*, assumed to be similar to

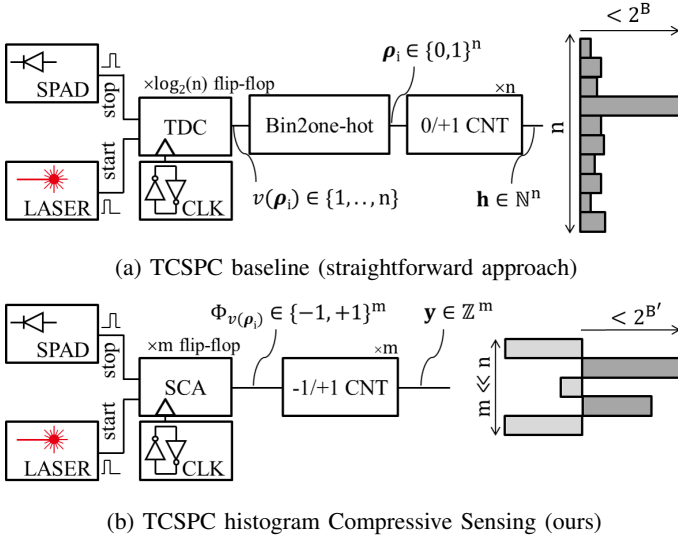


Fig. 4: System-level views of baseline (a) and its variant (b).

a 0-centered Gaussian distribution).

$$\mathbf{y} = \Phi \mathbf{h} \quad (2)$$

From the sensor point-of-view, a raw D-ToF measurement  $v(\rho_i) \in \{1, \dots, n\}$  provided by a TDC is a time index corresponding to a one-hot vector  $\rho_i \in \{0, 1\}^n$  for the  $i^{\text{th}}$  ToF acquisition,  $i \in \{1, \dots, N\}$ . In order to actually build the pixel-wise histogram  $\mathbf{h}$ , this vector  $\rho_i$  only needs to be summed over  $N$  TDC successive acquisitions thanks to its intrinsic position coding (*i.e.*,  $\mathbf{h} = \sum_{i=1}^N \rho_i$ ).  $\rho_i$  has a unique non-zero coefficient equal to 1 at the position  $v(\rho_i)$  and knowing that the multiplication by  $\Phi$  is distributive with respect to the addition; building  $\mathbf{y}$  is equivalent to calculating the sum over  $i$  of the columns  $\Phi_{v(\rho_i)}$  of  $\Phi$  at  $v(\rho_i)$  positions (cf. Eq. 3).

$$\mathbf{y} = \Phi \sum_{i=1}^N \rho_i = \sum_{i=1}^N \Phi_{v(\rho_i)} \quad (3)$$

It thus advantageously enables a direct acquisition of CS measurements without the need for an explicit representation of  $\mathbf{h}$  at any time, making the approach highly relevant in terms of hardware implementation. It means that the only requirement is to replace the one-hot encoding of the measured ToF by a said "chaotic" encoding (*i.e.*, the columns of  $\Phi$ ,  $\Phi_{v(\rho_i)}$ ). To our knowledge, the most appropriate solution is to use a basic Cellular Automaton (CA) replacing the commonly used TDC plus one-hot encoding. Indeed, CAs [38] have been employed in a wide range of systems and have recently become practical candidates to enable on-the-fly generations of CS matrices [39]–[41]. Note that a binary CA is composed of a finite number of cells that have a single binary state at each cell and discrete time step. For a regular Elementary CA (ECA), each cell state only depends on a logic rule taking as inputs the previous states of the cell itself and its two neighbors. The main advantage of CAs is that a complex global behavior can be obtained using only very few digital logic gates. Since

the CA is dedicated to the generation of a pseudo-random sequence, a slight modification has been made to the structure of the ECA rule 30. As depicted in Fig. 5, a simple routing for static shuffling is added to further increase the statistical independence of the binary states produced.

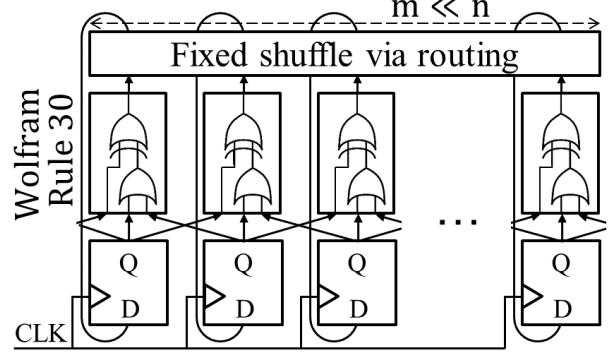


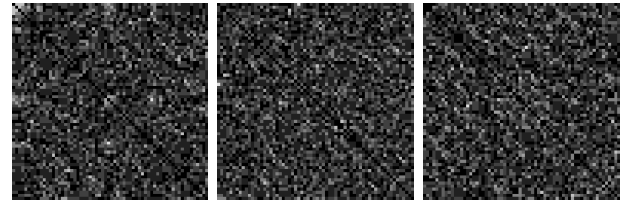
Fig. 5: Structure of a Shuffled Cellular Automaton (SCA).

