

A Principled Design of Image Representation: Towards Forensic Tasks

Shuren Qi, Yushu Zhang, *Member, IEEE*, Chao Wang,
Jiantao Zhou, *Senior Member, IEEE*, and Xiaochun Cao, *Senior Member, IEEE*

Abstract—Image forensics is a rising topic as the trustworthy multimedia content is critical for modern society. Like other vision-related applications, forensic analysis relies heavily on the proper image representation. Despite the importance, current theoretical understanding for such representation remains limited, with varying degrees of neglect for its key role. For this gap, we attempt to investigate the forensic-oriented image representation as a distinct problem, from the perspectives of theory, implementation, and application. Our work starts from the abstraction of basic principles that the representation for forensics should satisfy, especially revealing the criticality of robustness, interpretability, and coverage. At the theoretical level, we propose a new representation framework for forensics, called *dense invariant representation* (DIR), which is characterized by stable description with mathematical guarantees. At the implementation level, the discrete calculation problems of DIR are discussed, and the corresponding accurate and fast solutions are designed with generic nature and constant complexity. We demonstrate the above arguments on the dense-domain pattern detection and matching experiments, providing comparison results with state-of-the-art descriptors. Also, at the application level, the proposed DIR is initially explored in passive and active forensics, namely copy-move forgery detection and perceptual hashing, exhibiting the benefits in fulfilling the requirements of such forensic tasks.

Index Terms—Dense invariant representation, image forensics, orthogonal moments, covariance, fast Fourier transform.

1 INTRODUCTION

FAKE visual media are now a real-world threat due to their potential abuse to a number of critical scenarios, covering journalism, publishing, judicial investigations, and social networks. For such serious challenges, *image forensics*, aimed at verifying the authenticity and integrity, has attracted widespread attention [1].

Similar to other visual-understanding applications, forensic analysis, i.e., knowledge extraction for the traces of manipulation, relies heavily on proper *image representation*. This is not so surprising, as Herbert Simon asserted: “solving a problem simply means representing it so as to make the solution transparent” [2]. Therefore, a common practice in forensic algorithms is to inherit the representation techniques directly from other vision tasks.

We note that, despite the conceptual similarities, image representations for forensic tasks are quite different from the others in the design goals.

Let us begin the analysis from the basic principles that an efficient representation for forensics should satisfy: *Discriminability* – the representation is sufficiently informative to support the distinction between pristine and manipulated data; *Robustness* – the representation is not influenced by geometric or signal perturbations that may be introduced by the adversary; *Coverage* – the representation sufficiently covers the entire image plane, as manipulation may occur in any region including the background; *Interpretability* – the representation should have reliable theoretical guarantees, implying that causation is more important than correlation, due to the role in courts of judicature or in public discussion and debate; *Computational efficiency and accuracy* – the representation should have a reasonable implementation, especially for the widely-used dense processing under constraint of coverage.

1.1 State of the Art and Motivation

Starting from the above principles, one can note that the representations in popular vision tasks (e.g., image classification) may not be the best choice for forensic tasks, mainly due to their robustness, interpretability, and coverage properties.

The representations based on *deep learning*, e.g., Convolutional Neural Networks (CNN), have shown impressive performance on a variety of high-level vision tasks. Their success is mainly due to the powerful data-adaptive capability provided by the composition of multiple nonlinear transformations with trainable parameters [3]. Despite the

- S. Qi is with the College of Computer Science and Technology, Nanjing University of Aeronautics and Astronautics, Nanjing, China, and also with the Institute of Information Engineering, Chinese Academy of Sciences, Beijing, China (e-mail: shurenqi@nuaa.edu.cn).
- Y. Zhang is with the College of Computer Science and Technology, Nanjing University of Aeronautics and Astronautics, Nanjing, China, and also with the Guangxi Key Laboratory of Trusted Software, Guilin University of Electronic Technology, Guilin, China (e-mail: yushu@nuaa.edu.cn).
- C. Wang is with the College of Computer Science and Technology, Nanjing University of Aeronautics and Astronautics, Nanjing, China (e-mail: c.wang@nuaa.edu.cn).
- J. Zhou is with the State Key Laboratory of Internet of Things for Smart City, and also with the Department of Computer and Information Science, Faculty of Science and Technology, University of Macau, Macau, China (e-mail: jtzhou@umac.mo).
- X. Cao is with the School of Cyber Science and Technology, Shenzhen Campus of Sun Yat-sen University, Shenzhen, China (e-mail: caoxiaochun@mail.sysu.edu.cn).
- Corresponding author: Y. Zhang (e-mail: yushu@nuaa.edu.cn).

IEEE Transactions on Pattern Analysis and Machine Intelligence, 2022, doi: 10.1109/TPAMI.2022.3204971, link: ieeexplore.ieee.org/document/9881995/.

strong discriminability, two inherent properties make them inappropriate for many forensic tasks. The first is the difficulty in achieving satisfactory robustness [4], especially for geometric changes [5] and adversarial examples [6], which may be used to mislead forensic algorithm by the adversary, i.e., anti-forensics [1]. The second is the difficulty in understanding the decision mechanism [7], and such black-box nature reduces the credibility of the forensic results, especially for the critical scenarios like courts [1]. We must underline that the research on robustness and interpretability for deep learning is developing [8], [9] but still appears a long way from the requirements of forensic tasks.

