

CENICS 2025

The Eighteenth International Conference on Advances in Circuits, Electronics and Micro-electronics

ISBN: 978-1-68558-308-8

October 26th - 30th, 2025 Barcelona, Spain

CENICS 2025 Editors

Timm Bostelmann, FH Wedel (University of Applied Sciences), Germany

CENICS 2025

Forward

The Eighteenth International Conference on Advances in Circuits, Electronics and Micro-electronics (CENICS 2025), held between October 26th, 2025, and October 30th, 2025, in Barcelona, Spain, continued a series of events capturing advances in special circuits, electronics, and micro-electronics in both theory and practice, from fabrication to applications using these special circuits and systems. The topics covered

fundamentals of design and implementation, techniques for deployment in various applications, and advances in signal processing.

Innovations in circuits, electronics, and microelectronics enable a wide range of applications. The conference focused on advances in high-speed signal processing, micro/nano-electronics, low-power sensor systems, and specialized electronics for wearable, implantable, and e-health devices. These technologies drive new design methods, reconfigurable systems, and integration with Internet-based platforms, with special attention to devices operating close to or within the human body.

We take the opportunity to warmly thank all the members of the CENICS 2025 technical program committee, as well as all the reviewers. The creation of such a high-quality conference program would not have been possible without their involvement. We also kindly thank all the authors who dedicated much of their time and effort to contribute to CENICS 2025. We truly believe that, thanks to all these efforts, the final conference program consisted of top-quality contributions. We also thank the members of the CENICS 2025 organizing committee for their help in handling the logistics of this event.

We hope that CENICS 2025 was a successful international forum for the exchange of ideas and results between academia and industry for the promotion of progress in the field of circuits, electronics, and micro-electronics.

CENICS 2025 Chairs

CENICS 2025 Steering Committee

Junghee Lee, Korea University, Korea
David Cordeau, XLIM | University of Poitiers, France
Timm Bostelmann, FH Wedel (University of Applied Sciences), Germany
Kenneth Skovhede, Duplicati, Inc., Denmark

CENICS 2025 Publicity Chairs

Lorena Parra Boronat, Universidad Politécnica de Madrid, Spain Laura Garcia, Universidad Politécnica de Cartagena, Spain

CENICS 2025 Committee

CENICS 2025 Steering Committee

Junghee Lee, Korea University, Korea
David Cordeau, XLIM | University of Poitiers, France
Timm Bostelmann, FH Wedel (University of Applied Sciences), Germany
Kenneth Skovhede, Duplicati, Inc., Denmark

CENICS 2025 Publicity Chairs

Lorena Parra Boronat, Universidad Politécnica de Madrid, Spain Laura Garcia, Universidad Politécnica de Cartagena, Spain

CENICS 2025 Technical Program Committee

Naeem Abbasi, Qualcomm Technologies Inc., San Diego, USA

Francesco Aggogeri, University of Brescia, Italy

Ahmed Ammar, Ohio Northern University, USA

Amjad Anvari-Moghaddam, Aalborg University, Denmark

Mohammed A. Aseeri, King Abdulaziz City of Science and Technology (KACST), Kingdom of Saudi Arabia

M. Ali Aydin, Istanbul University, Turkey

Vincent Beroulle, Grenoble INP-UGA, France

Mahajan Sagar Bhaskar, Prince Sultan University (PSU), Saudi Arabia

Timm Bostelmann, FH Wedel (University of Applied Sciences), Germany

Manuel José Cabral dos Santos Reis, IEETA / University of Trás-os-Montes e Alto Douro, Portugal

Juan-Vicente Capella-Hernández, Universitat Politècnica de València, Spain

Spandonidis Christos, Prisma Electronics SA, Greece

Tales Cleber Pimenta, Universidade Federal de Itajuba, Brazil

David Cordeau, XLIM | University of Poitiers, France

Bappaditya Dey, imec, Leuven, Belgium

Said Drid, University of Batna 2, Algeria

Francisco Falcone, UPNA-ISC, Spain

Laurent Fesquet, TIMA Laboratory | Grenoble Institute of Technology, France

Nazila Fough, Robert Gordon University, UK

Kamoun Fourati Fourati, University of Sfax, Tunisia

Patrick Girard, LIRMM - University of Montpellier 2 / CNRS, France

Victor Grimblatt, Synopsys Chile R&D Center, Chile

Wenkai Guan, Marquette University, USA

Mohammad Haider, The University of Alabama at Birmingham, USA

Amir M. Hajisadeghi, Amirkabir University of Technology (Tehran Polytechnic), Iran

Petr Hanáček, Brno University of Technology, Czech Republic

Abdus Sami Hassan, Chosun University, Korea

Wen-Jyi Hwang, National Taiwan Normal University, Taipei, Taiwan

Malinka Ivanova, Technical University of Sofia, Bulgaria

Zhenge Jia, University of Pittsburgh, USA

Mouna Baklouti Kammoun, University of Sfax, Tunisia

Andrei Karatkevich, AGH University of Science and Technology, Krakow, Poland

Faiq Khalid, Technische Universität Wien, Austria

Kasem Khalil, Western Kentucky University, USA

Ioannis Kouretas, University of Patras, Greece

Junghee Lee, Korea University, South Korea

Kevin Lee, School of InformationTechnology | Deakin University, Melbourne, Australia

Samira Legrini, Badji Mokhtar University, Algeria

Shuai (Steven) Li, Swansea University, UK

Diego Liberati, National Research Council of Italy, Italy

Yo-Sheng Lin, National Chi Nan University, Taiwan

David Lizcano, Madrid Open University (UDIMA), Spain

Jose Manuel Molina Lopez, Universidad Carlos III de Madrid, Spain

Xuyang Lu, Shanghai Jiao Tong University, China

Weizhi Meng, Lancaster University, UK

Amalia Miliou, Aristotle University of Thessaloniki, Greece

Bartolomeo Montrucchio, Politecnico di Torino, Italy

Rafael Morales Herrera, University of Castilla-La Mancha, Spain

Ioannis Moscholios, University of Peloponnese, Greece

Umair Mujtaba Qureshi, City University of Hong Kong, Hong Kong

Arnaldo Oliveira, UA-DETI/IT-Aveiro, Portugal

Youssef Ounejjar, ETS, Montreal, Canada

Maria S. Papadopoulou, Aristotle University of Thessaloniki, Greece

Michalis Pavlidis, University of Brighton, UK

Ladislav Polak, Brno University of Technology, Czech Republic

Jorge Portilla, Universidad Politécnica de Madrid, Spain

Enrique Romero-Cadaval, University of Extremadura, Spain

Pedro Santana, ISCTE - University Institute of Lisbon, Portugal

Sergei Sawitzki, FH Wedel (University of Applied Sciences), Germany

Emilio Serrano Fernández, Technical University of Madrid, Spain

Mustafa M. Shihab, The University of Texas at Dallas, USA

Kenneth Skovhede, Duplicati, Inc. Denmark

Ivo Stachiv, Harbin Institute of Technology, Shenzhen China & Institute of Physics - Czech Academy of

Sciences, Prague, Czech Republic

Kenneth Stewart, University of California, Irvine

Viera Stopjakova, Slovak University of Technology, Bratislava, Slovakia

Carlos Travieso González, University of Las Palmas de Gran Canaria, Spain

Muhammad S. Ullah, Florida Polytechnic University, USA

Prajoona Valsala, Dhofar University, Salalah, Oman

John S. Vardakas, Iquadrat, Barcelona, Spain

Miroslav Velev, Aries Design Automation, USA

Gregg Vesonder, Stevens Institute of Technology, USA

Manuela Vieira, CTS/ISEL/IPL, Portugal

Aili Wang, ZJU-UIUC Institute | Zhejiang University, China

Hao Wang, MediaTek USA Inc., USA

Pengcheng Xu, UCLouvain, Belgium

Xiangxing Yang, pSemi Corporation, USA

Yintang Yang, Xidian University, China Fei Yuan, Ryerson University, Canada Piotr Zwierzykowski, Poznan University of Technology, Poland

Copyright Information

For your reference, this is the text governing the copyright release for material published by IARIA.

The copyright release is a transfer of publication rights, which allows IARIA and its partners to drive the dissemination of the published material. This allows IARIA to give articles increased visibility via distribution, inclusion in libraries, and arrangements for submission to indexes.

I, the undersigned, declare that the article is original, and that I represent the authors of this article in the copyright release matters. If this work has been done as work-for-hire, I have obtained all necessary clearances to execute a copyright release. I hereby irrevocably transfer exclusive copyright for this material to IARIA. I give IARIA permission or reproduce the work in any media format such as, but not limited to, print, digital, or electronic. I give IARIA permission to distribute the materials without restriction to any institutions or individuals. I give IARIA permission to submit the work for inclusion in article repositories as IARIA sees fit.

I, the undersigned, declare that to the best of my knowledge, the article is does not contain libelous or otherwise unlawful contents or invading the right of privacy or infringing on a proprietary right.

Following the copyright release, any circulated version of the article must bear the copyright notice and any header and footer information that IARIA applies to the published article.

IARIA grants royalty-free permission to the authors to disseminate the work, under the above provisions, for any academic, commercial, or industrial use. IARIA grants royalty-free permission to any individuals or institutions to make the article available electronically, online, or in print.

IARIA acknowledges that rights to any algorithm, process, procedure, apparatus, or articles of manufacture remain with the authors and their employers.

I, the undersigned, understand that IARIA will not be liable, in contract, tort (including, without limitation, negligence), pre-contract or other representations (other than fraudulent misrepresentations) or otherwise in connection with the publication of my work.

