

Mining DNA Sequences Based on Spatially Coded Technique Using Spatial Light Modulator

Masoome Fazelian, Sajjad AbdollahRamezani, Sima Bahrani, Ata Chizari, Mohammad Vahid Jamali, Pirazh Khorramshahi, Arvin Tashakori, Shadi Shahsavari, and Jawad A. Salehi, *Fellow, IEEE*
Optical Networks Research Laboratory (ONRL), Sharif University of Technology (SUT), Tehran, IRAN
Corresponding Author Email: jasalehi@sharif.edu

Abstract—In this paper, we present an optical computing method for string data alignment applicable to genome information analysis. By applying moiré technique to spatial encoding patterns of deoxyribonucleic acid (DNA) sequences, association information of the genome and the expressed phenotypes could more effectively be extracted. Such moiré fringes reveal occurrence of matching, deletion and insertion between DNA sequences providing useful visualized information for prediction of gene function and classification of species. Furthermore, by applying a cylindrical lens, a new technique is proposed to map two-dimensional (2D) association information to a one-dimensional (1D) column of pixels, where each pixel in the column is representative of superposition of all bright and dark pixels in the corresponding row. By such a time-consuming preprocessing, local similarities between two intended patterns can readily be found by just using a 1D array of photodetectors and post-processing could be performed on specified parts in the initial 2D pattern. We also evaluate our proposed circular encoding adapted for poor data alignment condition. Our simulation results together with experimental implementation verify the effectiveness of our dynamic proposed methods which significantly improve system parameters such as processing gain and signal to noise ratio (SNR).

Index Terms—String data alignment, moiré pattern, DNA sequencing, spatial light modulator.

I. INTRODUCTION

Emerging various widespread human diseases speeds up the growing rate of genomics. Accordingly, analysis of deoxyribonucleic acid (DNA) sequences, as a medium storing by far more important information about properties of an organism, has intrigued many researchers to extract significant knowledge about life sciences [1]–[3]. As a common event in evolution process, mutation would modify DNA data sequences comprising of a finite number of basic elements known as nucleotides, i.e., adenine (A), cytosine (C), guanine (G), and thymine (T), which are independent of each other. Since each sequence data conceals in a collection of one-dimensional (1D) strings forming a genome, the role of string data alignment or pattern matching against a sequence of genomes is much more critical for comparison and interpretation of DNA-based structures [4]. Due to rapidly evolving DNA-sequencing, investigating through highly extensive DNA databases to identify occurrence of exchange, deletion, and insertion of specific data, find target DNA strings or newly genes and classify species is becoming a costly and challenging problem for researchers [5], [6].

All recent sequencing technologies, including Roche/454, Illumina, SOLiD and Helicos, are able to produce data of the order of giga base-pairs (Gbp) per machine day [7]. However, with the emergence of such enormous quantities of data, even the fast digital electronic devices are not effective enough

to align capillary reads [8], [9]. Actually, today's electronics technology would not permit us to achieve high rate of analysis in sequence matching and information processing due to the time consuming nature of serial processing [5], [10], [11]. To keep pace with the throughput of sequencing technologies, many new alignment algorithms have been developed, but demands for faster alignment approaches still exist. As a result, the necessity of finding a novel implementation to provide high performance computational systems is undeniable [12], [13]. High data throughput, inherent parallelism, broad bandwidth and less-precise adjustment of optical computing provide highly efficient devices which can process information with high speed and low energy consumption. It is worth mentioning that visible light in optical computing systems realizes information visualization for human operators to more effectively carry out genome analysis. Employing a powerful technique to encode DNA information into an optical image besides optical computing capabilities would definitely guarantee to efficaciously perform genomes analysis [2], [12].

While recent implementations were static and relying on printed transparent sheets [14], [15], herein, we theoretically and experimentally present dynamic string data alignment based on a spatially coded moiré technique [16], [17] implemented on spatial light modulators (SLMs) which enables one to investigate useful hiding information in genomes. The remaining of the paper is organized as follows. In Section II, the principle of string data matching using the spatially coded technique is explained. In Section III, bar and circular patterns as an effective scheme for string data alignment will be discussed. In Section IV, the experimental optical architecture and obtained results will be appeared to verify practical feasibility of our proposed pattern, and Section V concludes the paper.