Turning now to the topic of hand-crafted feature descriptors. The representations based on *interest points*, e.g., Scale Invariant Feature Transform (SIFT), are very popular and competitive approaches for practical feature engineering. As a distinctive property, keypoint detectors and descriptors are commonly designed to yield a high repeatability under various geometric transformations [10]. Such interest points, however, are sparse on the image plane, leading to information loss in some regions. This is clearly detrimental to the tasks that require high coverage, such as fine-grained classification [11] and also the forensics [1]. In addition, current research shows that SIFT-like keypoints can be removed or injected through imperceptible pixel modifications [12]. This inherent vulnerability has potential applications in the anti-forensics, similar to the adversarial attacks for deep learning.

For these reasons, more forensic researchers have switched to the representation based on *dense sampling*, e.g., DAISY descriptor [13]. Such dense representation formed on regular sampling grid, possibly with a range of scales, thus ensuring not only the coverage but also the simple spatial relationships. This regular pattern may be useful in modeling Markov processes, constructing spatial-pooling representations, accelerating feature matching, etc.; while for interest points, the spatial relationship is more arbitrary, resulting in more complicated processing [14]. The main challenges, in dense approach, are geometric invariance and implementation efficiency. Existing methods are not stable under geometric distortions, even for common rotation, scaling, and flipping. Moreover, due to the dense nature, the implementation is generally time-consuming, which also limits the use of expensive-but-invariant features.

1.2 Contributions

Motivated by above facts, we attempt to present a principled study on the image representation for forensics, covering theory, implementation, and application. To the best of our knowledge, this is a very early work on the theoretical understanding of forensic-oriented image representation, and it was rarely considered as a distinct problem before.

Our main contributions are summarized as follows.

Theory. We propose a unified representation framework, named Dense Invariant Representation (DIR), by extending the definition of classical orthogonal moments from the global to the local with scale space. The mathematical analysis explains the important properties of DIR, e.g., the description stability based on covariance. Accordingly, our

TABLE 1
Theoretical Comparison With Related Methods

Method	CVPR'08 [19]	TSP'15 [20]	TPAMI'19 [21]	Ours
Generic Design				✓
Rotation Invariance	✓		✓	✓
Constant Complexity		✓	✓	✓

DIR is able to fulfill the core requirements of forensic tasks, especially for robustness, interpretability, and coverage.

Implementation. We derive an accurate and fast numerical calculation of the DIR. For the accuracy, it is dominated by numerical instability and integration error in the computation of basis functions. The corresponding solutions are recursive strategy and high-precision numerical integration method. For the efficiency, the convolution theorem and scaling theorem of Fourier transform are used to speed up the dense inner products. The resulting complexity does not depend on the window size and thereby is the constant complexity $\mathcal{O}(1)$. Note that such implementation is generic for arbitrary basis functions.

Application. We validate the effectiveness of DIR in various simulation experiments and real-world applications. Experiments are performed on dense-domain pattern detection and matching, exhibiting the state-of-the-art performance under challenging geometric transformations or signal corruptions. The direct applications to copy-move forgery detection (passive forensics) and perceptual hashing (active forensics) also demonstrate the accuracy and efficiency gains.

1.3 Related Works

We briefly review the topics of dense descriptor and image forensics that are closely related to our work.

Dense Descriptor. In the literature, the designs of dense features is typically based on frequency transform, texture, and orthogonal moments. For the frequency transform, well-known method like Walsh-Hadamard Transform [15] and Haar Wavelet [16] have been designed for dense pattern matching. Their research focuses on computational complexity and robustness to noise-like attacks, with few considerations on geometric invariance [17]. For the texture, the descriptors DAISY [13] and DASC [18] aim to model the image local structure through gradient distribution and self-similarity, respectively. We regarded them as competitive methods, since their robustness to variable imaging conditions, at the geometric level (e.g., viewing angle) and signal level (e.g., blurring), is highlighted in the application of stereo matching. For orthogonal moments, dense descriptors based on Fourier-Mellin transform [19], Tchebichef moments [20], and Zernike moments [21] have appeared in tasks of object detection and image forensics, for their desirable invariance and independence [22]. Compared to these theoretically relevant methods, the distinctive properties of our work are summarized in Table 1. As illustrated later, the discussion of DIR on properties and calculations is fully generic to a class of basis functions; also DIR exhibits in-form invariance to rotation and constant complexity to window size.

Image Forensics. For the fake content generation, previously, the manipulation was typically performed by Photoshop-like editing software, with the operations such as copy-move [23], splicing [24], and inpainting [25]. In recent years, as the rise of deep learning, the so-called deepfake [26] allows one to generate realistic fake content in a flexible way. For the fake content detection, the solutions are mainly divided into active and passive forensics. Active ones rely on specific information embedding (e.g., digital watermarking [27]) or extracting (e.g., perceptual hashing [28]) for the image prior to distribution. Owing to the side information, these methods have desirable property to detect any type of manipulations, but they obviously require additional implementation costs and cannot be used for the images that have been distributed. Passive forensic algorithms, on the contrary, do not rely on such prior processing; they work exclusively on the given image itself. It generally operates by seeking the inconsistencies of given image at the digital [29], physical [30], or semantic [31] level, which are inevitably introduced by certain manipulations. On the downside, this line of methods is unsatisfactory in terms of stability and generality, due to the lack of side information. Therefore, it is necessary to design a forensic-oriented image representation for fundamentally improving the accuracy and efficiency of active and passive approaches.

2 FOUNDATIONS

For the sake of completeness, we briefly remind some foundations of classical orthogonal moments (see [32] for a survey). In Table 2, we list core notations for this paper.