Exception to the above is made for work-for-hire performed while employed by the government. In that case, copyright to the material remains with the said government. The rightful owners (authors and government entity) grant unlimited and unrestricted permission to IARIA, IARIA's contractors, and IARIA's partners to further distribute the work.

Table of Contents

1

Adaptive Architectures for Forward Error Correction

| Sergei Sawitzki | |
|--|---|
| Optical-System-Aware Feature Extraction for Lithography Hotspot Detection Masahiro Yamamoto, Masato Inagi, and Shinobu Nagayama | 8 |

Adaptive Architectures for Forward Error Correction

Sergei Sawitzki
FH Wedel (University of Applied Sciences)
Wedel, Germany

e-mail: Sergei.Sawitzki@fh-wedel.de

Abstract—This work discusses adaptive designs of decoders for several forward error correction codes. In the first scenario, decoders capable of decoding the codes with different parameters belonging to the same family are introduced. The results suggest, that the hardware overhead caused by the additional flexibility in most cases is as low as 5-10% additional silicon footprint compared to the implementation based on the fixed code parameters. In the second scenario, reconfigurable designs of multi-family decoders are discussed. Since some parts of the data-path and internal memory can be reused for different decoders, silicon area savings of more than 40 % are achievable compared to the overall chip area costs of the individual decoder implementation per code family. The overall usefulness of such reconfigurable decoder designs depends on the application case, for instance, the throughput requirements or the necessity to process data streams encoded with different codes in parallel. The figures reported herein summarize and extend some previous work known from literature, as well as the research carried out by the author.

Keywords-forward error correction; FEC decoders; adaptive decoder architectures.

I. INTRODUCTION AND RELATED WORK

Most of the modern digital communication and data storage standards enforce the use of different error correcting codes. Since in most cases the error correction process is carried out by the receiver – without resending or rereading the data packet – these codes are usually referred to as Forward Error Correcting (FEC) codes and the respective hard- or software units are called FEC encoders and decoders. The FEC codes most widely used nowadays are belonging to one of the following code families: convolutional codes, Reed-Solomon (RS) codes, Low-Density Parity-Check (LDPC) codes, and turbo codes. Some standards specify varying codes for different data rates or combine codes from different families within one run to address specific channel properties like cross-talk and fading improving the overall bit error rate.

Over the last few years, mobile devices got significantly more flexible. Modern mobile phones and other wearable devices support numerous communication interfaces providing the user with almost unlimited connectivity. Even the low-price devices nowadays support wireless LAN, different mobile internet standards, and Bluetooth, just to name a few. At least some of them rely on the same kind of FEC algorithms, although the specific parameters like code generator polynomials, coding rates, constraint length, block sizes, and so on may vary even within the same communication standard. To address this issue, adaptive Viterbi and RS decoders will be discussed below and compared to the straight-forward implementation for the single

code instance. Such comparison will provide a glimpse of the overhead introduced by the additional flexibility.

On the other side, some computations carried out during the decoding process of different code families are quite similar. In addition, some data need to be temporarily stored either throughout the complete decoding procedure or, in case of LDPC and turbo codes, between adjacent decoding iterations resulting in specific memory requirements. This poses the question, if at least some parts of the silicon dedicated to the data-path and memory can be reused among different standards reducing the silicon footprint and improving the power balance compared to the straight-forward implementation of every single decoder. This question will be addressed in the discussion of combined Viterbi/turbo and LDPC/turbo decoders.

The general applicability and advantages of the adaptive FEC decoder implementation depend on the use case, e.g., the necessity to process several data streams encoded with the same code or with different codes in parallel or the option to switch between different data rates. This should be kept in mind during the discussion of the different design options.

In principle, FEC decoding algorithms can be implemented both in software and hardware. However, given the fact, that most of them are quite computationally intensive, softwarebased solutions are usually lagging behind the high throughput requirements imposed by the majority of contemporary communication and data storage standards. For instance, iterative decoding involved in the decoding process of turbo and LDPC codes may require as many as 1500 machine instructions per bit (or even more depending on the number of decoding iterations) [1]. Running such decoders at the data rates of several hundreds of megabits per second requires computational power which is beyond the capabilities of even most powerful microprocessors. For this reason, the scope of this paper is limited to the dedicated hardware architectures, i.e., direct mapping of the FEC decoders to silicon. Alternative approaches which are less flexible than the software solution but still offer at least some properties of the architectures built around the instruction set paradigm are Application-Specific Instruction-set Processors (ASIP) [2][3] and Networks-on-Chip (NoC) [4]. These are beyond the scope of this contribution. A comprehensive overview concerning the scalability issues and multi-standard capabilities of different FEC decoders is provided in [1][5]. Based upon these studies, some of the results are re-evaluated, supported by new findings and extended herein.

The rest of this paper is organized as follows. Section II introduces an adaptive implementation of the Viterbi decoder and compares several designs known from literature. In Sec-

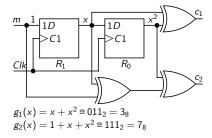


Figure 1. Simple convolutional encoder.

tion III, the additional costs of a reconfigurable architecture for a Reed-Solomon decoder are discussed. Section IV summarizes the general options of combining decoders for different code families within one design and analyzes the architectures of reconfigurable Viterbi-turbo and LDPC-turbo decoders in more detail. Finally, Section V provides a summary of this work and draws some conclusions concerning adaptive and multifamily FEC decoder implementation showing some directions for future research.

Due to the limited space, the detailed description of the coding schemes and decoder architectures is omitted. Instead, the key ideas are summarized and some references for the readers unfamiliar with the matter are provided where appropriate.

II. VITERBI DECODERS

Convolutional codes were introduced in 1955 [6] and became one the most widely applied FEC code family. The incoming data stream is stored in a shift register, from which several bits are combined using an XOR operator. The number of delay stages in the shift register including the input of the first delay stage is referred to as constraint length K. Larger K values provide better error correction capabilities (at the cost of the increasing decoder complexity). Constraint lengths between 5 and 11 are common. Particular XOR operators are selected to produce the output data streams. The result is the discrete convolution of the input data with encoder's impulse responses in the time domain, hence the name of the code [7, pp. 455–456]. The number n of the output data streams defines the code rate 1/n. For the codes used in most communication standards 2 to 4 output bits are produced for every input bit, so the code rate R lies between 1/2 and 1/4.

The particular codes differ depending on which bits from the sequence stored in the shift register are involved in the XOR calculation. This information is provided in the form of the so-called generator polynomials g, one for every output. Figure 1 shows an example of a simple convolutional encoder with K=3, R=1/2, $g_1=011$ and $g_2=111$. As can be seen, the generator polynomials are often specified by the binary or octal representation of their coefficients. The encoder is a Finite-State Machine (FSM) whose output depends on g.

The decoding process is based on the maximum-likelihood principle. Decoder's task is to find the state transition sequence, which encoder FSM most likely underwent based on its output sequence, which in most cases is corrupted by noise during the transmission. Afterwards, the state transitions can be mapped to the corresponding input data sequence. The most widely used decoding algorithm for convolutional codes is based on dynamic programming and was introduces in 1967 by Andrew Viterbi [8]. In the first stage of the decoding process, every received symbol is compared to all legal output symbols of the encoder. The result is called Branch Metric (BM). In the simplest case of the so-called hard-input decoding, BM is the Hamming distance. Due to modulation and analog-digital conversion the input symbols are usually integer numbers represented by several bits. This case is referred to as soft-input decoding, where BM is represented by other measures like squared euclidean distance. The number of BM values to be calculated depends on the code rate since every additional output bit increases the number of the legal output symbols by the factor of two and every received symbol needs to be compared with all of them.

In the second stage, branch metrics are used to calculate the Path Metrics (PM). One PM needs to be calculated for every state of the encoder. All path metrics are set to zero in the beginning and updated with every processed symbol. The corresponding operation is called Add-Compare-Select (ACS) and is the computational core of the decoding process. Every ACS unit consists of two adders, one comparator and one multiplexer. In addition, one register is required per unit to store the actual PM value. The number of ACS units is equal to 2^{K-1} , i.e., it increases exponentially with the constraint length of the code. The result of the ACS calculation is the new path metric and the binary decision, which state transition the encoder underwent while the corresponding output symbol was generated. This decision is stored for every processed symbol and every ACS unit in form of a single bit called decision bit. The decision bits represent several paths through the possible state transition sequences of the encoder with every path corresponding to the certain input sequence. It can be shown, that under normal conditions all paths merge to a single one after a certain number of steps, which ist called Trace-Back Depth (TBD). This path represents the most probable sequence of the encoder's state transitions. As a consequence, decoder needs to store the decision bits for at least the number of symbols equal to TBD. Given the fact, that one decision bit is produced by every ACS unit, the memory capacity for a single trace-back run equals to TBD $\times 2^{K-1}$. As the decision bits cannot be overwritten during the trace-back, this amount needs to be at least doubled. The exact value of the TBD depends on the channel conditions. For AGWN channels TBD values of $5 \times K$ are sufficient, but they may be higher than $20 \times K$ for channels with fading and multi-path propagation.

Trace-back is the final stage of the decoding process, in which the decision bits corresponding to the most probable path are traversed and the corresponding input bits are stored in a Last-In First-Out (LIFO) buffer. Finally, the LIFO memory is read to produce the decoded input sequence (which in optimal case will be the same as the original input to the encoder).