As presented in Fig. 4 and according to Eq. 3, the construction of the histogram  $\mathbf{h}$  can thus be formally replaced by the direct construction of a CS measurement vector  $\mathbf{y}$ . This way, the number of clock cycles between the laser shot and the SPAD trigger (*i.e.*, ToF index) is equal to the previously denoted  $v(\rho_i)$  and the state vector of the SCA is considered to be equal to  $\Phi_{v(\rho_i)}$  in its signed representation.



(a) Rule 30 with canonical initialization. (b) Rule 30 with random initialization. (c) Rule 30 SCA with random initialization.

Fig. 6: Cell states for various CA configurations ( $\sim \Phi$ ).



(a)  $\mu_M = 0.875$ ,  $\mu_A = 0.1959$  (b)  $\mu_M = 0.75$ ,  $\mu_A = 0.1956$  (c)  $\mu_M = 0.75$ ,  $\mu_A = 0.1938$

Fig. 7: Normalized Gram matrix ( $\sim \Phi^T \Phi$  with zero diagonal for proper rendering) of the sensing matrices as presented in Fig. 6, where  $\mu_M$  corresponds to the mutual coherence and  $\mu_A$  is the average of the Gram matrix except the diagonal.

The same way as a TDC starts counting from zeros, the SCA starts updating its states from the initial state triggered by the laser shot and stops on the SPAD trigger. However, this deterministic process (*i.e.*, the state sequence (columns of  $\Phi$ ) only depends on the logic rule, the shuffle map and the initial states, see Fig. 5) requires a proper choice of the Wolfram's rule, the initialization and the shuffle pattern,

to ensure having the longest possible cycle to provide a full rank  $\Phi$ . Even if the CS community mainly focused on the Class 3 cellular automaton rule 30 which exhibits a proper chaotic behavior [42], using a shuffle stage does not imply any additional hardware component (only based on CA cells interconnections) while improving the properties of the generated pseudo-random vectors (see Fig. 6 for cell states ( $\sim \Phi$ ) and Fig. 7 for Gram matrices ( $\sim \Phi^T \Phi$ )).

#### IV. HARDWARE SYNTHESIS RESULTS

Both hardware architectures presented in Fig. 4 were synthesized using SYNOPSIS<sup>®</sup> Design-Compiler<sup>™</sup> over a 40nm standard cell technology to obtain the area results reported in Tab. II and Tab. III. For the sake of clarity, synthesis-estimated digital areas are given in equivalent pixel pitches ( $\mu\text{m}$ ) instead of areas ( $\mu\text{m}^2$ ). Our synthesis results are to be put in perspective against the exceedingly well optimized design of [32] that presents a tiled full custom designed histogram builder with a pitch of  $36.72\mu\text{m}$  targeting 16 bins of  $B=14$  bitwidth each. Those results which are obtained using an equivalent technology node show a size difference factor of 1.2 (*i.e.*,  $44.9\mu\text{m}$  for the same configuration as in [32]). This gap would be bridged by redesigning our architecture using a full custom design methodology with specific attention to area minimization. The main takeaway is that the relative results between Tab. II and Tab. III demonstrate that SCA configurations imply a marginal size increase (lower than 10%) compared to an uncompressed TCSPC baseline that explicitly builds a ToF histogram for the same number of measurements (*i.e.*,  $m=n$ ) and with identical counters bitwidths.

$n \setminus B$	4	5	6	7	8	9	10	11	12	13	14
8	19.1	20.8	22.4	23.9	25.3	26.6	27.9	29.1	30.3	31.4	32.5
16	25.3	27.9	30.3	32.5	34.5	36.5	38.3	40.1	41.8	43.4	44.9
256	97.8	108	118	127	136	143	151	158	165	171	178

TABLE II: Equivalent pixel pitch ( $\mu\text{m}$ ) required obtained from TCSPC baseline architecture synthesis (Fig. 4 (a)).

$m \setminus B$	4	5	6	7	8	9	10	11	12	13	14
8	15.3	19.9	22.0	23.9	25.6	27.3	28.8	30.3	31.6	33.0	34.3
16	21.6	28.1	31.1	33.7	36.1	38.5	40.7	42.8	44.8	46.7	48.6

TABLE III: Equivalent pixel pitch ( $\mu\text{m}$ ) required obtained from SCA architecture synthesis (Fig. 4 (b)).