II. PRINCIPLES

In this section, the principles of string alignment by moiré technique are outlined. Consider two data sequences. The goal of string alignment is evaluation of similarities and differences between them. In particular, we are interested in distinguishing insertion and deletion of elements in any strings with respect to each other. Moiré technique applies high speed parallel processing of light to perform string alignment. In this approach, four components of strings, namely {A, G, C, T} are encoded as {1000, 0100, 0010, 0001}, respectively. Based on this coding, the strings are spatially coded into images where each component corresponds to four narrow stripes with one bright stripe as “1” and three dark stripes as “0” (see Fig. 1). The coded images are then overlapped with a small relative

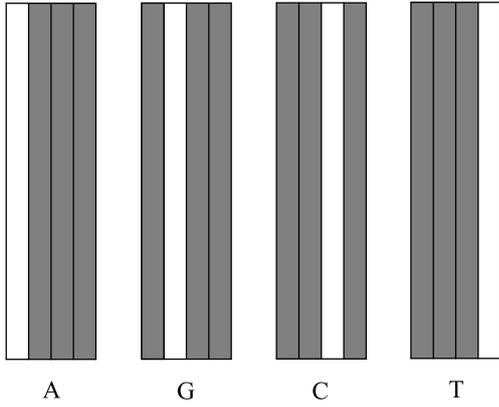


Fig. 1. Graphical patterns for DNA bases

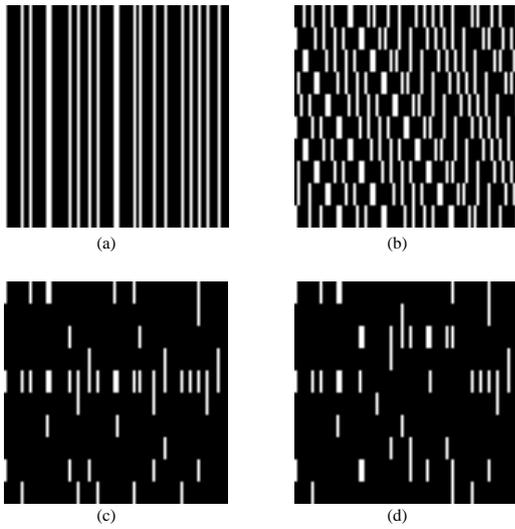


Fig. 2. Spatial code patterns for (a) S_2 , (b) subsequent shifts of initial string S_1 , (c) output pattern by overlapping (a) and (b), (d) output pattern by overlapping S_3 and (b).

angle, and by using this technique, correlating segments of the second string in various shifts of the first one can evidently be distinguished. The subsequent matched elements will be appeared as a bright line in the observed pattern of overlapped images.

As an example, consider two strings S_1 of length 40 and S_2 of length 20. Now, we want to search for S_2 in S_1 . Fig. 2(a) shows $S_2 = \{ACGTATCCGTACAGGTCGAA\}$ with respect to the codes appeared in Fig. 1, and each row in Fig. 2(b) shows subsequent shifts of initial string $S_1 = \{TCCGTACGTATCCGTACAGGTCGAATGCGTACATCGACCT\}$; for example first row shows $S_1(1:20)$, second row shows $S_1(2:21)$, up to the last row. Overlapping Fig. 2(a) and (b) results in the pattern shown in Fig. 2(c); the bright line in the fourth row illustrates that a correlation has happened for a shift of 6, i.e., S_2 and $S_1(6:25)$ are matched.

The insertion and deletion of elements lead to a vertical shift in some parts of the bright line in the overlapping pattern. Each break point indicates the location where insertion or deletion is occurred. The positive and negative vertical shifts correspond to insertion and deletion of some

TABLE I
CORRESPONDING CODES FOR POLARIZED SPATIAL PATTERNS IN FIGS. 3 AND 4.