A Viterbi decoder for a specific convolutional code has fixed design parameters like the number of the BM and the ACS units, TBD, and so on. It is quite different in the case of the adaptive decoder. Figure 2 shows the ACS calculation

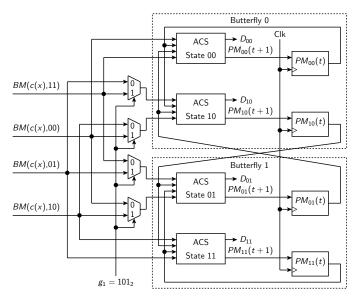


Figure 2. Path metric calculation for a reconfigurable Viterbi decoder.

unit for a decoder which can decode the code produced by the encoder in Figure 1, as well as the code generated by $g_1 = 101$ and $g_2 = 111$. While for each particular code the interconnect between the BM calculation outputs and PM calculation inputs is fixed, the adaptive decoder needs to switch between two different patterns depending on the code to be decoded. This functionality is realized by the multiplexers placed at the inputs of the particular ACS units. If the decoder has to process codes with different constraint lengths, additional flexibility is required in the interconnect pattern between the registers storing the actual PM values and the inputs of the ACS units as well. This means more and larger multiplexers at the inputs reducing the clock frequency and the throughput. A detailed design of the interconnect network for an adaptable Viterbi Decoder including the particular interconnect patterns for different settings is described in [9].

In general, varying the constraint length and/or the generator polynomials has the strongest impact on the additional chip area and the critical path of the decoder. For the trace-back stage, the memory architecture has to be calculated for the worst case, i.e., the largest K and TBD values. In addition, the addressing scheme has to be adapted to every single case, which imposes the need for programmable address counters. Compared to other modifications this overhead is almost neglectable. Figure 3 shows the trace-back unit of the reconfigurable Viterbi decoder for variable values of the constraint length. The building blocks affected by the reconfigurable nature of the design are shaded. In principle, the address generators could be dimensioned for the worst case, just like the memories. However, larger TBD increases the latency of the decoder, which would be an unnecessary overhead for smaller K values. At the same time, making the address generators pre-loadable and scalable adds almost nothing to the overall silicon footprint of the decoder. The trace forward block depicted in Figure 3 is used to calculate the initial state of the new trace-back cycle and is

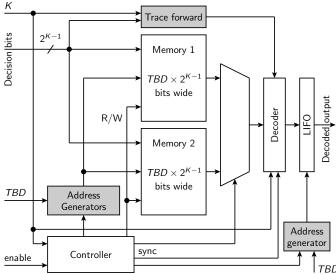


Figure 3. Trace back unit for a reconfigurable Viterbi decoder with clock and reset signals not shown for the sake of clarity.

TABLE I. COMPARISON OF ADAPTABLE AND SINGLE-CODE VITERBI DECODERS (*b* REFERS TO THE INPUT BIT-WIDTH).

| Architecture | Metric | Reference | Result |
|--|-------------------------|-----------------------------------|-------------------------|
| [11] $K = 3 \dots 7$, q variable, b fixed | silicon area | K = 7, g fixed | +2,9% |
| g variable, o ince | throughput | K = 7, g fixed | -1,5% |
| VITURBO [12] $K = 3 \dots 9$, | silicon area | K = 9, g variable | +9 % |
| g variable, b fixed | max. clock frequency | K = 9, g variable | -30 % |
| [13] $K \in \{7, 9\}$, g according to | silicon area | K = 9, g fixed (CDMA2000) | overhead is neglectable |
| EDGE, CMDA2000 and WCDMA, b fixed | critical path | K = 9, g fixed (CDMA2000) | +4% |
| [14] $K = 3 \dots 11$, g variable, | silicon area | K = 11, n = 4, $b = 5, g fixed$ | +50,1% |
| $n = 2 \dots 4, \\ b = 1 \dots 5$ | throughput | K = 11, n = 4, b = 5, g fixed | -26, 1% |

described in detail in [10]. One memory block is used to store the actual decision bits (write access) while the other is used to trace back the decision bits from the previous cycle (read access). The memories switch their roles after every cycle, which explains the need for the multiplexer at their read ports.

Table I summarizes the figures of different adaptable Viterbi decoder designs. The major findings can be generalized as follows:

• The overhead increases with the degree of flexibility. The Viterbi decoders designed to support a (small) fixed amount of codes introduce the least – in some cases even neglectable – overhead. In contrast, the full flexibility regarding the constraint length, the code rate as well as the generator polynomials requires up to 50 % more silicon area and reduces the throughput by about 26 %. Although looking diminishing at the first glance, it may still be quite a good price to pay considering the alternative to implement a dedicated decoder for every single code.

• Varying the constraint length K has the highest impact on the area overhead of the adaptive Viterbi decoder implementation. Two factors can be identified as being responsible for this matter. On one hand, the number of ACS units increases exponentially with K, which means, the number and the size of the multiplexers required to implement different interconnect patterns increases accordingly. On the other hand, the TBD value grows linearly with K, so the memory depth of the trace-back grows accordingly. At the same time, the memory width grows exponentially with K. These consequences become crucial for higher values of K: the area overhead for K = 9 is about 9 %, while for K = 11 it is more than 50 %. For Viterbi decoders with smaller constraint lengths the area overhead is almost neglectable.

III. REED-SOLOMON DECODERS

The basic idea of the error correcting linear cyclic codes nowadays knows as Reed-Solomon codes can be traced back to 1952 [15][16]. In its current form the codes were introduced in the seminal paper by Irving S. Reed and Gustave Solomon, hence the name [17]. RS codes are defined over a finite field $GF(q^m)$, which is an extension field of GF(q), where q is a prime number and m is a natural number different from zero. For all practically applied codes q = 2 is chosen, which means, that the number of the elements in the field is a power of two. Binary representation of the elements of $GF(2^m)$ requires mbits per element. This is the reason, why almost all RS codes used in communication and storage systems are defined over $\mathrm{GF}(2^8)$: each symbol can be encoded by exactly one byte. In contrast to the most other code families, the error correcting properties of the RS codes are defined symbol-wise, i.e., it does not matter, if a single bit or any other number of bits within one symbol are corrupted – the decoder treats it as a single symbol error. This explains, why RS codes are very powerful at correcting burst errors. A (255,239) RS code has the block size of 255 bytes, with 239 bytes of original message data and 2t = 255 - 239 = 16 bytes of the checksum added during encoding. t = 16/2 = 8 is the number of corrupted symbols which can be corrected by the code. If 8 adjacent symbols (64 adjacent bits in the worst case) would be corrupted during the transmission, the decoder would still be able to correct all of them. The general notation used to specify RS codes is (n, k), where n is the block size and k is the original message size before encoding (both specified in symbols of m bits width). As stated in the example above n - k = 2t, where t is the number of correctable symbol errors and R = k/n is the code rate.

The t parameter can be chosen depending on the required error correction capacity of the code. Once fixed, it defines the generator polynomial for the code, which has the degree 2t:

$$\mathbf{g}(x) = g_0 + g_1 x + g_2 x^2 + \ldots + g_{2t-1} x^{2t-1} + x^{2t}. \tag{1}$$

The coefficients g_0, \ldots, g_{2t-1} are all elements of $GF(2^m)$ and the polynomial itself is defined in such a way, that its 2t roots correspond to the 2t consecutive powers of the single element

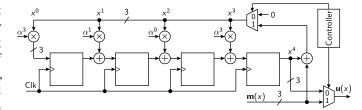


Figure 4. Reed-Solomon encoder for a (7,3) code over $GF(2^3)$.

of $GF(2^m)$. For practical reasons, usually the primitive element of the extension field $GF(q^m)$ is chosen as the first root (being the generator of the multiplicative group of this field).

The encoding procedure appends the rest $\mathbf{p}(x)$ of the division of the polynomial representing the original message $\mathbf{m}(x)$ shifted to the left by n-k symbols by the generator polynomial:

$$x^{n-k}\mathbf{m}(x) = \mathbf{q}(x)\mathbf{g}(x) + \mathbf{p}(x), \tag{2}$$

which is equivalent to

$$\mathbf{p}(x) = x^{n-k}\mathbf{m}(x) \bmod \mathbf{g}(x). \tag{3}$$

This operation can be easily carried out using feedback shift-register as illustrated in Figure 4 for a very simple (7,3) RS code. Note, that addition and multiplication operations are defined on the elements of the corresponding finite field (i.e., not in terms of common arithmetic).

Given the fact, that the encoding process can be represented by means of polynomial multiplication over $GF(q^m)$, it becomes clear, how the transmission errors can be recognized (and corrected). Since multiplication operation preserves the roots of the generator polynomial, the resulting encoded message can be inspected by the decoder by checking if 2t consecutive powers of the primitive element of $GF(q^m)$ are the roots of the polynomial corresponding to the received message. This is the first step of the decoding process which is called syndrome calculation. If at least one of the syndromes is not equal to zero, then the message was corrupted.

The syndrome values are used in the next step of the decoding process to calculate the number of the corrupted symbols. The most intuitive way to accomplish this is - in the first step - to assume the number of errors to be one and try to solve the corresponding system of equations based on the co-called error location polynomial. The solution indicates the position of the corrupted symbol. In case the system of equations for a single error is not solvable, the assumption of two corrupted symbols has to be made resulting in a different system of equations. The process is repeated until a solvable system of equations is found. This procedure is called Peterson-Zierler-Gorenstein algorithm [18][19]. For the hardware implementation the Berlekamp-Massey algorithm [20] is better suited, since it can be better parallelized and modified to reuse the same hardware to calculate both error location and error value polynomials [21]. Its detailed description, however, would go far beyond the scope of this paper. Independent of the choice of the algorithm, finite field arithmetic is heavily involved in the search for the error positions and values.