#### V. EXPERIMENTAL RESULTS USING DNN COMBINED WITH DATA AGNOSTIC COMPRESSIVE SENSING

This last section presents the CS design performances in comparison to existing, SPAD-optimized sensing designs, *i.e.*, the SiFH [13] and FiFH [12] methods. For the sake of fair comparisons, we propose to replace the canonical argmax peak detection usually used for FiFH and SiFH by our own DNN reconstruction algorithm in order to properly characterize the CS scheme itself, decorrelated from the effect of the reconstruction strategy. Note that one CS input channel is replaced by a dummy photon counter. The proposed DNN model was therefore trained with FiFH [12], SiFH

SBR	Fifh (m=16)		Sifh (m=16)		CS (m=8)		CS (m=16)	
	B=5	B=7	B=5	B=7	B=5	B=7	B=5	B=7
40 : 2000	0.40	0.04	0.33	0.03	0.08	0.09	0.05	0.05
40 : 4000	0.64	0.15	0.38	0.14	0.11	0.11	0.06	0.06
120 : 4000	0.62	0.14	0.38	0.12	0.07	0.05	0.03	0.04
Avg.	0.51	0.09	0.34	0.07	<b>0.07</b>	0.06	<b>0.03</b>	0.04
Pitch ( $\mu\text{m}$ )	28	32	28	32	<b>20</b>	27	<b>28</b>	34

TABLE IV: Depth comparisons of several acquisition methods under various SBR, reported as an average RMSE in m.

SBR	Fifh (m=16)		Sifh (m=16)		CS (m=8)		CS (m=16)	
	B=5	B=7	B=5	B=7	B=5	B=7	B=5	B=7
40 : 2000	12.2	33.1	28.2	10.7	28.5	25.5	35.1	29.5
40 : 4000	10.7	20.6	27.5	12.6	27.1	26.5	37.4	30.6
120 : 4000	10.7	20.8	27.6	12.9	29.9	26.5	38.5	31.7
Avg.	11.4	26.2	28.0	13.4	<b>29.2</b>	26.9	<b>36.9</b>	31.1
Pitch ( $\mu\text{m}$ )	28	32	28	32	<b>20</b>	27	<b>28</b>	34

TABLE V: Intensity comparisons of several acquisition methods under various SBR, reported as an average PSNR in dB.

[13] as well as with CS data inputs, in the same way as described in sec. II, except for the considered SBRs. The data augmentation here consists in the generation of 9 samples for each image of the NYU dataset having a dynamic range below 3m (under 9 SBRs, *i.e.*, an average of 40, 80, 120, 200, 400 signal photon counts and 2000, 4000 noise photon counts).

In case of the absence of counter overflows, Tabs. IV and V reports that FiFH and SiFH [12], [13] can outperform our solution in some conditions. However, considering a practical bit depth, our proposed design shows better performances than [12], [13]. Even with same reconstruction algorithm, histogram CS allows to reach a highly accurate reconstruction, while being able to reduce both the size of measurement vectors and their bitwidths ( $B=5$ ), therefore reducing the total memory needs (acting on the pixel pitch, cf. section IV). On the other hand, when decreasing the bitwidth (*i.e.*, to  $B=5$ ), the performances of [12], [13] are significantly downgraded cf., Fig. 8 and Fig. 9. These Tabs. also highlight that our CS approach provides a better depth estimation RMSE compared to [12], [13] ( $B=7$ ), with an estimated pitch reduction of 39%. It leads to the conclusion that –in terms of depth estimation accuracy and for the considered configurations– our proposed system improves by a 11x factor the depth RMSE at iso-surface (a pitch of  $28\mu\text{m}$ ) or a pixel surface reduction by a 2x factor at iso-performance (a RMSE of 0.07m). Similar conclusions can be drawn when considering intensity reconstruction results.

Finally, Fig. 10 puts into perspectives the results with a learned-compression pattern under low photon counts with the ones obtained using the proposed CS scheme, under high photon counts. Although the SCA implementation highly constrains the measurement vector with a small bitwidth of  $B=5$ , counters do not much overflow even under high photon counts. Fig. 10 demonstrates that the proposed DNN combined with the CS provides more accurate depth/intensity reconstructions than the DNN topology integrating learned-compression pattern (if saturation issues are not taken into consideration directly during the training stage).