DNA bases	A	G	C	T
Type I	1000	0100	0010	0001
Type II	$H00H$	$V0V0$	$0V0V$	$0HH0$

H : Horizontal, V : Vertical

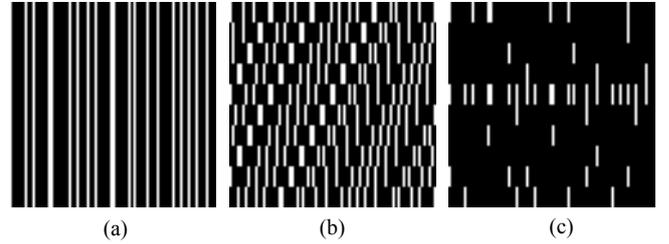


Fig. 3. Spatial code patterns of (a) S_1 , (b) S_2 , and (c) corresponding correlation for type I.

elements, respectively. As an example, consider the string $S_3 = \{ACGTATAGCCGTACAGGTCGAA\}$ generated by insertion of “AG” between the sixth and seventh element of S_2 and deletion of the fourteenth and fifteenth element of S_2 . Figure 2(d) depicts the output pattern obtained by multiplying S_3 and Fig. 2(b).

III. PROPOSED METHODS

In this section, we propose several practically feasible moiré patterns for string data alignment applications. Wave nature of light provides enough degrees of freedom, i.e., amplitude, phase, and polarization manipulation for sequence data processing.

The first coding approach is based on correlation in which a sequence is simply symbol-by-symbol compared to another sequence. In DNA sequence data processing, each symbol denotes a DNA base. In comparing two symbols with each other, similar symbols generate a bright spot; hence, a correlated set realizes a bright line. This line is fragmented in the case of insertion and deletion in which vertical distance between fragmented lines identifies the number of deleted or inserted elements in that place. In this method, SNR can easily be calculated; in the case of two independent and identical distributed sequences, the probability of such a random similarity and hence number of bright spots with respect to full matching is 0.25, leading to 6 dB SNR. It is notable that system’s SNR is proportional to the ratio of bright line intensity to the average intensity of other rows.

Another coding technique is based on concatenating two subsequent elements, for example $S(i : i + 1)$ and $S(i + 1 : i + 2)$, as a group. Subsequent groups have a common element which ensures an easier detection procedure of insertion or deletion. Coding sequences in overlapped pairs not only does increase the SNR but also makes correlated elements more distinguishable even in the cases of insertion and deletion. In this method, a 12 dB SNR can be expected in that the probability of random similarity for a word of two symbols is 0.0625.

TABLE III

CHARACTERIZATION OF THE PROPOSED METHOD AND THE CORRESPONDING SNRS OBTAINED VIA SIMULATION, WHERE N STANDS FOR THE NUMBER OF SLM PIXELS.

	Type I	Type II	Type III	Type IV
Processing Gain	$N/4$	$N/4$	$N/8$	$N/16$
SLM Modulation Capacity	Intensity or Polarization	Intensity and Polarization	Intensity and Polarization	Intensity or Polarization
SNR (dB)	6.8854	6.4648	12.2260	12.0715

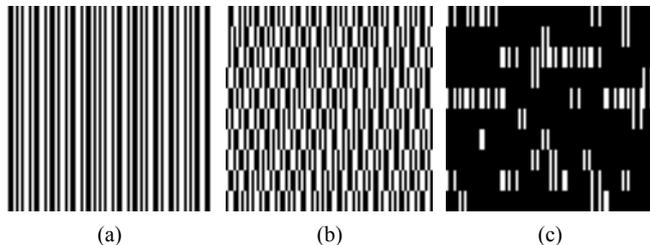


Fig. 4. Spatial code patterns of (a) S_1 , (b) S_3 , and (c) corresponding correlation for type II.

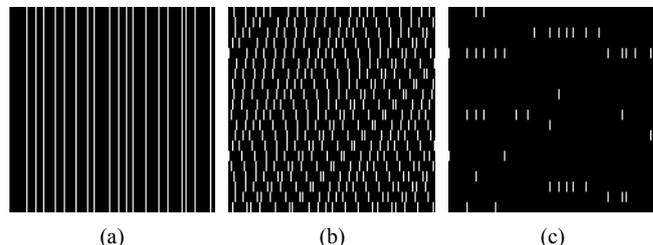


Fig. 5. Spatial code patterns of (a) S_1 , (b) S_3 , and (c) corresponding correlation for type III.