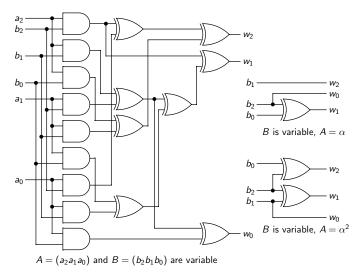


Figure 5. Variable-variable versus constant-variable multiplier for $GF(2^3)$.

As just stated, since RS codes work on symbols of m bits width, knowing the symbol error position is not enough. An additional step is required to determine the error values. This can be done by using the syndromes computed in the first step, since every set of syndromes can be mapped to certain errors. Alternatively, the so-called error value polynomial can be calculated [22]. In the final decoding step, the error values are added to the received symbols at the calculated error positions restoring the original message. To accomplish this, the received (corrupted) data is stored in the First-In First-Out (FIFO) memory making sure the right values appear at the output of the decoder at the right time. The finite filed addition corresponds to a bitwise XOR operation, which scales well and is trivial from the implementation point of view.

As mentioned above, decoding the RS codes requires numerous calculations based on the finite field arithmetic. In an RS decoder designed for a specific code, many of the GF-multipliers have one constant input (variable-constant multipliers). For adaptive decoding of different codes, a variable-variable multiplier is required. The corresponding overhead is quite high as Figure 5 illustrates. However, since only multipliers involved in the syndrome and error location computation are affected, the impact on the area of the whole RS decoder is limited.

Table II summarizes the area and throughput figures of different RS decoder designs. For the adaptive decoding, n, t, and m parameters can be changed, although changing m does not make a lot of sense, since almost all real-world RS codes are based on m=8 (ITU-T J.83, IEEE802.3bp, and IEEE802.3bj standards being an exception by specifying codes based on $\mathrm{GF}(2^7)$, $\mathrm{GF}(2^9)$, and $\mathrm{GF}(2^{10})$ correspondingly). In addition, different polynomials can be chosen as the finite field generators. Varying the GF generator polynomial $\mathbf{f}(x)$ results in a different encoding of the particular elements of the finite field. However, it has been shown, that all finite fields for the fixed values of q and m are isomorphic.

TABLE II. COMPARISON OF ADAPTABLE AND SINGLE-CODE REED-SOLOMON DECODERS.

| Comparison | Decoder architecture | | | | | | |
|--------------------|----------------------|-----------|---|----------------------------|----------------------------|--|--|
| metric | [24] | [25] [26] | | [23] | [23] | | |
| | | | | 1st variant | 2nd variant | | |
| Code | fixed | n and t | | Universal deco | | | |
| parameter | | variable | $n, t, m \text{ and } \mathbf{f}(x) \text{ variable}$ | | | | |
| \overline{m} | 8 | 8 | 18 | 110 | 18 | | |
| t | 8 | 18 | 18 | 18 | 116 | | |
| Erasure correction | no | no | no | ≤ 16 | ≤ 16 | | |
| Throughput in Gb/s | 1,6 | 0,8 | 0,048 | 2,2 | 2,4 | | |
| Gate equivalents | 21 000 | 34 000 | 44 000 | 75 000 + 35 Kbit RAM | 39 000 + 15 Kbit RAM | | |

The conclusions from the study of the different RS decoder designs are as follows:

- The area overhead of the adaptive RS decoder implementation can be quite significant. Compared to the fixed-code decoder with m=t=8, a fully reconfigurable design can consume about 85% more silicon area in terms of gate equivalents in addition to 15 KBit more SRAM. Just like in the case of adaptive Viterbi decoder this is still quite a low price to pay, since the error correction capabilities are much better and the additional erasure correcting feature is useful in many application scenarios (erasures are errors whose position is known in advance).
- Throughput is usually not an issue with RS decoding, since the symbol-per-second decoding rate is multiplied by m to obtain the bit-per-second values. For adaptive decoding it means, that the reduced clock rates caused by the additional overhead of the adaptive implementation are not significant in most cases.

IV. MULTI-FAMILY DECODERS

As can be seen from previous sections, the arithmetic involved in decoding of convolutional and RS codes is quite different which limits the possibilities of senseful hardware reuse. However, it may be an option to combine decoders for the other code families within one design with Viterbi decoder or with each other.

Turbo codes were introduced in 1993 [27]. The basic idea is to use several (usually two) recursive convolutional encoders to process the same input data. First encoder receives the data directly, the second one gets an interleaved data stream. The size of the interleaver is usually in the order of several kilobits. This approach makes the same data appearing as two statistically independent messages. The decoding is done in an iterative manner with two decoders which process the received data and pass the output to each other (with interleaving and deinterleaving steps inbetween). The decoders are based on the Soft-Input Soft-Output (SISO) principle, i.e., their output is a metric providing the likelihood for a bit to be 0 or 1 instead of the final bit decision (soft-input was already discussed in the Section II). The idea is, that a SISO decoder uses the output of its predecessor as a-priori information to improve its own

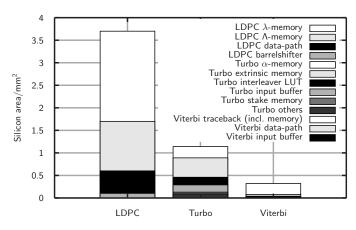


Figure 6. Chip area distribution for different FEC decoder designs (adapted from [1]).

decisions, which are then passed back to the first decoder and so on. The certainty of the calculated metrics improves over time, so that the decoding process usually can be stopped after several iterations.

LDPC codes were described in detail in the dissertation of Robert Gallager in 1963 [28]. Due to the high computational requirements of the decoding process, first practical applications came more that 30 years later as LDPC were proposed as channel codes for several communication standards. In a nutshell, LDPC codes are linear block codes with very large, very sparsely filled generator matrix. As in the case of turbo codes, the decoding process is iterative. In one iteration, parities are checked to determine, which bits of the information block are responsible for most unsatisfied parity equations. These bits are inverted and the next iteration starts. The decoding ends, as soon as all equations are satisfied or a certain number of iterations is reached. As in the case of turbo codes, most applications use soft metrics during the decoding process, which are converted to the hard bit decisions at the end.

According to this brief description of the iterative decoding, it should be clear that a lot of local memory is required to store the information bits and intermediate results. Figure 6 shows the silicon area distribution between different parts of the design for LDPC, turbo, and Viterbi decoders. The extreme disproportion in the required silicon area between LDPC and Viterbi decoders suggests, that a combination of both within the same design with the goal to reuse some of the resources does not make sense. Even if the complete Viterbi decoder area could be reused for LDPC (which is virtually impossible), the overall gain would sum up to less than 10% of the silicon footprint.

At the same time, the SISO module required for turbo decoding can be implemented based on the Soft-Output Viterbi-Algorithm (SOVA), which can also be used to decode the convolutional codes. This increases the potential for hardware reuse at the cost of slightly reduced bit error ratio for turbo decoder compared to the SISO module based on the Max-log-MAP algorithm [29].

A closer look at the internal data-path of the turbo and LDPC

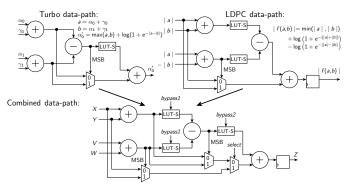


Figure 7. Combining Turbo and LDPC data-paths within a single FFU (adapted from [30], LUT-S are Look-Up-Tables storing the approximated logarithmic values).

TABLE III. COMPARISON OF MULTI-FAMILY FEC DECODERS.

| Architecture | Metric | Reference | Results |
|------------------|--------------|---------------------|------------------|
| | Viterbi | -Turbo Decoders | |
| VITURBO | silicon | Viterbi decoder | +5 % logic |
| [12] | area | $(K=3\dots 9)$ | +25 % memory |
| | silicon | separate turbo and | -14,5% |
| [29] | area | Viterbi decoder | |
| (fully parallel) | throughput | separate turbo and | no change |
| | 0 1 | Viterbi decoder | |
| [29] | silicon | separate turbo and | -28,1% |
| (time- | area | Viterbi decoder | |
| multiplex) | throughput | separate turbo and | same for Viterbi |
| manipiex) | unougnput | Viterbi decoder | 1/2 for turbo |
| [31] | silicon area | Turbo decoder | +20 % |
| | Turbo- | LDPC Decoders | |
| [1] | silicon | separate turbo and | -10% |
| area | | LDPC decoder | 1070 |
| | | turbo decoder | +10+20% |
| [30] | silicon | LDPC decoder | +15+20% |
| [20] | area | reused vs. separate | -3942% |
| | | data-paths | |

decoders reveals some similarities as well. Figure 7 shows, how the calculations involved in both decoding algorithms can be combined within the same Flexible Functional Unit (FFU) [30]. As in the case of multi-standard Viterbi decoder, some additional multiplexers are required to switch between the codes, which slightly increases the critical path.

Accordingly, the only options for multi-family decoders are Viterbi-Turbo or Turbo-LDPC designs. Some case studies of such designs are known from literature. Table III summarizes the results. As expected, some portions of the data-path and memory can be reused resulting in silicon area savings of 10–42 % while throughput is not affected in most cases of combined Viterbi-turbo decoder. If adaptable designs are compared with a single family decoder, the overhead of introducing an additional code family is very low. For instance, adding the LDPC functionality to the existing turbo design comes at only 10–20 % additional silicon area [30].

V. CONCLUSION AND FUTURE WORK

This paper discussed several options of the adaptive FEC decoder design. The findings based on literature study and own research suggest that overhead of extending a decoder

towards other codes of the same family comes at moderate cost. As expected, this overhead increases with the amount of flexibility and code complexity. At the same time, decoders covering several code families can take advantage of the resource sharing reducing the overall silicon footprint compared to the dedicated implementation of one particular decoder per family. One aspect not discussed herein is the impact of the adaptive decoder implementation on the power consumption. On one hand, reconfigurability comes at the cost of additional area and thus should result in increased energy-per-bit figures. On the other hand, it can be expected that reusing parts of the data-path and memory for different decoding algorithms leads to overall power savings. Supporting these assumptions by concrete numbers is a promising direction for further research.