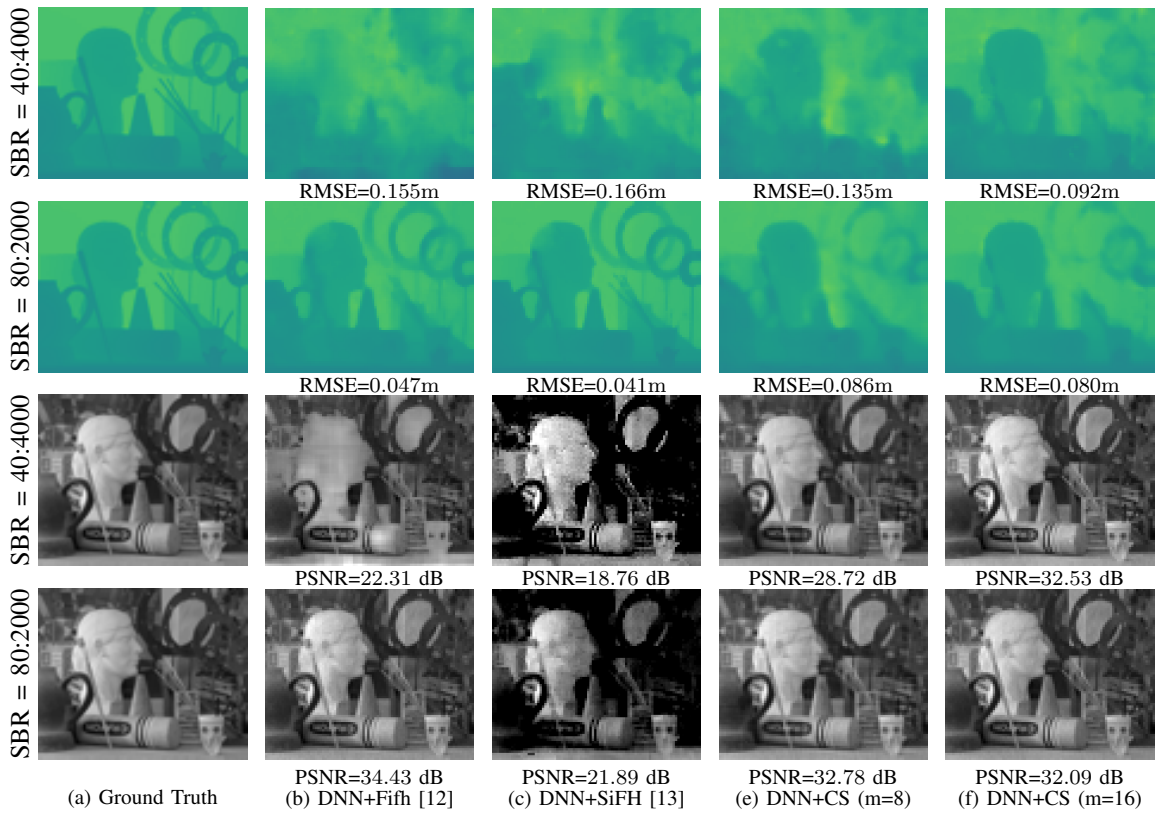


Fig. 8: Intensity-Depth reconstructions under hardware constraints with counters bitwidth of  $B=7$  and high photon counts.