Mathematically, the general theory of *image moments* is based on the definition of the following inner product $\langle f, V_{nm} \rangle$ [22]:

$$\langle f, V_{nm} \rangle = \iint_D V_{nm}^*(x, y) f(x, y) dx dy, \quad (1)$$

with the image function f and the basis function V_{nm} of order $(n, m) \in \mathbb{Z}^2$ on the domain $D \in \mathbb{R}^2$, where “*” denotes the complex conjugate.

As far as image representation tasks are concerned, two constraints are typically imposed on the explicit forms of basis functions.

Orthogonality is a core for achieving various beneficial properties in signal analysis, e.g., information preservation and decorrelation. It occurs when any two basis functions V_{nm} and $V_{n'm'}$ satisfy the condition:

$$\begin{aligned} \langle V_{nm}, V_{n'm'} \rangle &= \iint_D V_{nm}(x, y) V_{n'm'}^*(x, y) dx dy \\ &= \delta_{nm} \delta_{m'm'}, \end{aligned} \quad (2)$$

where $\delta_{\alpha\beta}$ is the Kronecker delta function: $\delta_{\alpha\beta} = [\alpha = \beta]$. For a set of orthogonal functions, the *completeness* in a Hilbert space means its linear span is dense in the space.

Invariant in form w.r.t. rotations of axes about the origin $(x, y) = (0, 0)$ is a desirable structure in geometric analysis. It means that the basis function V_{nm} is expressed in polar coordinates (r, θ) and is of the form:

$$V_{nm}(\underbrace{r \cos \theta}_x, \underbrace{r \sin \theta}_y) \equiv V_{nm}(r, \theta) = R_n(r) A_m(\theta), \quad (3)$$

TABLE 2
Notations and Definitions

Notation	Definition
f	The image function
$(x, y)/(r, \theta)/(i, j)$	The Cartesian/polar/pixel coordinates
V	The basis function
$(n, m)/(u, v)/w$	The order/position/scale parameters
D	The domain of V
A/R	The angular/radial basis function
\mathcal{R}	The image representation
\mathcal{D}	The image degradation
$f_{T/R/VF/HF/S}$	The translated/rotated/vertical-flipped/horizontal-flipped/scaled versions of f
$S_{nm/uv/w}$	The sampling sets of $(n, m)/(u, v)/w$
$\#_{nm/uv/w}$	The sizes of $S_{nm/uv/w}$
K	The constraint on S_{nm}
(N, M)	The size of f
h	The integral value of V over a valid pixel region
H	The kernel w.r.t. h for dense representation
$\langle \cdot, \cdot \rangle$	The inner product
$*$	The conjugate of complex
\angle	The phase of complex
\exp	The natural exponential function
j	The imaginary unit
\mathcal{F}	The Fourier transform
T	The matrix transpose
$ \cdot $	The absolute value
$\ \cdot\ _p$	The p -norm
\otimes	The convolution
\odot	The point-wise (Hadamard) product

with *angular basis function* $A_m(\theta) = \exp(jm\theta)$ ($j = \sqrt{-1}$) and *radial basis function* $R_n(r)$ could be of any form [33].

By imposing the above two constraints, we now write down a rotation-invariant and orthogonal version of (1) in polar coordinates:

$$\langle f, V_{nm} \rangle = \iint_D R_n^*(r) A_m^*(\theta) f(r, \theta) r dr d\theta, \quad (4)$$

where $R_n(r)$ should satisfy the weighted orthogonality condition: $\int_0^1 R_n(r) R_{n'}^*(r) r dr = \frac{1}{2\pi} \delta_{nn'}$ and the domain is generally the unit disk: $D = \{(r, \theta) : r \in [0, 1], \theta \in [0, 2\pi)\}$.

In Appendix A, we include some common definitions of $R_n(r)$ for (4), with orthogonality and completeness. Note that the main discussion in this paper is generic to all such definitions.

3 DENSE INVARIANT REPRESENTATION: THEORY

This section is dedicated to our continuous-domain representation framework for forensics. Firstly, we explicitly give the mathematical definition of DIR and clarify its theoretical relationship with classical orthogonal moments. Secondly, the properties of DIR regarding geometric transformations are analyzed on the basis of covariance.

3.1 Definition Extension and Basic Formula

One important nature of the classical orthogonal moments with form (1) is the *global definition*: the image function f and the basis function V_{nm} share the same coordinate system (Cartesian or polar), more precisely, the same origin and similar scale.

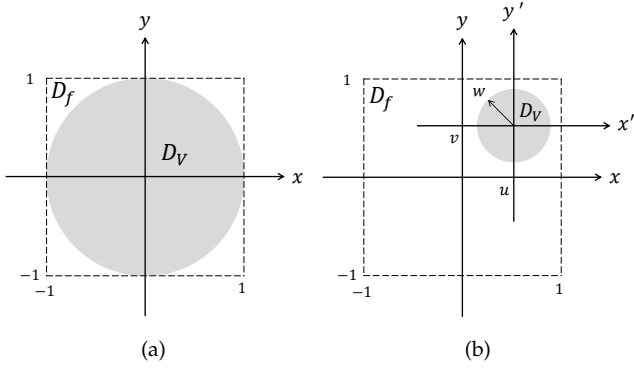


Fig. 1. An illustration of the definition extension from the global (a) to the local with scale space (b).

Theoretically, such fact renders the classical orthogonal moments useless for local descriptions, and in turn prevents their application in many practical vision problems. For the case of image forensics, especially the forgery localization tasks, a pixel-by-pixel local representation is often a fundamental composition as the coverage requirement. Obviously, the classical definition (1) does not have such a capability.