TABLE II
CORRESPONDING CODES FOR SPATIAL PATTERNS IN FIGS. 5 AND 6.

DNA bases	Type III	Type IV
AA	H0000000	1000000000000000
GA	0H000000	0100000000000000
CA	00H00000	0010000000000000
TA	000H0000	0001000000000000
AG	0000H000	0000100000000000
GG	00000H00	0000010000000000
CG	000000H0	0000001000000000
TG	0000000H	0000000100000000
AC	V0000000	0000000010000000
GC	0V000000	0000000001000000
CC	00V00000	0000000000100000
TC	000V0000	0000000000010000
AT	0000V000	0000000000001000
GT	00000V00	0000000000000100
CT	000000V0	0000000000000010
TT	0000000V	0000000000000001

A. Bar Pattern

We examined two different sets of symbols in a bar moiré pattern. While the first one employs pulse position modulation (PPM), the second comprises of a set of four orthogonal codes using both intensity and polarization (see Table I). Since there is no useful information in shifts that are not an integer product of symbol length, different rows are shifted by an integer product of four slots that form a symbol. This is by far more efficient than horizontal tilting of the second pattern and consequently compatible with finite resolution of SLM.

Simulation results are depicted in Figs. 3 and 4. By comparing the results, it is clear that using type II increases the intensity of both noise and signal but does not improve the SNR. Measuring symbol-by-symbol correlation, we see the SNR does not go further than 6 dB.

In the second approach, the codes in Table II are applied which means that for tilted pattern different rows are shifted by $8k$ in type III (word length is eight here) and $16k$ in type IV; k is a positive integer. Figs. 5 and 6 illustrate the simulation results. As it can be seen, the horizontal straight line is more vivid in types III and IV since the probability of random similarity for a word of two symbols is 0.0625; therefore, we can expect a SNR about 12 dB. In type III, we need a

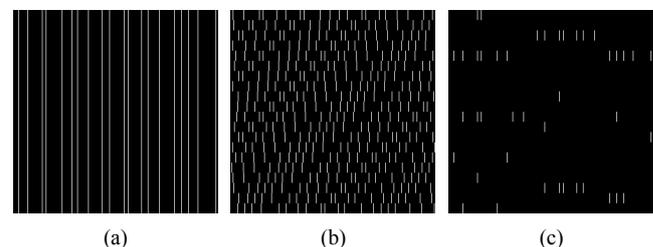


Fig. 6. Spatial code patterns of (a) S_1 , (b) S_3 , and (c) corresponding correlation for type IV.

SLM with independent intensity and polarization modulation while in type IV only intensity or polarization modulation is required. Polarization modulation can easily be converted to intensity modulation via a polarizer. Since word length for type IV is twice type III, for an equal number of SLM surface pixels, processing gain, the number of DNA bases that the setup is able to compare in each run, of type III is twice type IV. Types III and IV offer better detection capability facing insertion and deletion. It is notable that each insertion or deletion changes two words. In case of n subsequent deletion or insertion, $n+1$ words differ from initial pattern. Moreover, if we increase the word length to code the DNA bases in group of length L , $L-1$ elements are needed to be overlapped in order to detect insertion and deletion. When misalignment and other types of errors are addressed, the maximum performance of such a system could be achieved. In this case, the number of pixels of SLMs to compare two sequences of length N follows Table III. It also reports the SNR values of different types for a random sequence of length 48.

B. Circular Pattern

Optical alignment could be quite problematic in implementing bar patterns. In correlating two bar patterns, the dimension precision required should be about d/N , where d is the transverse length of a pixel and N is the total number of vertical pixels on SLM surface. On the other hand, circular moiré patterns are basically easier to be adjusted in experimental setups since only the center of circles should be aligned. Besides, it is sensitive to neither rotation nor

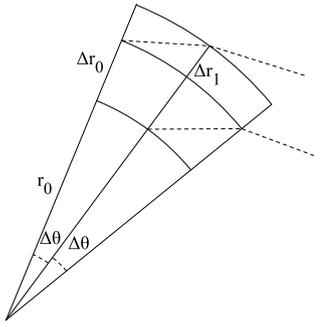


Fig. 7. Design of a sector for circular moiré pattern. Δr_1 is chosen such that $r_0 \Delta \theta \Delta r_0 = r_1 \Delta \theta \Delta r_1$.