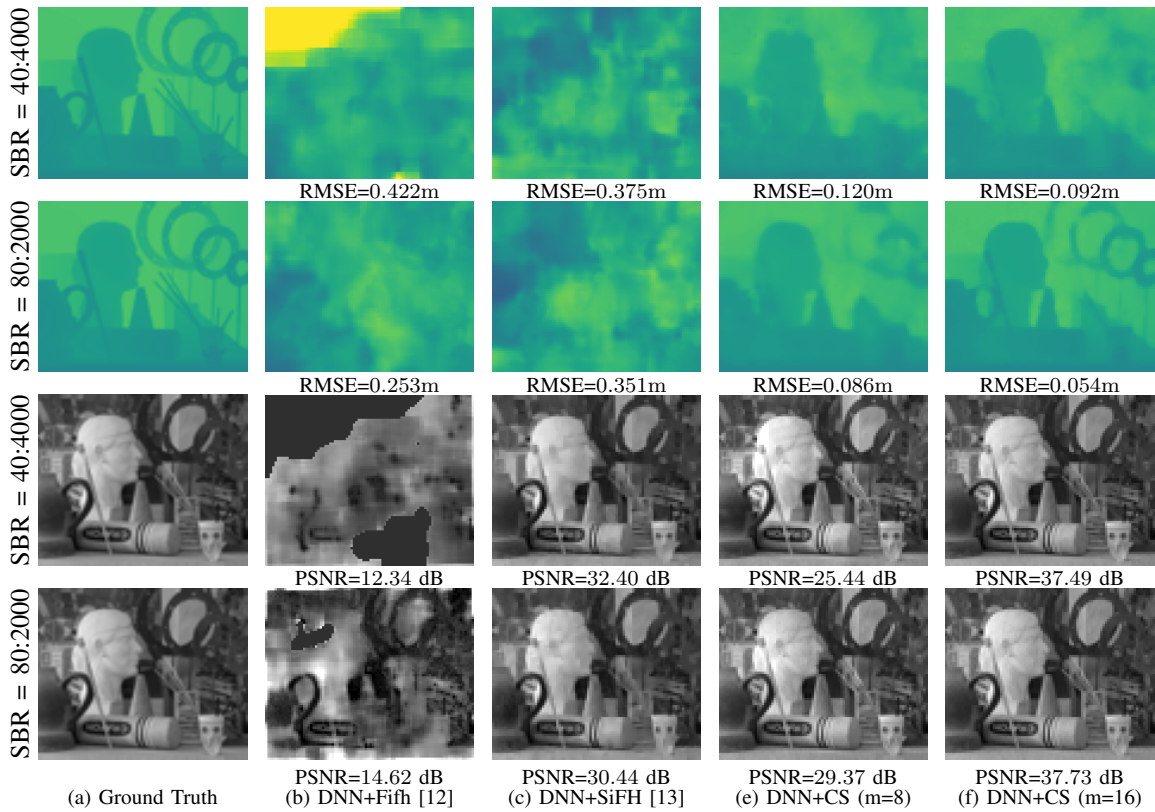


Fig. 9: Intensity-Depth reconstructions under hardware constraints with counters bitwidth of  $B=5$  and high photon counts.