Here, we aim to meet the coverage requirement in forensic tasks by generalizing the classical definition of image moments. Considering a local coordinate system (x', y') for basis function V_{nm} , which is a translated and scaled version of the global coordinate system (x, y) with translation offset (u, v) and scale factor w . Their explicit relationship is expressed as follows:

$$(x', y') = \frac{(x, y) - (u, v)}{w}, \quad (5)$$

and hence the corresponding polar coordinates are:

$$\begin{cases} r' = \sqrt{(x')^2 + (y')^2} = \frac{1}{w} \sqrt{(x-u)^2 + (y-v)^2} \\ \theta' = \arctan(\frac{y'}{x'}) = \arctan(\frac{y-v}{x-u}) \end{cases} \quad (6)$$

Definition 1. With the local polar coordinates (6), a local definition of the rotation-invariant and orthogonal moments (4) can be derived as follows:

$$\langle f, V_{nm}^{uvw} \rangle = \iint_D \underbrace{R_n^* \left(\frac{\sqrt{(x-u)^2 + (y-v)^2}}{w} \right)}_{(V_{nm}^{uvw}(x,y))^*} \underbrace{A_m^* \left(\arctan \left(\frac{y-v}{x-u} \right) \right)}_{\theta'} f(x, y) dx dy, \quad (7)$$

where the domain is a disk with center (u, v) and radius w : $D = \{(x, y) : (x-u)^2 + (y-v)^2 \leq w^2\}$.

The resulting new definition (7) plays a foundational role in constructing the DIR, and therefore it is treated as the **basic formula** throughout this paper. We would like to make a comment on the notation V_{nm}^{uvw} : the superscript and subscript denote the parameters for spatial and frequency domains, respectively. As illustrated in Fig. 1, this new definition allows the domain of image D_f and the domain of basis function D_V to be built in the separate coordinate

systems, providing greater flexibility than the classical orthogonal moments.

Two interesting propositions can be directly observed from this definition, i.e., *generic nature* and *local representation capability*.

Proposition 1. The classical form (4) of orthogonal moments is a special case for the basic formula (7), with fixed parameters $(u, v) = (0, 0)$ and $w = 1$.

Hence, the basic formula (7) of DIR is regarded as a unified mathematical framework for the research of moments and moment invariants.

Proposition 2. With the basic formula (7), the center and scale of local description region D can be controlled by adjusting the parameters (u, v) and w , respectively.

This local representation capability is also a distinct characteristic of DIR that the classical orthogonal moments do not have, leading to potential use in the local behavior based image processing and visual understanding.

From an analytical perspective, the related descriptors in [19], [21] can be derived from the basic formula (7) by giving a fixed definition of $R_n(r)$. This fact means that our approach is a more generic design, and following discussion on the properties and calculations of (7) is also fully generic. As for related work [20], its Cartesian definition of basis functions is inconsistent with (3) and therefore inconsistent with (7). As illustrated later, such difference limits it to achieve the desirable rotation invariance.

3.2 Important Properties and Representation Formulas

The representation robustness is quite critical for a range of visual forensic algorithms, due to its intrinsic *two-player* nature [1]. Skilled adversary may introduce the trace removal operations such as geometric transformations and signal corruptions to interfere with forensic analysis.

In Section 3.1, a local definition of image moments has been derived, satisfying the coverage requirement. Also, some beneficial properties of classical moment theory, such as completeness and orthogonality, are inherited in (7). However, the local robustness of (7) for geometric transformations has not been investigated, which is clearly distinct from the classical theory based on global assumption.

Next, we aim to meet the robustness requirement in forensic tasks by analyzing the important representation properties under various image transformations.

Our analysis relies on three terms of the representation: *invariance*, *equivariance*, and *covariance* [34], [35]. Considering a representation \mathcal{R} and a degradation \mathcal{D} , such terms correspond to the following three identities:

- invariance – $\mathcal{R}(\mathcal{D}(f)) \equiv \mathcal{R}(f)$,
- equivariance – $\mathcal{R}(\mathcal{D}(f)) \equiv \mathcal{D}(\mathcal{R}(f))$,
- covariance – $\mathcal{R}(\mathcal{D}(f)) \equiv \mathcal{D}'(\mathcal{R}(f))$,

where \mathcal{D}' is a composite function of \mathcal{D} . In fact, covariance is a generalized expression of invariance and equivariance. Thus, in general, invariant/equivariant representation can be constructed under the premise that covariance holds.

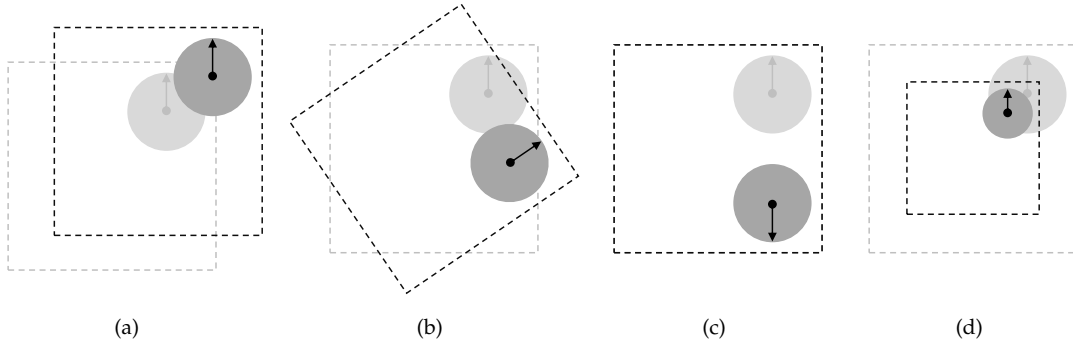


Fig. 2. An illustration of the geometric transformations: translation (a), rotation (b), flipping (c), and scaling (d).