REFERENCES

- [1] J. T. M. H. Dielissen, N. Engin, S. Sawitzki, and C. H. van Berkel, "FEC Decoders for Future Wireless Devices: Scalability Issues and Multi-Standard Capabilities", in *Circuits and Systems for Future Generations of Wireless Communications*, A. Tasić, W. A. Serdijn, L. E. Larson, and G. Setti, Eds., ser. Series on Integrated Circuits and Systems, Berlin: Springer, 2009, pp. 271–297.
- [2] T. Vogt and N. Wehn, "A Reconfigurable ASIP for Convolutional and Turbo Decoding in an SDR Environment", *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, vol. 16, no. 10, pp. 1309–1320, Oct. 2008.
- [3] P. M. Velayuthan, "Towards Optimized Flexible Multi-ASIP Architectures for LDPC/Turbo Decoding", PhD thesis, Université de Bretagne-Sud, Dec. 2012.
- [4] M. Scarpellino, A. Singh, E. Boutillon, and G. Masera, "Reconfigurable Architecture for LDPC and Turbo Decoding: A NoC Case Study", in *Proceedings of IEEE 10th International Symposium on Spread Spectrum Techniques and Applications*, Bologna, Italy: IEEE, Aug. 2008, pp. 671–676.
- [5] J. T. M. H. Dielissen, N. Engin, S. Sawitzki, and C. H. van Berkel, "Multi-Standard FEC Decoders for Wireless Devices", *IEEE Transactions on Circuits and Systems II*, vol. 55, no. 3, pp. 284–288, May 2008.
- [6] P. Elias, "Coding for Noisy Channels", *IRE Convention Record*, vol. 3, no. 4, pp. 37–46, 1955.
- [7] T. K. Moon, Error Correction Coding. Mathematical Methods and Algorithms. John Wiley & Sons, 2005.
- [8] A. J. Viterbi, "Error Bounds for Convolutional Codes and an Asymptotically Optimum Decoding Algorithm", *IEEE Transactions on Information Theory*, vol. 13, no. 2, pp. 260–269, Apr. 1967.
- [9] W. Tang, N. Engin, F. A. Steenhof, M. Klaassen, A. Hekstra, and S. Sawitzki, *Multi-Standard Viterbi Processor*, US Patent 8,904,266 B2, NXP B. V., Eindhoven, The Netherlands, Dec. 2014.
- [10] P. J. Black and T. H. Meng, "Hybrid Survivor Path Architectures for Viterbi Decoders", in *Proceedings of International Conference on Acoustics, Speech, and Signal Processing*, Minneapolis, MN, USA, Apr. 1993, pp. 433–436.
- [11] K. Chadha and J. R. Cavallaro, "A Reconfigurable Viterbi Decoder Architecture", in *Conference Record of the Thirty-Fifth Asilomar Conference on Signals, Systems & Computers*, M. B. Matthews, Ed., vol. 1, Pacific Grove, CA, USA: IEEE, Nov. 2001, pp. 66–71.
- [12] J. R. Cavallaro and M. Vaya, "VITURBO: A Reconfigurable Architecture for Viterbi and Turbo Decoding", in *Proceedings* of *International Conference on Acoustics, Speech, and Signal* Processing, Hong Kong: IEEE, Apr. 2003, pp. 497–500.

- [13] T. Vogt, N. Wehn, and P. Alves, "A Multi-Standard Channel-Decoder for Base-Station Applications", in 17th Symposium on Integrated Circuits and Systems Design (SBCCI 2004), Porto de Galinhas, Brazil: IEEE, Sep. 2004, pp. 192–197.
- [14] S. Sawitzki, "Embedded Reconfigurable Computing: Architekturen, Anwendungen und Entwurfswerkzeuge (Architectures, Applications, Tools)", Habilitation Thesis, TU Dresden, 2024.
- [15] K. A. Bush, "Orthogonal Arrays of Index Unity", The Annals of Mathematical Statistics, vol. 23, no. 3, pp. 426–434, Sep. 1952.
- [16] W. C. Huffman and V. Pless, Fundamentals of Error-Correcting Codes. Cambridge University Press, 2003.
- [17] I. S. Reed and G. Solomon, "Polynomial Codes over Certain Finite Fields", *Journal of the Society for Industrial and Applied Mathematics*, vol. 8, no. 2, pp. 300–304, Jun. 1960.
- [18] W. Peterson, "Encoding and Error-Correction Procedures for the Bose-Chaudhuri Codes", *IRE Transactions on Information Theory*, vol. 6, no. 4, pp. 459–470, Sep. 1960.
- [19] D. E. Gorenstein and N. Zierler, "A Class of Error Correcting Codes in p^m Symbols", *Journal of the Society of Industrial and Applied Mathematics*, vol. 9, no. 2, pp. 207–214, Jun. 1961.
- [20] E. R. Berlekamp, Algebraic Coding Theory, revised. Aegean Park Press, Jun. 1984.
- [21] H.-J. Kang and I.-C. Park, "A High-Speed and Low-Latency Reed-Solomon Decoder based on a Dual-Line Structure", in IEEE International Conference on Acoustics, Speech, and Signal Processing, Orlando, FL, USA, May 2002, pp. 3180–3183.
- [22] G. D. Forney, "On Decoding BCH Codes", IEEE Transactions on Information Theory, vol. 11, no. 4, pp. 549–557, Oct. 1965.
- [23] F.-K. Chang, C.-C. Lin, H.-C. Chang, and C.-Y. Lee, "Universal Architectures for Reed-Solomon Error-and-Erasure Decoder", in *Asian Solid-State Circuits Conference*, Hsinchu, Taiwan: IEEE, Nov. 2005, pp. 229–232.
- [24] A. G. M. Strollo, N. Petra, D. De Caro, and E. Napoli, "An Area-Efficient High-Speed Reed-Solomon Decoder in 0.25 μm CMOS", in *Proceedings of the IEEE 30th European Solid* State Circuits Conference, Leuven, Belgium: IEEE, Sep. 2004, pp. 479–482.
- [25] H.-Y. Hsu and A.-Y. Wu, "VLSI Design of a Reconfigurable Multi-Mode Reed-Solomon Codec for High-Speed Communication Systems", in *IEEE Asia-Pacific Conference on ASIC Proceedings*, Taipei, Taiwan: IEEE, Aug. 2002, pp. 359–362.
- [26] J.-C. Huang, C.-M. Wu, M.-D. Shieh, and C.-H. Wu, "An Area-Efficient Versatile Reed-Solomon Decoder for ADSL", in *Proceedings of the 1999 IEEE International Symposium on Circuit and Systems (ISCAS'99)*, vol. 1, Orlando, FL, USA: IEEE, May 1999, pp. 517–520.
- [27] C. Berrou, A. Glavieux, and P. Thitimajshima, "Near Shannon Limit Error-Correcting Coding and Decoding: Turbo Codes", in Proceedings of the IEEE International Conference on Communication, Geneva, Switzerland: IEEE, May 1993, pp. 1064–1070.
- [28] R. G. Gallager, "Low-Density Parity-Check Codes", Sc. D. thesis, Massachusetts Institute of Technology, Cambridge, MA, USA, 1963.
- [29] G. Krishnaiah, N. Engin, and S. Sawitzki, "Scalable Reconfigurable Channel Decoder Architecture for Future Wireless Handsets", in *Design, Automation and Test in Europe Conference and Exhibition, DATE 2007*, Nice, France: European Design and Automation Association, Apr. 2007, pp. 1563–1568.
- [30] Y. Sun and J. R. Cavallaro, "A Flexible LDPC/Turbo Decoder Architecture", *Journal of Signal Processing Systems*, vol. 64, no. 1, pp. 1–16, Jul. 2011.
- [31] I. Ahmed and T. Arslan, "A Reconfigurable Viterbi Decoder for a Communication Platform", in *Proceedings of the 2005 International Conference on Field Programmable Logic and Applications (FPL)*, T. Rissa, S. Wilton, and P. Leong, Eds., Tampere, Finland: IEEE, Aug. 2005, pp. 435–440.

Optical-System-Aware Feature Extraction for Lithography Hotspot Detection

Masahiro Yamamoto, Masato Inagi, Shinobu Nagayama Graduate School of Information Sciences Hiroshima City University Hiroshima, Japan

e-mail: m_yamamoto@lcs.info.hiroshima-cu.ac.jp, {inagi|s_naga}@hiroshima-cu.ac.jp

Abstract—This paper proposes a new feature vector for machine learning-based hotspot detection in lithography for Large-Scale Integration (LSI) fabrication, which incorporates the optical characteristics of exposure systems. Unlike existing features that focus only on local layout sub-patterns, the proposed feature takes into account optical behavior essential to accurate pattern transfer. In LSI fabrication, a hotspot is a region in the layout where an undesired open or short circuit may occur, even if the design rules are satisfied. Hotspots can significantly reduce manufacturing yield, and the cost of reworking after fabrication begins is substantial. Therefore, it is crucial to detect and remove hotspots at the pre-fabrication stage. Although several feature vectors have been developed for hotspot detection, most of them ignore the optical system's influence, which is critical in the lithography process. By incorporating optical characteristics, our proposed feature aims to improve detection accuracy and reduce the need for time-consuming lithography simulations.

Keywords-lithography; hotspot; feature vector; optical system.