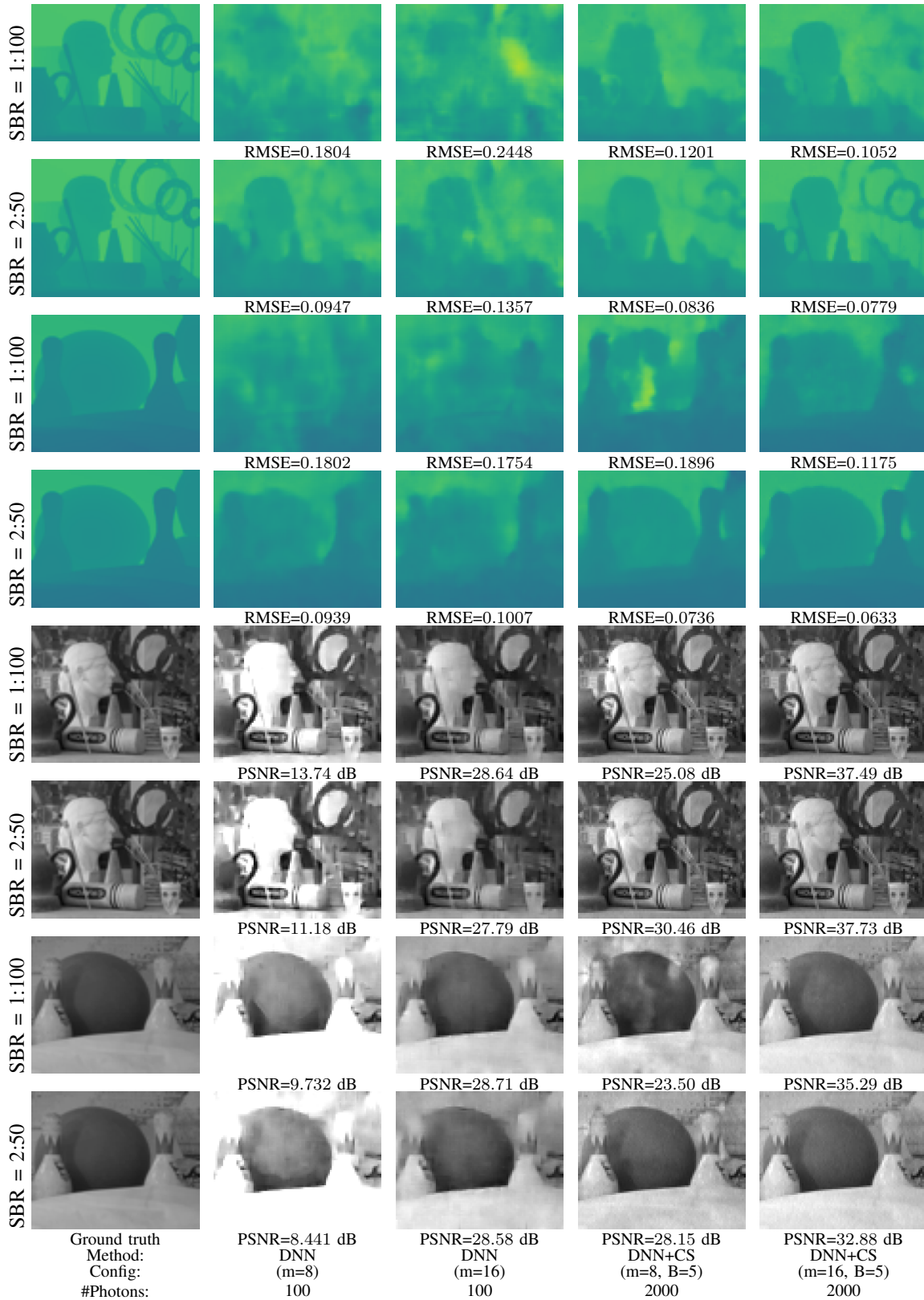


Fig. 10: Intensity-Depth reconstructions. Note that the 2<sup>nd</sup> and 3<sup>rd</sup> columns reconstructions are under low photon counts (an average photon counts of 100) due to hardware constraints related to the TCSPC data format. While the 4<sup>th</sup> and 5<sup>th</sup> columns reconstructions are under high photon counts (an average photon counts of 2000) thanks to the CS hardware design.

## VI. CONCLUSION

This paper introduces a Deep Neural Network (DNN) with a linear compression front-end whose goal is to demonstrate a possible input data dimensionality reduction without too much information loss. Under the low photon counts SPAD operating mode, this work typically reports state-of-the-art results with an average RMSE depth reconstruction loss of only 0.016 m and an average PSNR intensity reconstruction increase of 10 dB in comparison to prior works.

Subsequently, a Compressive Sensing (CS) hardware implementation scheme replacing the DNN learned-compression pattern layer is proposed, enabling to relax hardware constraints on the SPAD sensor, at the pixel level. The CS design consequently reduces the number of memory words as well as the number of bits per words, thereby reducing the pixel pitch required for an in-pixel implementation with the use of a pixel-wise shuffled Cellular Automaton (CA). Based on the same reconstruction algorithm, the proposed acquisition method thus allows a higher reconstruction performance compared to existing, most efficient sensing methods [12], [13] with an estimated pitch reduction of approximately 40%. Finally, the CS scheme avoids counters overflows even under high photon counts, which, in combination with the proposed DNN, provides a higher reconstruction accuracy than the best-in-class remote processing work that limits its mode of operation to low photon counts.