3.2.1 Equivariance to Translation

Suppose f_T is a translated version of image function f with offset $(\Delta x, \Delta y)$, i.e., $f_T(x, y) = f(x + \Delta x, y + \Delta y)$, as shown in Fig. 2.

Property 1. By plugging f_T into (7), it can be checked that image translation operation only affects the representation parameters (u, v) , as follows:

$$\langle f_T(x, y), V_{nm}^{uvw}(x, y) \rangle = \langle f(x, y), V_{nm}^{(u+\Delta x)(v+\Delta y)w}(x, y) \rangle, \quad (8)$$

where the same offset $(\Delta x, \Delta y)$ also appears in representation $\langle f(x, y), V_{nm}^{(u+\Delta x)(v+\Delta y)w}(x, y) \rangle$, implying the equivariance w.r.t. translation.

Proof. The proof of (8) is given in Appendix B. \square

Proposition 3. With the translation equivariance (8), it is feasible to retrieve the translated version of a given pattern over the domain of (u, v) by brute-force feature matching. Moreover, the translation-invariant features can be generated by permutation-invariant mapping, such as average or maximum pooling, over the domain of (u, v) .

As seen later, this translation equivariance is also crucial in the derivation of the properties for rotation, flipping, and scaling.

3.2.2 Invariance to Rotation and Flipping

Suppose f_R is a rotated version of image function f with angle ϕ around the center, i.e., $f_R(r, \theta) = f(r, \theta + \phi)$, where polar coordinates are used for convenience.

Considering any pair of corresponding circular regions in f_R and f , their geometric relationship can be modeled as a composite of center-aligned rotation and translation, as shown in Fig. 2. Since the translation equivariance has been confirmed, the rest of the analysis will focus only on the center-aligned rotation, i.e., we can restrict the parameters $(u, v) = (0, 0)$ without loss of generality.

Property 2. By plugging f_R into (7) with $(u, v) = (0, 0)$ and writing down in polar form, it can be checked that image rotation operation only affects the phase of representation, as follows:

$$\langle f_R(r, \theta), V_{nm}^{uvw}(r', \theta') \rangle = \langle f(r, \theta), V_{nm}^{uvw}(r', \theta') \rangle A_m^*(-\phi), \quad (9)$$

where the same angle ϕ also appears in phase of the representation $\langle f(r, \theta), V_{nm}^{uvw}(r', \theta') \rangle A_m^*(-\phi)$, implying the covariance w.r.t. rotation.

Proof. The proof of (9) is given in Appendix C. \square

Proposition 4. Considering the covariance (9) of phase, the rotation invariance will hold straightforwardly in the magnitude domain: $|\langle f_R(r, \theta), V_{nm}^{uvw}(r', \theta') \rangle| = |\langle f(r, \theta), V_{nm}^{uvw}(r', \theta') \rangle|$.

Towards rotation invariance, more generalized and information-preserving strategy is based on the proper phase cancellation, which eliminates the effect of the term $A_m^*(-\phi)$. Here, the resulting feature is typically complex-valued containing phase information. In addition to the invariants, retrieving the rotation angle from the phase is easy to solve and robust to noise. For more details on phase cancellation and angle estimation, we address the reader to [36].

Proposition 5. By analogy to above derivation, the representations of horizontally flipped version $f_{HF}(r, \theta) = f(r, \pi - \theta)$ and vertically one $f_{VF}(r, \theta) = f(r, -\theta)$ can be written as $(-1)^m (\langle f(r, \theta), V_{nm}^{uvw}(r', \theta') \rangle)^*$ and $(\langle f(r, \theta), V_{nm}^{uvw}(r', \theta') \rangle)^*$, respectively. Therefore, the flipping invariance will also hold in the magnitude domain.

3.2.3 Covariance to Scaling

Suppose f_S is a scaled version of image function f with factor s around the center, i.e., $f_S(x, y) = f(sx, sy)$.

Similarly, the geometric relationship of any pair of corresponding circular regions in f_S and f is a composite of center-aligned scaling and translation, as shown in Fig. 2. Taking advantage of translation equivariance, the following derivation is carried out under the setting $(u, v) = (0, 0)$.

Property 3. By plugging f_S into (7) with $(u, v) = (0, 0)$, it can be checked that image scaling operation only affects the representation parameters w , as follows:

$$\langle f_S(x, y), V_{nm}^{uvw}(x, y) \rangle = \langle f(x, y), V_{nm}^{uv(ws)}(x, y) \rangle, \quad (10)$$

where the same factor s also appears in the representation $\langle f(x, y), V_{nm}^{uv(ws)}(x, y) \rangle$, implying the covariance w.r.t. scaling.

Proof. The proof of (10) is given in Appendix D. \square

Proposition 6. With the scaling covariance (10), it is feasible to retrieve the scaled version of a given pattern over the domain of w

by brute-force feature matching. Moreover, the scaling-invariant features can be generated by permutation-invariant mapping, such as average or maximum pooling, over the domain of w .