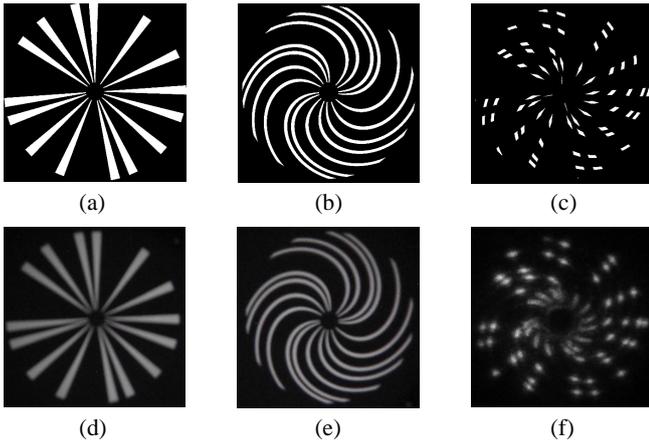


Fig. 8. Images of first, second, and output patterns obtained from simulation ((a), (b), and (c), respectively) and experiment ((d), (e), and (f), respectively) in the case of exact matching.

divergence during propagation in free-space. Despite circular pattern could only process sparse sequence data, it simplifies the optical alignment complexity when transceivers are distant.

Our proposed approach is based on encoding the strings into circular images. In this method, instead of rectangular stripes, narrow sectors are applied, as shown in Fig. 8(a). To realize the shifted versions of another string, we use curved pattern, as depicted in 8(b). Each curved sector in this pattern is designed such that the area of its segments at different radii are approximately equal. Defining r_0 and Δr_0 as our initial values for the first segment of the curved sector, we have (see Fig. 7);

$$r_0 \Delta \theta \Delta r_0 = (r_0 + \Delta r_0) \Delta \theta \Delta r_1. \quad (1)$$

From the above equation, Δr_1 is given by;

$$\Delta r_1 = \frac{r_0 \Delta r_0}{r_0 + \Delta r_0}. \quad (2)$$

Furthermore, it is obvious from Fig. 7 that $r_1 = r_0 + \Delta r_0$. Hence, we can define the following recursive relation to obtain Δr_i s and r_i s as;

$$\begin{aligned} \Delta r_i &= \frac{r_{i-1} \Delta r_{i-1}}{r_{i-1} + \Delta r_{i-1}}, \\ r_i &= r_{i-1} + \Delta r_{i-1}, \end{aligned} \quad (3)$$

respectively. Fig. 8(c) depicts the simulation and the experimental patterns after overlapping two images of Figs. 8(a) and

Equipment	Description
Reflective SLM	Holoeye, LC R-720 and LC R-2500
Biconvex lens	Thorlabs, $f = 75$ and 100 mm, $d = 40$ mm
Laser diode	1 mW, Green (532 nm), Polarized
CCD camera	Teviscom
Achromatic objective lens	Thorlabs, $NA = 0.25$
Cylindrical lens	$f = 75$ mm, $h = 50.8$ mm, $l = 53$ mm
Holder	Standa
Analyzer and polarizer	Wire grid

(b). As can be seen, a bright ring appears at the intersection of matched elements.

IV. EXPERIMENTAL SETUP AND RESULTS

In this paper, the optical architecture is implemented by two separate programmable reflective SLMs. The pixel pitch of the liquid-crystal display of each SLM is $20 \mu\text{m}$, and the pixels number is 1280×768 . The larger number of elements it comprises, the less resolution in the output plane we achieve. Further details about equipments can be found in Table III. Fig. 9 shows the architecture of multiplier realizing the proposed optical processing method for genome analysis based on spatial coded technique. The output light of a spatially coherent source such as a laser diode source or a non-coherent source like LED is collimated and then impinges the first SLM containing S_1 . A linear polarizer in front of the first SLM sets the incoming polarization state. Since the laser emits elliptical polarized light, the intensity is dependent on both the SLM and the first polarizer's states. To remove this ambiguity, the intensity of the light leaving the first polarizer has to be set to be independent of the angle.