I. Introduction

In the lithography process, which is one of the key steps in semiconductor manufacturing, laser light from the exposure system is projected onto a photomask, which serves as the master template of circuit patterns, and the pattern is transferred onto a semiconductor wafer coated with a photosensitive material. In this process, due to optical diffraction, some areas may fail to be correctly transferred even if they comply with the design rules. Such regions are referred to as hotspots. Since photomask fabrication is highly expensive, it is necessary to detect these hotspots prior to manufacturing and revise the layout patterns accordingly.

Lithography simulation, which computes phenomena, such as light diffraction and the behavior of the photosensitive material on the wafer, is a common method used to detect hotspots before mask or product fabrication. However, applying this simulation across the entire layout is extremely time-consuming. If hotspots can be rapidly detected through methods other than simulation, allowing prompt initiation of pattern revision, the overall cost in terms of simulation coverage, frequency, and runtime can be significantly reduced.

Therefore, several studies have explored hotspot detection using machine learning techniques [1]–[5]. These methods train classifiers using known hotspot and non-hotspot layout patterns and detect unseen hotspot patterns based on learned features. However, false detections still occur, and higher detection accuracy is desired. In machine learning-based approaches,

the design of features that effectively capture characteristics strongly related to hotspots is crucial. A widely used pixel-based feature is Density Based Layout Feature (DBLF) [1][2], which represents the local wiring density in layout patterns. Other proposed pixel-based features include Histogram of Oriented Light Propagation (HOLP) [3], which approximates optical diffraction by smoothing layout images, and Line Width and Separation (LiWS) [6][7], which considers wire widths and the spacing between wires in the layout.

While several features have been proposed for machine learning-based hotspot detection, actual hotspots vary depending on the behavior of light on the wafer surface, which in turn is influenced by the optical characteristics of the exposure system. Since hotspots are caused by light diffraction from the exposure source to the wafer, it is important to consider the optical characteristics (i.e., source characteristics) of the exposure system. These source characteristics are indispensable in lithography simulation and are already known to those who perform hotspot detection. Thus, this information is potentially applicable outside simulation-based approaches.

Some studies do make use of source characteristics for hotspot determination [4][8], but such approaches remain close to machine learning-based lithography simulation. For example, the method from [4] is also regarded as considering optical characteristics, but it is based on the idea of training a model using intensity images generated by lithography simulation. Therefore, it does not directly incorporate the optical parameters of the exposure system into the learning process.

In this study, we propose a new feature that considers the optical characteristics of the exposure system, which have not been taken into account in existing features. This work is an extended and revised version of our previously published technical report [9]. Because hotspots are induced by diffraction of light as it travels from the light source to the wafer, incorporating source characteristics is essential. By leveraging information already available from lithography simulators, our method enables effective hotspot detection in a machine learning framework. Through comparative experiments with existing features, we confirmed that our proposed feature improves detection accuracy. Note that the effectiveness of the proposed feature may depend on the availability and precision of source characteristics provided by the exposure system.

The remainder of this paper is organized as follows. Section II provides an overview of lithography, hotspots, and

machine learning. Section III defines the hotspot detection problem and describes machine learning-based hotspot detection methods and existing features. Section IV explains the proposed feature incorporating source characteristics. Section V presents the experimental results, and Section VI concludes the paper.

II. PRELIMINARIES

In this section, we first explain the mechanism of lithography and the concept of lithography hotspots. We then describe the machine learning framework used for hotspot detection, focusing on supervised learning for classification, and the process of feature extraction required when applying machine learning.

A. Mechanism of Lithography

Lithography is the pattern transfer process in semiconductor fabrication [10]. In lithography, ultraviolet (UV) light is projected onto a photomask, which serves as the master template for semiconductor chips, and the circuit pattern (layout pattern) is transferred onto a silicon wafer through the photomask, as illustrated in Figure 1.

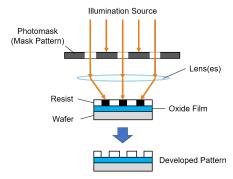


Figure 1. Lithography

B. Lithography Hotspots

With the continued scaling of semiconductor devices, it has become increasingly difficult to accurately transfer designed layout patterns. Examples of degraded pattern fidelity include corner rounding, necking, and line-end shortening, which are caused by Optical Proximity Effects (OPE). To improve the fidelity of pattern transfer, technologies, such as Optical Proximity Correction (OPC) and Sub-Resolution Assist Features (SRAF), have been developed [11]. However, patterns that cannot be accurately transferred still emerge, even with these technologies. Such patterns are referred to as hotspots, and they are one of the factors that degrade yield and reliability of semiconductor products. Figure 2 shows an example of short and open circuits caused by lithography.

Because photomasks, which serve as the masters for layout patterns, are extremely expensive, it is essential to eliminate hotspots at the design stage to avoid costly rework.

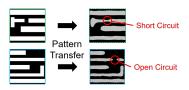


Figure 2. Short and open circuits caused by lithography

C. Optical Simulation for Hotspot Detection

In lithography, UV light is projected onto a photomask, which acts as the master template for semiconductor chips, and the pattern is transferred onto the silicon wafer through the photomask.

In optical simulation for lithography, the exposure process is simulated by calculating the light intensity distribution as the light emitted from the source passes through the photomask and projection lens system and reaches the photoresist. In this simulation, the optical characteristics of the exposure system (i.e., source characteristics) are represented by a matrix called the Sum of Coherent Systems (SOCS) kernel [12]. The post-exposure light intensity distribution is computed as the squared magnitude of the convolution between the SOCS kernel in the spatial domain and the layout pattern.

Let ϕ_j be the j-th kernel matrix and M the layout pattern matrix. The resulting light intensity distribution I is given by:

$$I(x,y) = \sum_{j=1}^{n} \sigma_j |(\phi_j * \mathbf{M})(x,y)|^2, \tag{1}$$

where the symbol * denotes convolution and σ_i is a constant.

D. Machine Learning and Feature Extraction

It is time-consuming and labor-intensive for humans to analyze large amounts of data and derive rules or conditions related to specific phenomena. To address this, machine learning [13] has attracted attention as a technique that enables computers to learn from large-scale data and automatically construct models or algorithms for tasks, such as classification and prediction.

Machine learning can be categorized into supervised and unsupervised learning. Supervised learning involves estimating a mapping function based on given input data and corresponding labeled outputs. Among supervised learning tasks, classification aims to predict the class label of a given instance.

In classification problems, raw data alone often fails to achieve sufficient prediction accuracy. To address this, relevant elements are extracted from raw data in a process known as feature extraction. If two such elements are denoted as x and y, then the vector $\mathbf{f}=(x,y)$ is referred to as a feature vector. In this paper, we refer to both the feature vector itself and its components as features.

III. HOTSPOT DETECTION USING MACHINE LEARNING

In this section, we define the hotspot detection problem and describe a method for detecting hotspots using machine learning. We also introduce two widely used features: DBLF, which considers layout density, and HOLP, which approximates optical diffraction effects.

A. Hotspot Detection Problem

A hotspot refers to a pattern in lithography that poses a risk of causing an open or short circuit. Whether a given pattern is a hotspot can only be determined through lithography simulation or after actual semiconductor fabrication. The hotspot detection problem is to identify such hotspot patterns from a layout, based on known hotspot and non-hotspot examples.

B. Detection Method Using Machine Learning

Hotspot detection using machine learning consists of two main phases: the training phase and the testing phase. In this study, layout patterns are assumed to be represented as bitmap images, where wiring areas are white (pixel value: 1) and empty areas are black (pixel value: 0). The flow of machine learning-based hotspot detection is illustrated in Figure 3.

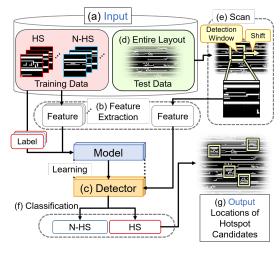


Figure 3. Flow of hotspot detection using machine learning [7]

In the training phase, as shown in Figure 3(a), a set of known hotspot and non-hotspot pattern images are provided as training data. Feature extraction is performed on each image to obtain features (Figure 3(b)). These features, along with the corresponding class labels indicating whether the image is a Hotspot (HS) or Non-Hotspot (N-HS), are input to a machine learning algorithm. The model is then trained to construct a hotspot classifier (Figure 3(c)).

In the testing phase, as shown in Figure 3(d), an image of the entire layout is given as test data. A region of interest used to determine whether a hotspot is present is referred to as a detection window. The detection window is scanned over the layout image (as in Figure 3(c)), and for each region corresponding to the detection window, feature extraction is performed and the extracted features are input into the trained classifier. The classifier outputs predicted labels (Figure 3(f)), enabling hotspot detection.

C. Existing Feature: DBLF

The feature DBLF considers the density of wiring in the layout, i.e., the proportion of area occupied by wires. The procedure to compute DBLF is as follows. The image corresponding to a detection window of $i \times i$ pixels is divided into $N \times N$ subregions, each consisting of $k \times k$ pixels (Figure 4). These subregions are referred to as local regions.

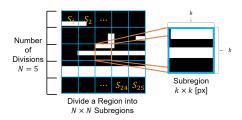


Figure 4. Division of image and local regions

For each local region s_l , the wire area ratio d_l is calculated. The DBLF feature is then represented as the vector of these values, as shown in (2):

$$\mathbf{F}_{\text{DBLF}} = (d_1, d_2, \dots, d_{N^2}).$$
 (2)

D. Existing Feature: HOLP

The feature HOLP approximates the effect of optical diffraction in lithography by smoothing layout images.