We remark that, for the discrete case, this continuous-domain invariance will degenerate to a certain tolerance, as the dense sampling of the scales is difficult in practice. In addition to the permutation-invariant mapping, *scale selection* is also a common path to achieve invariance. Here, the state-of-the-art works in dense scale selection [37] are easy to combine with our DIR.

3.2.4 Orthogonality and Completeness

Following the classical moment theory, the set of V_{nm}^{uvw} forms an orthogonal and complete basis for Hilbert space, which in turn ensures the uniqueness and independence of the moments [22], [32]. Such information-preservation architecture can thus provide complementary information that improves the discriminability.

Additionally, following the classical moment theory, the representation with low-order (n, m) is more stable under the signal corruptions such as noise, blur, and sampling/quantization effect. The reason is based on the following fact: the lossy signal operations generally work on the high-frequency components of the image, i.e., they mainly affect the high-order moments [22], [32]. In Appendix E, we provide a formal analysis of the robustness to signal corruption from a frequency-domain perspective.

In order to allow the robust low-frequency information preferentially used for the image representation, the set of orders (n, m) is able to be constrained by a specific norm and a integer constant K :

$$S_{nm}(K) = \{(n, m) : \|(n, m)\|_p \leq K\}, \quad (11)$$

where $\|\cdot\|_p$ denotes the p -norm and the value of p was typically taken as 1 or ∞ in the literature.

Such properties (8) – (11) provide interpretable guidelines for constructing robust local features, and therefore they are treated as the **representation formulas** throughout this paper. Note that it is not wise to cover all these invariance types, due to the concerns about discriminability and efficiency, the best way is to take what we need in the application.

In related works [19], [20], [21], the properties w.r.t. translation, flipping, scaling, and signal corruptions have hardly been discussed. The [19], [21] exhibit the rotation invariance due to similar covariance in (9); while [20] does not have such property as a result of the Cartesian structure.

4 DENSE INVARIANT REPRESENTATION: IMPLEMENTATION

The basic and representation formulas have been described in Sections 3.1 and 3.2, both of them in the continuous domain. Now, we move away from such analytic formulas and concentrate on the discrete implementation of DIR for the digital images. Here, the accuracy and efficiency problems encountered in the discrete implementation of DIR are discussed, and corresponding solutions are provided.

4.1 Accurate Computation

The digital-level forensic analysis relies on high-order statistics, e.g., PRNU-like method [29]. In such scenarios, computational errors can even dominate the forensic performance. As illustrated later, the direct implementation of DIR is not suitable for the above scenarios.

Next, we aim to meet the computational accuracy requirement in forensic tasks by reducing the approximation error and representation error.

Considering a digital image over the discrete Cartesian grid $\{f(i, j) : (i, j) \in \{1, 2, \dots, M\} \times \{1, 2, \dots, N\}\}$, we introduce a continuous version $(x, y) \in [1, M] \times [1, N]$ of discrete variables (i, j) for convenience. Here, a (i, j) -centered pixel region is thus defined as $D_{ij} = \{(x, y) \in [i - \frac{\Delta i}{2}, i + \frac{\Delta i}{2}] \times [j - \frac{\Delta j}{2}, j + \frac{\Delta j}{2}]\}$; the value of f over this region is constant and equal to $f(i, j)$.

Definition 2. With above notations, the basic formula (7) can be rewritten into discrete form, as follows:

$$\langle f, V_{nm}^{uvw} \rangle = \sum_{(i,j) \text{ s.t. } D_{ij} \cap D \neq \emptyset} h_{nm}^{uvw}(i, j) f(i, j), \quad (12)$$

where $h_{nm}^{uvw}(i, j)$ is the integral value of $(V_{nm}^{uvw})^*$ over the intersection of (i, j) -centered pixel region D_{ij} and the domain of definition D , i.e.,

$$h_{nm}^{uvw}(i, j) = \iint_{D_{ij} \cap D} (V_{nm}^{uvw}(x, y))^* dx dy. \quad (13)$$

From a practical perspective, the calculation accuracy of DIR is dominated by (13), due to the continuous integration of complicated functions. More specifically, it normally requires i) determining a numerical integration strategy and ii) calculating the values of the basis functions at the sampling points. Next, we will discuss such two aspects separately.

For the strategy of numerical integration, the *Zero-Order Approximation* (ZOA) is a popular algorithm that directly sets $h_{nm}^{uvw}(i, j) \simeq (V_{nm}^{uvw}(i, j))^* \frac{\Delta i \Delta j}{w^2}$. Mathematically, its error depends on the frequency of the function, and thus the error may be very significant when the order (n, m) is large or the scale parameter w is small. Note that both cases are common in DIR. Therefore, we suggest that *high-precision numerical integration* methods, e.g., pseudo up-sampling [38] and Gaussian quadrature rule [39], should be introduced depending on the application requirements. Such strategies can be uniformly formulated as follows.

Definition 3. With the L -dimensional cubature formulas, a numerical approximation of the $h_{nm}^{uvw}(i, j)$ in (13) can be derived as:

$$h_{nm}^{uvw}(i, j) \simeq \sum_{(a,b)} c_{ab} (V_{nm}^{uvw}(x_a, y_b))^* \frac{\Delta i \Delta j}{w^2}, \quad (14)$$

where $(x_a, y_b) \in D_{ij}$ is the sampling point with corresponding weight c_{ab} , and $\#\{(a, b)\} = L$.