The reflected light from the first SLM meets the second SLM which implements S_2 . Since horizontal tilt angle is so small ($< 5^\circ$) for a reflective SLM, the two SLMs have to be placed at further distance to realize an appropriate setup. Consequently, the high resolution pattern in SLM₁ would be damaged in that it convolves with free-space Green's function. Using a lens at a distance of twice the focal length ($2f$) between two SLMs makes it possible to have the exact sharp pattern of S_1 on SLM₂. Moreover, fine adjustment of the first polarizer and the analyzer ensures the maximum contrast in the plane of SLM₂. Each SLM consists of rectangular pixels where each pixel corresponds to the programmed generated binary element. A pixel with binary element "1" allows light to reflect with the same impinging polarization, ideally without any attenuation, corresponding to the white string and a pixel with binary element "0" rotates the incoming light polarization by 90° corresponding to the black string. A two-dimensional array of photodetectors can be employed to capture the output pattern; then digital processing would be done by a host computer to extract precise matching. Alternatively, analyzing the output pattern can be realized by visual inspection or using a CCD camera.

In order to verify our proposed method, we firstly show bar strings alignment between two DNA-simulated sequences; then the circular one will be demonstrated as our proposed new encoded pattern. One-dimensional strings to be aligned are illustrated in Figs. 10 and 11 in which S_1 , S_2 , and S_3 were introduced earlier. To more straightforwardly realize string

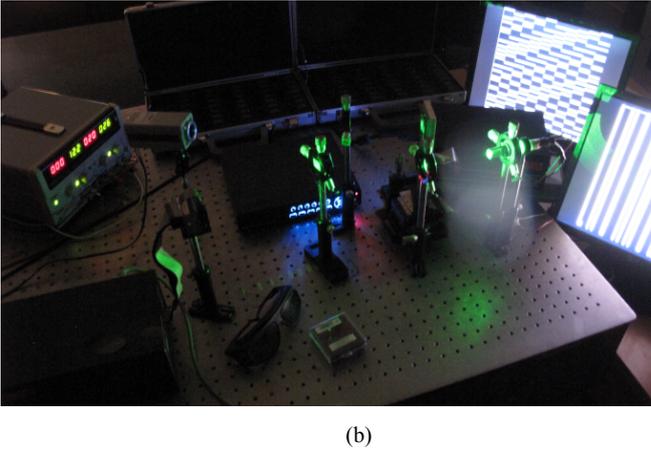
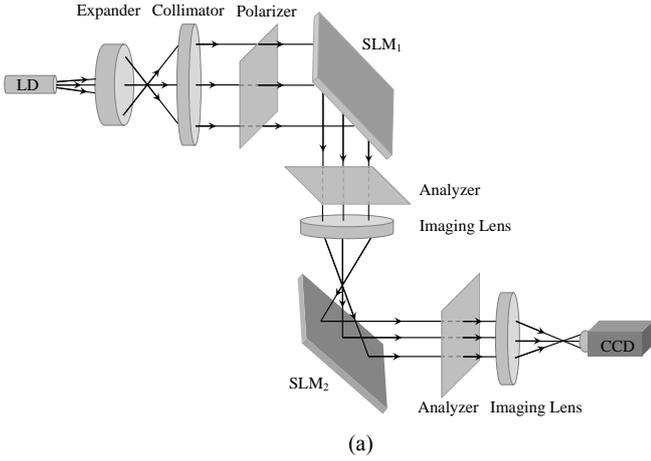


Fig. 9. (a) Schematic block diagram and (b) experimental setup for the proposed optical sequence data processing.