First, a Gaussian filter is applied to the image M corresponding to a detection window to produce a smoothed image M_S as shown in Figure 5(a). For each local region of M_s , intensity gradients are computed, as illustrated in Figure 5(b), and a histogram is constructed using the intensity gradients. Gradient

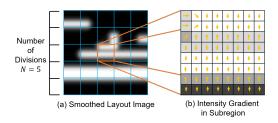


Figure 5. Gradient computation from smoothed image [3]

angles are quantized into bins as shown in Figure 6(a). In the example with 8 bins, the illustrated gradient falls into bin 4, and weighted voting is applied using gradient magnitude as the weight, as shown in Figure 6(b). Next, each local histogram is normalized so that the sum of all bin values equals 1. For a local region s_l , the histogram with B bins is represented as a vector $(g_1^l, g_2^l, \ldots, g_B^l)$. The overall HOLP feature is obtained by concatenating the histograms from all local regions, as shown in (3):

$$\mathbf{F}_{\text{HOLP}} = (g_1^1, g_2^1, \dots, g_B^1, \dots, g_1^{N^2}, \dots, g_B^{N^2}). \tag{3}$$

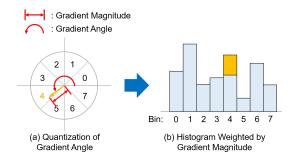


Figure 6. Construction of gradient histograms [3]

IV. PROPOSED FEATURE

Since hotspots are caused by optical effects as light travels from the exposure system's light source to the wafer, it is important to consider information about the optical system. Therefore, we propose a novel feature vector, named Optical System-Aware Mapping (OSAM), which incorporates the optical characteristics of the exposure system (hereafter, source characteristics).

A. Feature Considering Optical System Characteristics

As the source characteristics can be represented using SOCS kernels, we propose a feature based on these kernels. The proposed feature vector is derived from a simplified light intensity distribution calculated using reduced versions of both the mask pattern and the kernels, as well as by truncating the number of kernel components. By reducing the mask and kernel sizes in advance and limiting the kernel order, the computation time becomes significantly shorter compared to full lithography simulations.

B. Computation Procedure of the Proposed Feature

The computation procedure of the proposed feature is described below. We assume that the pattern image and the kernels have the same size.

First, the mask pattern image \mathbf{M} is divided into $N \times N$ regions, where N is a user-defined constant. Next, for each $k \times k$ -pixel local region s_l in the pattern image, the proportion of area occupied by wires is calculated as d_l , and a new $N \times N$ matrix \mathbf{M}' is formed using these values in the same way as in DBLF, as illustrated in Figure 7.

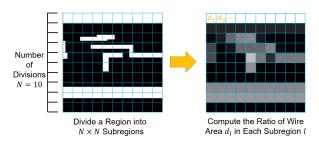


Figure 7. Downsampling of pattern image

Then, the kernels ϕ_j $(j = 1, 2, \dots n')$ in the spatial domain are divided into $k \times k$ -pixel local regions such that the center

of the central local region aligns with the center of the kernels, where n'(< n) is a user-defined integer constant and n is the original number of SOCS kernels used in the full optical simulation.

That is, to avoid splitting the central part of the kernels, each kernel is divided into $(N-1)\times (N-1)$ local regions when N is even, and into $N\times N$ regions when N is odd. When N is even, the peripheral areas of the kernels are ignored.

For each kernel, the average value e_l of the pixels in each local region s_l is calculated, forming a matrix ϕ'_j of size $(N-1)\times (N-1)$ or $N\times N$ depending on whether N is even or odd, as shown in Figure 8.

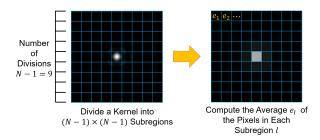


Figure 8. Downsampling of kernel

The simplified light intensity distribution I' is calculated using the following equation:

$$I'(x,y) = \sum_{j=1}^{n'} \sigma_j |(\phi'_j * \mathbf{M}')(x,y)|^2.$$
 (4)

The central $N \times N$ submatrix C from I' is then flattened to form the proposed feature vector \mathbf{F}_{OSAM} , as shown in (5):

$$\mathbf{F}_{\text{OSAM}} = (C_{1,1}, C_{1,2}, \dots, C_{N,N-1}, C_{N,N}). \tag{5}$$

Note that C can be obtained without computing the full convolution results, by restricting the computation to the region of interest

The overall flow of computing the proposed feature is illustrated in Figure 9.

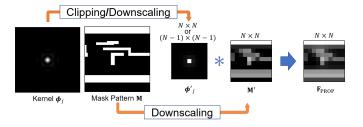


Figure 9. Flow of proposed feature computation

In our experiments, the proposed feature OSAM is used in combination with DBLF to enhance detection performance.

V. EXPERIMENTAL RESULTS

In this section, we compare the hotspot detection accuracy between the existing feature DBLF and the proposed feature that considers optical system characteristics. Experiments were conducted on a Linux server equipped with an Intel Xeon E5-2620 v4 2.2GHz processor and 128GB of memory. The experimental programs were implemented in Python 3.6.10, using OpenCV, Cython, scikit-learn, and TensorFlow as libraries.

To evaluate the proposed method, a dataset with specified optical conditions is required. We used data relabeled using an optical simulator [14] based on the ICCAD 2012 CAD contest dataset [10], as shown in Table I.

TABLE I. RELABELED ICCAD 2012 CONTEST DATASET

| Circuit | #Samples | Process | #HS | #N-HS |
|---------|----------|---------|------|-------|
| data1 | 545 | 32nm | 246 | 299 |
| data2 | 4644 | 28nm | 1417 | 3227 |
| data3 | 5349 | 28nm | 2930 | 2419 |
| data4 | 3563 | 28nm | 1100 | 2463 |
| data5 | 2152 | 28nm | 625 | 1527 |

A. Comparison Using AdaBoost

To compare the proposed feature with DBLF, we used AdaBoost [15], a commonly used machine learning algorithm in feature-based hotspot detection. The implementation used scikit-learn, and decision trees were used as weak learners in the ensemble model.

To focus on the fundamental performance of each feature, we evaluated hotspot classification (not full detection). In the experiments, 70% of the HS and N-HS regions in each dataset were used for training, and the remaining 30% were used for testing. The classification results were categorized into four types, as shown in Table II.

TABLE II. CLASSIFICATION ACCURACY CATEGORIES

| | Hotspot | Non-hotspot |
|--------------------------|---------------------|---------------------|
| Predicted as Hotspot | True Positive (TP) | False Positive (FP) |
| Predicted as Non-hotspot | False Negative (FN) | True Negative (TN) |

As an evaluation metric, we used the F1 score, which is the harmonic mean of Precision and Recall, defined as follows:

$$Precision = \frac{TP}{TP + FP}$$

$$TP$$

$$TP$$

$$TP$$

$$TP$$

$$Recall = \frac{TP}{TP + FN} \tag{7}$$

$$Recall = \frac{TP}{TP + FN}$$
(7)

$$F1_score = 2 \times \frac{Precision \times Recall}{Precision + Recall}$$
(8)

For DBLF, the pattern image was divided into N=10, resulting in a 100-dimensional feature vector. For the proposed feature, the pattern was also divided with N=10, while the kernel was divided with N-1=9 to align the region center. To match region sizes, the outer region of the kernel was excluded. Although the optical simulator's kernel order n was 24, we set n'=1 for the proposed feature. Thus, the proposed feature had 100 dimensions, and since it was concatenated with DBLF, the total dimensionality became 100 + 100 = 200.

To ensure a fair comparison, we explored all 225 combinations of the following hyperparameters for each feature, selecting those that yielded the highest F1 scores:

- Number of weak learners: 2, 4, 6, ..., 1000
- Maximum tree depth: 2, 3, 4
- Learning rate: 0.95, 0.96, 0.97, 0.98, 0.99

These hyperparameters follow prior studies [7]. Each dataset was tested 5 times, and the average F1 score was computed. The best and average F1 scores across all combinations are shown in Table III. In all tables presented in this section, PROP denotes our proposed feature, OSAM.

TABLE III. F1 Scores with All Parameter Combinations (DBLF 10×10

| | Best | | | Average | | |
|---------|--------|--------|----------|---------|--------|----------|
| Dataset | DBLF | PROP | (Ratio) | DBLF | PROP | (Ratio) |
| data1 | 0.8640 | 0.8769 | (1.0149) | 0.8111 | 0.8271 | (1.0197) |
| data2 | 0.6199 | 0.6090 | (0.9824) | 0.5332 | 0.5370 | (1.0071) |
| data3 | 0.8399 | 0.8456 | (1.0067) | 0.7907 | 0.8171 | (1.0334) |
| data4 | 0.8281 | 0.8240 | (0.9950) | 0.6672 | 0.7475 | (1.1204) |
| data5 | 0.7993 | 0.8260 | (1.0334) | 0.6544 | 0.7515 | (1.1423) |
| Average | 0.7902 | 0.7963 | (1.0077) | 0.6913 | 0.7360 | (1.0647) |

To compare under more similar dimensionality, we also tested DBLF with $14 \times 14 = 196$ dimensions. Results are shown in Table IV.