Proposition 7. Regarding the approximation error of (14) with given (n, m) , it is in order $\mathcal{O}((\frac{\Delta i \Delta j}{w^2})^{L+1})$. Hence, when larger values of L are used, higher accuracy can be achieved w.r.t. the error order $\mathcal{O}((\frac{\Delta i \Delta j}{w^2})^2)$ of ZOA.

For the calculation of basis functions, numerical instability is a common concern, mainly due to the factorial/gamma

terms in the Jacobi polynomial based basis functions. Here, we recommend the *recursive strategy* [40] to derive the high-order basis functions directly from several low-order ones, without the factorial/gamma of large number.

Here, the discretization (12) – (14) are treated as the **accurate computation formulas** throughout this paper.

4.2 Fast Computation

Starting from the coverage and robustness for forensic analysis, we prefer DIR to work in a dense manner, i.e., the parameters (u, v) and w are sufficiently sampled. In this scenario, the direct computation from definition will exhibit considerable complexity, and thus efficient implementation of DIR is desired in practice.

Next, we aim to meet the computational efficiency requirement in forensic tasks by introducing some useful theorems and data structures.

Definition 4. Considering a dense sampling of (u, v) over the discrete grid $\{1, 2, \dots, M\} \times \{1, 2, \dots, N\}$, an equivalent version of (12) can be derived as follows:

$$\langle f, V_{nm}^{uvw} \rangle = f(i, j) \otimes (H_{nm}^w(i, j))^T, \quad (15)$$

where \otimes denotes the convolution operation and $(\cdot)^T$ indicates the matrix transpose; H_{nm}^w is a kernel defined by h_{nm}^{uvw} with following form:

$$H_{nm}^w(i, j) = \{h_{nm}^{uvw}(i, j) : u, v = w, (i, j) \text{ s.t. } D_{ij} \cap D \neq \emptyset\}. \quad (16)$$

Here, the dense inner product in (12) is converted to the convolution in (15). Note that we use “convolution” instead of the similar term “cross-correlation” is to facilitate subsequent discussion.

Regarding the computational complexity of (15) with given (n, m) , it is determined by the size of kernel H_{nm}^w (i.e., $4w^2$), size of sample set S_{uv} for parameter (u, v) (denoted as $\#_{uv}$), and size of sample set S_w for parameter w (denoted as $\#_w$); hence a total of $\Theta(4 \sum_{w \in S_w} w^2 \#_{uv})$ or $\mathcal{O}(w_{\max}^2 \#_{uv} \#_w)$ multiplications, where w_{\max} is the maximum in S_w . As for (16), the complexity is almost negligible, because the kernel only needs to calculate once for given n, m and w , without dense processing.

Now, we show that there are efficient ways to conduct the DIR.

Definition 5. Let us introduce the convolution theorem of the Fourier transform, such that the spatial-domain convolution in (15) can be converted to a frequency-domain product as [41]:

$$\langle f, V_{nm}^{uvw} \rangle = \mathcal{F}^{-1}(\mathcal{F}(f) \odot \mathcal{F}((H_{nm}^w)^T)), \quad (17)$$

where \mathcal{F} denotes the Fourier transform and \odot indicates the point-wise multiplication.

Proposition 8. If the Fast Fourier Transform (FFT) algorithm is used for the implementation of (17), the multiplication complexity will become $\Theta(\#_w(3\#_{uv} \log \#_{uv} + \#_{uv}))$ or $\mathcal{O}(\#_w \#_{uv} \log \#_{uv})$.

It should be highlighted that the scale parameter w has no role in the complexity, meaning a constant-time calculation w.r.t. the kernel/window size. Ideally, when w_{\max} is

large enough such that w_{\max}^2 is greater than $\log(\#_{uv})$, this FFT-based method will take less time than the convolution-based one in (15). This observation is crucial due to the fact that the sampling of w generally involves some large values, for covering the scale variations in the application.

Further, we point out where the above FFT-based algorithm can be accelerated.

Definition 6. Let us introduce the scaling theorem of the Fourier transform, allowing the frequency coefficients $\mathcal{F}((H_{nm}^w)^T)$ to be derived directly from the calculated ones with scale parameter w_0 :

$$\mathcal{F}((H_{nm}^w)^T) = \left(\frac{w}{w_0}\right)^2 \mathcal{F}((H_{nm}^{w_0})^T) \left(\frac{w\xi_i}{w_0}, \frac{w\xi_j}{w_0}\right), \quad (18)$$

where (ξ_i, ξ_j) are frequency variables.

Proposition 9. In practice, therefore, for given (n, m) and a pair of w and w_0 , the direct FFT can be replaced with the computationally inexpensive interpolation as (18), ideally saving $\Theta(\#_{uv} \log \#_{uv})$ multiplications.

In addition to convolution and scaling theorems, we also note an implementation trick for some practical scenarios.

Proposition 10. Going back to (17), a valuable observation is that the term $\mathcal{F}((H_{nm}^w)^T)$ is independent of the content of the input f and is only relevant to its size. Thus, a lookup table with indexes (n, m, w) can be pre-calculated, which contains the ready-to-use terms $\mathcal{F}((H_{nm}^w)^T)$ of fixed size $M_0 \times N_0$. Also, the size of the digital image f will be normalized to $M_0 \times N_0$. Therefore, for given (n, m) and w , the lookup table strategy ideally saves $\Theta(\#_{uv} \log \#_{uv})$ multiplications.

Here, the implementation (15) – (18) are treated as the **fast computation formulas** throughout this paper.