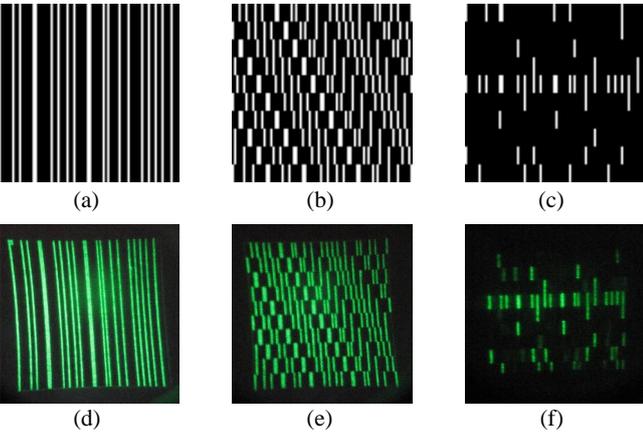


Fig. 10. Image of (a) S_1 , (b) S_2 , (c) output pattern achieved via simulation. Photograph of (d) S_1 , (e) S_2 , (f) output pattern obtained from experiment.

alignment, a cylindrical lens could be employed between the third polarizer and the output display. It is well known that such a lens transforms plane wave to an ultra-thin line. As a result, each horizontal bright line in the output pattern right behind the lens is mapped to a luminous point on the display which enables us to use a simple one-dimensional array of photodetectors to detect the occurrence of exact matching and the number of deleted or inserted elements. Figs. 12(a) and

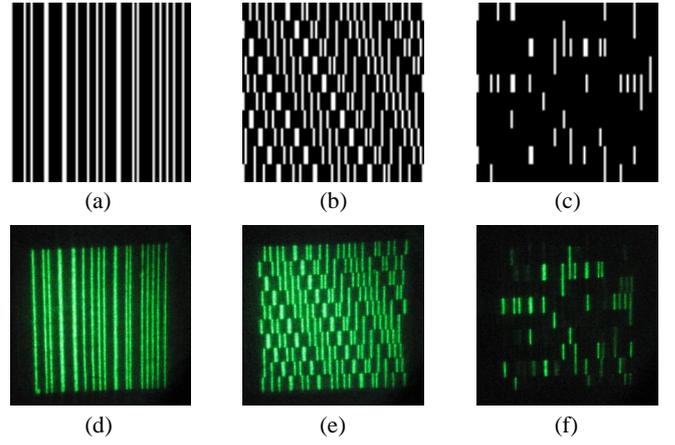


Fig. 11. Image of (a) S_1 , (b) S_3 , (c) output pattern achieved via simulation. Photograph of (d) S_1 , (e) S_3 , (f) output pattern obtained from experiment.

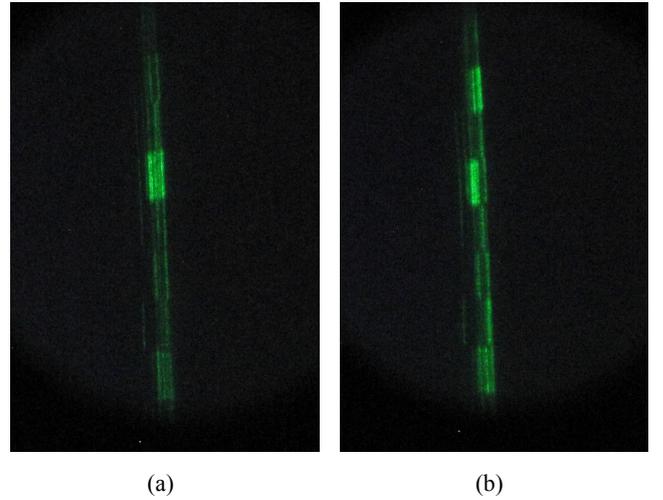


Fig. 12. Transformed output patterns of Figs. 10(f) and 11(f) on the display using cylindrical lens.

12(b) respectively illustrate the transformed versions of the output patterns in Figs. 10(f) and 11(f) at the focused plane of the cylindrical lens. Additionally, simulated and experimental results for circular patterns presented in Fig. 8 are in good agreement.