## REFERENCES

- [1] A. Carimatto, S. Mandai, E. Venialgo, T. Gong, G. Borghi, D. R. Schaart, and E. Charbon, "11.4 A 67,392-SPAD PVTB-compensated multi-channel digital SiPM with 432 column-parallel 48ps 17b TDCs for endoscopic time-of-flight PET." *IEEE*, Feb. 2015, pp. 1–3.
- [2] A. R. Ximenes, P. Padmanabhan, M. Lee, Y. Yamashita, D. N. Young, and E. Charbon, "A 256x256 45/65nm 3d-stacked SPAD-based direct TOF image sensor for LiDAR applications with optical polar modulation for up to 18.6db interference suppression," in *2018 IEEE International Solid - State Circuits Conference - (ISSCC)*, Feb. 2018, pp. 96–98.
- [3] L. Zhang, D. Chitnis, H. Chun, S. Rajbhandari, G. Faulkner, D. O'Brien, and S. Collins, "A Comparison of APD- and SPAD-Based Receivers for Visible Light Communications," *Journal of Lightwave Technology*, vol. 36, no. 12, pp. 2435–2442, Jun. 2018.
- [4] J. Tachella, Y. Altmann, X. Ren, A. McCarthy, G. S. Buller, J.-Y. Tourneret, and S. McLaughlin, "Bayesian 3d Reconstruction of Complex Scenes from Single-Photon Lidar Data," *SIAM Journal on Imaging Sciences*, vol. 12, no. 1, pp. 521–550, Jan. 2019, arXiv: 1810.11633.
- [5] D. Shin, A. Kirmani, V. K. Goyal, and J. H. Shapiro, "Photon-Efficient Computational 3-D and Reflectivity Imaging With Single-Photon Detectors," *IEEE Transactions on Computational Imaging*, vol. 1, no. 2, pp. 112–125, Jun. 2015.
- [6] K. Dabov, A. Foi, V. Katkovnik, and K. Egiazarian, "Image denoising by sparse 3-d transform-domain collaborative filtering," *IEEE Transactions on Image Processing*, vol. 16, no. 8, pp. 2080–2095, 2007.
- [7] J. Rapp and V. K. Goyal, "A Few Photons Among Many: Unmixing Signal and Noise for Photon-Efficient Active Imaging," *IEEE Transactions on Computational Imaging*, vol. 3, no. 3, pp. 445–459, Sep. 2017, arXiv: 1609.07407.
- [8] J. Peng, Z. Xiong, X. Huang, Z.-P. Li, D. Liu, and F. Xu, "Photon-Efficient 3d Imaging with A Non-Local Neural Network," Aug. 2020.
- [9] G. Yao, Y. Chen, Y. Liu, X. Hu, and Y. Pan, "Robust photon-efficient imaging using a pixel-wise residual shrinkage network," *arXiv preprint arXiv:2201.01453*, 2022.
- [10] E. Charbon, C. Bruschini, and M. Lee, "3d-Stacked CMOS SPAD Image Sensors: Technology and Applications," in *2018 25th IEEE International Conference on Electronics, Circuits and Systems (ICECS)*, Dec. 2018, pp. 1–4.
- [11] F. Arvani and T. C. Carusone, "Direct Time-of-Flight TCSPC Analytical Modeling Including Dead-Time Effects," in *2018 IEEE International Symposium on Circuits and Systems (ISCAS)*, May 2018, pp. 1–4.
- [12] I. Vornicu, A. Darie, R. Carmona-Galán, and Á. Rodríguez-Vázquez, "Compact Real-Time Inter-Frame Histogram Builder for 15-Bits High-Speed ToF-Imagers Based on Single-Photon Detection," *IEEE Sensors Journal*, vol. 19, no. 6, pp. 2181–2190, Mar. 2019.
- [13] I. Vornicu, A. Darie, R. Carmona-Galan, and Á. Rodríguez-Vázquez, "Tof estimation based on compressed real-time histogram builder for spad image sensors," in *2019 IEEE International Symposium on Circuits and Systems (ISCAS)*. IEEE, 2019, pp. 1–4.
- [14] B. Kim, S. Park, J.-H. Chun, J. Choi, and S.-J. Kim, "7.2 A 48x40 13.5mm Depth Resolution Flash LiDAR Sensor with In-Pixel Zoom Histogramming Time-to-Digital Converter," in *2021 IEEE International Solid-State Circuits Conference (ISSCC)*, vol. 64, Feb. 2021, pp. 108–110.
- [15] S. Park, B. Kim, J. Cho, J.-H. Chun, J. Choi, and S.-J. Kim, "An 80x 60 flash lidar sensor with in-pixel histogramming tdc based on quaternary search and time-gated  $\delta$ -intensity phase detection for 45m detectable range and background light cancellation," in *2022 IEEE International Solid-State Circuits Conference (ISSCC)*, vol. 65. IEEE, 2022, pp. 98–100.
- [16] L. B. O'Toole, Matthew, and Wetzstein, Gordon, "Single-photon 3d imaging with deep sensor fusion," *ACM Transactions on Graphics (TOG)*, Jul. 2018.
- [17] T.-W. Hui, C. C. Loy, and X. Tang, "Depth Map Super-Resolution by Deep Multi-Scale Guidance," in *Computer Vision – ECCV 2016*, ser. Lecture Notes in Computer Science. Springer, Cham, Oct. 2016, pp. 353–369.
- [18] D. Ferstl, C. Reinbacher, R. Ranftl, M. Ruether, and H. Bischof, "Image Guided Depth Upsampling Using Anisotropic Total Generalized Variation," in *2013 IEEE International Conference on Computer Vision*, Dec. 2013, pp. 993–1000.
- [19] A. Ruget, S. McLaughlin, R. K. Henderson, I. Gyongy, A. Halimi, and J. Leach, "Robust super-resolution depth imaging via a multi-feature fusion deep network," *arXiv.org*, Nov. 2020.
- [20] E. J. Candes and M. B. Wakin, "An Introduction To Compressive Sampling," *IEEE Signal Processing Magazine*, vol. 25, no. 2, pp. 21–30, Mar. 2008.
- [21] A. Colaço, A. Kirmani, G. A. Howland, J. C. Howell, and V. K. Goyal, "Compressive depth map acquisition using a single photon-counting detector: Parametric signal processing meets sparsity," in *2012 IEEE Conference on Computer Vision and Pattern Recognition*, Jun. 2012, pp. 96–102, iSSN: 1063-6919.
- [22] G. A. Howland, D. J. Lum, M. R. Ware, and J. C. Howell, "Photon counting compressive depth mapping," *Optics Express*, vol. 21, no. 20, pp. 23 822–23 837, Oct. 2013.
- [23] Q. Sun, X. Dun, Y. Peng, and W. Heidrich, "Depth and Transient Imaging with Compressive SPAD Array Cameras," in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, Jun. 2018, pp. 273–282.
- [24] A. Farina, A. Farina, A. Candeo, A. D. Mora, A. Bassi, A. Bassi, R. Lussana, F. Villa, G. Valentini, G. Valentini, S. Arridge, C. D'Andrea, and C. D'Andrea, "Novel time-resolved camera based on compressed sensing," *Optics Express*, vol. 27, no. 22, pp. 31 889–31 899, Oct. 2019.
- [25] Z. Sun, D. B. Lindell, O. Solgaard, and G. Wetzstein, "SPADnet: deep RGB-SPAD sensor fusion assisted by monocular depth estimation," *Optics Express*, vol. 28, no. 10, pp. 14 948–14 962, May 2020.
- [26] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *International Conference on Medical image computing and computer-assisted intervention*. Springer, 2015, pp. 234–241.
- [27] S. Isaak, M. Pitter, S. Bull, and I. Harrison, "Design and characterisation of 16x1 parallel outputs SPAD array in 0.18 um CMOS technology," in *2010 IEEE Asia Pacific Conference on Circuits and Systems*, Dec. 2010, pp. 979–982.
- [28] S. Jahromi, J. Jansson, P. Keränen, and J. Kostamovaara, "A 32 x 128 SPAD-257 TDC Receiver IC for Pulsed TOF Solid-State 3-D Imaging," *IEEE Journal of Solid-State Circuits*, vol. 55, no. 7, pp. 1960–1970, Jul. 2020.
- [29] N. Silberman, D. Hoiem, P. Kohli, and R. Fergus, "Indoor Segmentation and Support Inference from RGBD Images," in *Computer Vision –*