TABLE IV. F1 Scores with All Parameter Combinations (DBLF 14×14

| | Best | | | Average | | |
|---------|--------|--------|----------|---------|--------|----------|
| Dataset | DBLF | PROP | (Ratio) | DBLF | PROP | (Ratio) |
| data1 | 0.8711 | 0.8769 | (1.0067) | 0.8111 | 0.8271 | (1.0197) |
| data2 | 0.6247 | 0.6090 | (0.9749) | 0.5332 | 0.5370 | (1.0071) |
| data3 | 0.8279 | 0.8456 | (1.0214) | 0.7907 | 0.8171 | (1.0334) |
| data4 | 0.8306 | 0.8240 | (0.9921) | 0.6672 | 0.7475 | (1.1204) |
| data5 | 0.8191 | 0.8260 | (1.0084) | 0.6543 | 0.7515 | (1.1423) |
| Average | 0.7947 | 0.7963 | (1.0020) | 0.6913 | 0.7360 | (1.0647) |

As shown in Table III, the proposed feature performed better than DBLF in both best and average F1 scores. Computation times were comparable. Table IV further shows that even with nearly equal dimensionality, the proposed feature still performed better than DBLF. These results indicate the effectiveness of incorporating source characteristics in hotspot detection.

We also investigated whether the proposed feature can be further improved by maximizing the optical detail, ignoring computation time. We used simulation images directly as feature vectors, resizing them to control dimensionality. These were concatenated with DBLF as in the proposed method. Results are shown in Table V.

TABLE V. F1 Scores Using Simulation Images as Features

| Dataset | 10×10 | 15×15 | 20×20 | 30×30 | 50×50 | 100×100 |
|---------|----------------|----------------|----------------|----------------|----------------|------------------|
| data1 | 0.8839 | 0.8885 | 0.8849 | 0.8848 | 0.8772 | 0.8771 |
| data2 | 0.7164 | 0.7284 | 0.7176 | 0.7003 | 0.6797 | 0.6768 |
| data3 | 0.8735 | 0.8823 | 0.8726 | 0.8688 | 0.8665 | 0.8617 |
| data4 | 0.8611 | 0.8987 | 0.8864 | 0.8823 | 0.8811 | 0.8714 |
| data5 | 0.8441 | 0.8613 | 0.8329 | 0.8372 | 0.8267 | 0.8225 |
| Average | 0.8358 | 0.8518 | 0.8389 | 0.8347 | 0.8263 | 0.8219 |

These results suggest that while the proposed feature can be further improved by adjusting kernel order or partition size, its performance is already strong at 10×10 . Unexpectedly, the best score occurred at 15×15 , not at higher resolutions, indicating that the relabeled dataset is challenging, and that AdaBoost may perform better with lower-dimensional input. Similar findings have been reported in [16].

B. Comparison Using CNN

Based on the previous findings, we also performed experiments using Convolutional Neural Networks (CNNs) instead of AdaBoost. We used TensorFlow for implementation.

We compared DBLF and the proposed feature using the same parameter (N=10). The CNN consisted of five layers: conv1-pool1-conv2-pool2-dense. Each convolutional layer used ReLU activation, with filters of size 3×3 and stride 1. The number of filters was 16 in conv1 and 32 in conv2. Each pooling layer performed max pooling with a 2×2 filter. We experimented with all combinations of epochs $\{10, 20, 30, 40, 50\}$ and batch sizes $\{16, 32, 64, 128, 256\}$. The best F1 scores are shown in Table VI.

TABLE VI. BEST F1 SCORES USING CNN (FEATURE SIZE 10×10)

| Dataset | DBLF | PROP | (Ratio) |
|---------|--------|--------|----------|
| data1 | 0.8774 | 0.8662 | (0.9872) |
| data2 | 0.6552 | 0.6714 | (1.0247) |
| data3 | 0.8531 | 0.8697 | (1.0195) |
| data4 | 0.8471 | 0.8405 | (0.9922) |
| data5 | 0.7880 | 0.8126 | (1.0312) |
| Average | 0.8042 | 0.8120 | (1.0098) |

From Tables III and VI, both features improved in F1 score using CNN. From Table VI, the proposed feature performed better than DBLF by approximately 1%, suggesting its potential effectiveness. We further experimented using simulation images as features in CNN, comparing 10×10 and 20×20 sizes. Results are shown in Table VII.

TABLE VII. BEST F1 SCORES USING CNN WITH DIFFERENT FEATURE

| | Feature Size 10×10 | | | Feature Size 10×10 Feature Size 20×20 | | |
|---------|-----------------------------|--------|----------|---|--------|----------|
| Dataset | DBLF | Sim. | (Ratio) | DBLF | Sim. | (Ratio) |
| data1 | 0.8774 | 0.8774 | (1.0000) | 0.8533 | 0.8701 | (1.0197) |
| data2 | 0.6552 | 0.6583 | (1.0047) | 0.6910 | 0.7477 | (1.0821) |
| data3 | 0.8531 | 0.8659 | (1.0150) | 0.8755 | 0.8972 | (1.0248) |
| data4 | 0.8471 | 0.8513 | (1.0050) | 0.8649 | 0.9356 | (1.0817) |
| data5 | 0.7880 | 0.8486 | (1.0769) | 0.8563 | 0.9065 | (1.0586) |
| Average | 0.8042 | 0.8203 | (1.0200) | 0.8282 | 0.8714 | (1.0522) |

These results indicate that the proposed feature has room for improvement, but already performs well at 10×10 , supporting the practicality and effectiveness of our approach.

VI. CONCLUSION AND FUTURE WORK

In this study, we proposed a novel feature for machine learning-based hotspot detection that incorporates the optical characteristics of the exposure system, which are typically overlooked in existing approaches. Experimental comparisons with existing features showed that the proposed feature consistently improved detection performance.

As future work, we aim to improve the runtime efficiency of the proposed approach, for example, by performing convolutions in the frequency domain.

REFERENCES

- [1] T. Matsunawa, J.-R. Gao, B. Yu, and D. Z. Pan, "A new lithography hotspot detection framework based on AdaBoost classifier and simplified feature extraction," in *Design-Process-Technology Co-optimization for Manufacturability IX*, vol. 9427, SPIE, Mar. 2015, pp. 94270S1–11.
- [2] Y.-T. Yu, G.-H. Lin, I. H.-R. Jiang, and C. Chiang, "Machine-learning-based hotspot detection using topological classification and critical feature extraction," *IEEE Trans. Comput.-Aided Design Integr. Circuits Syst.*, vol. 34, no. 3, pp. 460–470, Jan. 2015
- [3] Y. Tomioka, T. Matsunawa, C. Kodama, and S. Nojima, "Lithography hotspot detection by two-stage cascade classifier using histogram of oriented light propagation," in 2017 22nd Asia and South Pacific Design Automation Conference (ASP-DAC), Jan. 2017, pp. 81–86.
- [4] J. W. Park, A. Torres, and X. Song, "Litho-aware machine learning for hotspot detection," *IEEE Trans. Comput.-Aided Design Integr. Circuits Syst.*, vol. 37, no. 7, pp. 1510–1514, Jul. 2018.
- [5] T. Zhou *et al.*, "Mining lithography hotspots from massive sem images using machine learning model," in *2021 China Semiconductor Technology Int. Conf. (CSTIC)*, Mar. 2021, pp. 1–3.
- [6] G. Kataoka, M. Inagi, S. Nagayama, and S. Wakabayashi, "Novel feature vectors considering distances between wires for lithography hotspot detection," in 2018 21st Euromicro Conference on Digital System Design (DSD), Aug. 2018, pp. 85–90.
- [7] G. Kataoka, M. Yamamoto, M. Inagi, S. Nagayama, and S. Wakabayashi, "Feature vectors based on wire width and distance for lithography hotspot detection," *IPSJ Trans. System* and LSI Design Methodology, vol. 16, pp. 2–11, Feb. 2023.
- [8] T. Matsunawa, T. Kimura, and S. Nojima, "Lithography hotspot candidate detection using coherence map," in *Design-Process-Technology Co-optimization for Manufacturability XIII*, J. P. Cain, Ed., vol. 10962, SPIE, Mar. 2019, pp. 109620Q1–8.
- [9] M. Yamamoto, M. Inagi, and S. Nagayama, "A feature vector considering characteristics of optical system for lithography hotspot detection," in *IEICE Technical Report (VLD2022-81)*, in Japanese, vol. 122, Mar. 2023, pp. 49–54.
- [10] J. A. Torres, "ICCAD-2012 CAD contest in fuzzy pattern matching for physical verification and benchmark suite," in 2012 IEEE/ACM Int. Conf. Computer-Aided Design (ICCAD), Nov. 2012, pp. 349–350.
- [11] T. Matsunawa, B. Yu, and D. Z. Pan, "Optical proximity correction with hierarchical Bayes model," *J. Micro/Nanolithography*, *MEMS*, and *MOEMS*, vol. 15, no. 2, pp. 021009-1–8, Mar. 2016.
- [12] N. Cobb, "Sum of coherent systems decomposition by SVD," University of California, Berkeley, Technical Report, Sep. 1995, pp. 1–7.
- [13] S. Raschka, Y. H. Liu, and V. Mirjalili, Machine Learning with PyTorch and Scikit-Learn: Develop machine learning and deep learning models with Python. Packt Publishing Ltd., Feb. 2022.
- [14] S. Banerjee, Z. Li, and S. R. Nassif, "ICCAD-2013 CAD contest in mask optimization and benchmark suite," in 2013 IEEE/ACM Int. Conf. Computer-Aided Design (ICCAD), Nov. 2013, pp. 271–274.
- [15] D. Solomatine and D. Shrestha, "AdaBoost.RT: A boosting algorithm for regression problems," in 2004 IEEE Int. Joint Conf. Neural Networks, vol. 2, Jul. 2004, pp. 1163–1168.
- [16] X. Liu, Y. Dai, Y. Zhang, Q. Yuan, and L. Zhao, "A preprocessing method of adaboost for mislabeled data classification," in 2017 29th Chinese Control And Decision Conf. (CCDC), May 2017, pp. 2738–2742.