V. CONCLUSION

In conclusion, a simple and practical method based on spatially coded moiré matching technique has been proposed for string alignment processing. Easy interpretation and inherent parallelism with almost real-time processing are the main specifications of our approach which is compatible with digital devices. The processing gain and SNR of the proposed patterns, i.e., bar and circular patterns, have numerically been calculated to show the effectiveness of our method. Moreover, a preprocessing stage which remarkably decreases post-processing time needed for interpretation of output pattern has been introduced. The capability of our proposed method in DNA sequence matching has been shown via simulation. Finally, experimental results verify the performance of the method in genomics processing applications based on optical computing.

REFERENCES

- [1] J. M. Kinser, "Mining DNA data in an efficient 2d optical architecture," in *2000 International Topical Meeting on Optics in Computing (OC2000)*. International Society for Optics and Photonics, 2000, pp. 104–110.
- [2] J. Shendure and H. Ji, "Next-generation DNA sequencing," *Nature biotechnology*, vol. 26, no. 10, pp. 1135–1145, 2008.
- [3] A. C. Rajan, M. R. Rezapour, J. Yun, Y. Cho, W. J. Cho, S. K. Min, G. Lee, and K. S. Kim, "Two dimensional molecular electronics spectroscopy for molecular fingerprinting, DNA sequencing, and cancerous DNA recognition," *ACS nano*, vol. 8, no. 2, pp. 1827–1833, 2014.
- [4] J. Eid, A. Fehr, J. Gray, K. Luong, J. Lyle, G. Otto, P. Peluso, D. Rank, P. Baybayan, B. Bettman *et al.*, "Real-time DNA sequencing from single polymerase molecules," *Science*, vol. 323, no. 5910, pp. 133–138, 2009.
- [5] J. M. Rothberg, W. Hinz, T. M. Rearick, J. Schultz, W. Mileski, M. Davey, J. H. Leamon, K. Johnson, M. J. Milgrew, M. Edwards *et al.*, "An integrated semiconductor device enabling non-optical genome sequencing," *Nature*, vol. 475, no. 7356, pp. 348–352, 2011.
- [6] S. K. Min, W. Y. Kim, Y. Cho, and K. S. Kim, "Fast DNA sequencing with a graphene-based nanochannel device," *Nature nanotechnology*, vol. 6, no. 3, pp. 162–165, 2011.
- [7] M. L. Metzker, "Sequencing technologies the next generation," *Nature reviews genetics*, vol. 11, no. 1, pp. 31–46, 2010.
- [8] Z. Ning, A. J. Cox, and J. C. Mullikin, "Ssaha: a fast search method for large dna databases," *Genome research*, vol. 11, no. 10, pp. 1725–1729, 2001.
- [9] W. J. Kent, "Blat the blast-like alignment tool," *Genome research*, vol. 12, no. 4, pp. 656–664, 2002.
- [10] J. Tanida, "String data alignment by a spatial coding and moiré technique," *Optics letters*, vol. 24, no. 23, pp. 1681–1683, 1999.
- [11] J. Tanida and K. Nitta, "String data matching based on a moiré technique using 1d spatial coded patterns," in *2000 International Topical Meeting on Optics in Computing (OC2000)*. International Society for Optics and Photonics, 2000, pp. 16–23.
- [12] E. R. Mardis, "Next-generation DNA sequencing methods," *Annu. Rev. Genomics Hum. Genet.*, vol. 9, pp. 387–402, 2008.
- [13] J. L. Merklings, "Sequence matching in holographically stored genetic strings," Ph.D. dissertation, Texas Tech University, 2005.
- [14] J. Tanida, K. Nitta, and A. Yahata, "Spatially coded moire matching technique for genome information visualization," in *Photonics Asia 2002*. International Society for Optics and Photonics, 2002, pp. 26–33.
- [15] K. Niita, H. Togo, A. Yahata, and J. Tanida, "Genome information analysis using spatial coded moire technique," in *Lasers and Electro-Optics, 2001. CLEO/Pacific Rim 2001. The 4th Pacific Rim Conference on*, vol. 2. IEEE, 2001, pp. II–II.
- [16] I. Amidror, *The theory of the moiré phenomenon*. Springer, 2000, no. LSP-BOOK-2000-001.
- [17] E. Gabrielyan, "The basics of line moiré patterns and optical speedup," *arXiv preprint physics/0703098*, 2007.