- ECCV 2012*, ser. Lecture Notes in Computer Science. Springer, Berlin, Heidelberg, Oct. 2012, pp. 746–760.
- [30] H. Hirschmuller and D. Scharstein, “Evaluation of Cost Functions for Stereo Matching,” in *2007 IEEE Conference on Computer Vision and Pattern Recognition*, Jun. 2007, pp. 1–8.
- [31] D. P. Kingma and J. Ba, “Adam: A Method for Stochastic Optimization,” *arXiv.org*, Dec. 2014.
- [32] S. W. Hutchings, N. Johnston, I. Gyongy, T. Al Abbas, N. A. W. Dutton, M. Tyler, S. Chan, J. Leach, and R. K. Henderson, “A Reconfigurable 3-D-Stacked SPAD Imager With In-Pixel Histogramming for Flash LIDAR or High-Speed Time-of-Flight Imaging,” *IEEE Journal of Solid-State Circuits*, vol. 54, no. 11, pp. 2947–2956, Nov. 2019.
- [33] A. Gupta, A. Ingle, and M. Gupta, “Asynchronous Single-Photon 3d Imaging,” in *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, Oct. 2019, pp. 7908–7917.
- [34] E. Sarbazi, M. Safari, and H. Haas, “Statistical Modeling of Single-Photon Avalanche Diode Receivers for Optical Wireless Communications,” *IEEE Transactions on Communications*, vol. 66, no. 9, pp. 4043–4058, Sep. 2018.
- [35] V. Poisson, T. V. Nguyen, W. Guicquero, and G. Sicard, “Luminance-depth reconstruction from compressed time-of-flight histograms,” *IEEE Transactions on Computational Imaging*, pp. 1–1, 2022.
- [36] F. Gutierrez-Barragan, A. Ingle, T. Seets, M. Gupta, and A. Velten, “Compressive single-photon 3d cameras,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2022, pp. 17 854–17 864.
- [37] P. L. Bartlett and S. Mendelson, “Rademacher and gaussian complexities: Risk bounds and structural results,” *Journal of Machine Learning Research*, vol. 3, no. Nov, pp. 463–482, 2002.
- [38] K. Bhattacharjee, N. Naskar, S. Roy, and S. Das, “A survey of cellular automata: types, dynamics, non-uniformity and applications,” *Natural Computing*, vol. 19, no. 2, pp. 433–461, Jun. 2020.
- [39] W. Benjlali, W. Guicquero, L. Jacques, and G. Sicard, “Hardware-Compliant Compressive Image Sensor Architecture Based on Random Modulations and Permutations for Embedded Inference,” *IEEE Transactions on Circuits and Systems I: Regular Papers*, vol. 67, no. 4, pp. 1218–1231, Apr. 2020.
- [40] W. Guicquero, A. Dupret, and P. Vandergheynst, “An algorithm architecture co-design for cmos compressive high dynamic range imaging,” *IEEE Transactions on Computational Imaging*, vol. 2, no. 3, pp. 190–203, 2016.
- [41] M. Trevisi, A. Akbari, M. Trocan, Á. Rodríguez-Vázquez, and R. Carmona-Galán, “Compressive Imaging Using RIP-Compliant CMOS Imager Architecture and Landweber Reconstruction,” *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 30, no. 2, pp. 387–399, Feb. 2020.
- [42] “Theory of cellular automata: A survey,” *Theoretical Computer Science*, vol. 334, no. 1-3, pp. 3–33, Apr. 2005.