In related works [19], [20], [21], the accurate computation has hardly been discussed. The [20], [21] derive fast computation strategy with the constant complexity w.r.t. kernel/window size, but both rely heavily on specific basis functions. In contrast, our constant-time calculation is built on a generic framework, regardless of the specific definition of basis functions.

It is worth to remark that accurate strategy (14) is mainly for small w , while the fast strategy (17) is more beneficial with large w .

5 PRACTICAL INTUITIONS AND GUIDELINES

In this section, we answer several practical questions that could be raised in the applications of DIR.

5.1 Practical Intuitions

For the practical intuitions of Sections 3 and 4, we illustrate the DIR by considering the analysis of a texture image, as shown in Fig. 3¹.

Computation. Going back to (13), (14), and (16), one can estimate a set of kernels H_{nm}^w , which are derived from the

1. The real/imaginary parts of DIR kernels and the phases of DIR coefficients are displayed in a colormap ranging from blue to white to red, where 0 corresponding to the white; the magnitudes of DIR coefficients are displayed in a colormap ranging from black to yellow, where 0 corresponding to the black.

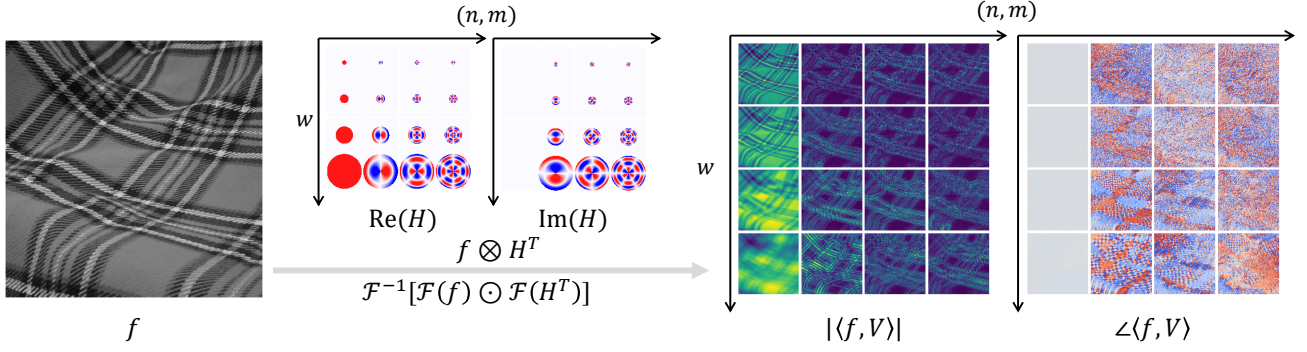


Fig. 3. A high-level intuition of the proposed dense invariant representation.

basis functions V_{nm}^{uvw} with different orders (n, m) and scales w . The real and imaginary parts of H_{nm}^w are given in Fig. 3. Then, by applying spatial-domain convolution (15) of image f with the transpose of such kernels H_{nm}^w , or its equivalent frequency-domain counterpart (17), the image f is decomposed into dense coefficients $\langle f, V_{nm}^{uvw} \rangle$ along multiple orders and scales.

Representation. As illustrated in Fig. 3, the magnitudes of such coefficients capture the main structural information (e.g., the line patterns), while the phases provide very detailed properties (e.g., the rich textures). For robust representation, the invariants can be simply built by magnitude (w.r.t. rotation and flipping) and pooling (w.r.t. translation and scaling) functions over the coefficients, referring to properties (8) – (11). Typically, this magnitude-only strategy is sufficient in practice. As for more informative representation, the expressive phases can be further encoded into the features, based on appropriate phase cancellation.

5.2 Practical Guidelines

Here, we give some general principles for setting parameters and radial basis functions, which support the applications of DIR for practical forensic problems.

Considering the clear physical meaning of DIR parameters (n, m) , (u, v) , and w , their settings can be directly derived from the mathematical properties of the forensic problem.

- Regarding the (n, m) , it controls the *frequency* of the represented image information. For the forensics that reveal the semantic/physical/digital artifacts w.r.t. low/mid/high-frequency information (following the taxonomy of [1]), the pertinent forensic analysis can be achieved by low/mid/high-order DIR.
- Regarding the (u, v) , it controls the *position* of the represented image information. In general, the dense sampling is a good choice for avoiding false negatives (under the coverage principle). In practice, such dense strategy may sometimes be compromised to interval sampling for reducing complexity.
- Regarding the w , it controls the *scale* of the represented image information. For a given forensic task that does not require scale invariance, w can be directly taken as a fixed value. While for the task that requires such invariance, w should be sampled sufficiently. Note that if the prior knowledge w.r.t. the

scale of the pattern under analysis, such knowledge should be introduced in the setting to avoid unnecessary sampling.

As for the definition of radial basis function in DIR (as listed in Appendix A), the choice mainly refers to the representation capability and computational accuracy/efficiency.

- Regarding the representation capability, it is mainly influenced by the distribution of the zeros of radial basis function. Theoretically, a uniform distribution of zeros on $[0, 1]$ is optimal. Another alternative path is to combine radial basis functions with complementary distributions of zeros for capturing complementary image information [40].
- Regarding the computational accuracy/efficiency, it is mainly influenced by the basic mathematical properties of radial basis function. Here, the complexity is related to whether factorial/gamma terms, summation/series operations, and root-finding processes are involved. The numerical stability is related to whether the factorial/gamma terms and very high absolute values are involved.

For the radial basis function, interested readers can access relevant knowledge from our survey paper [32].

We also would like to clarify that when the DIR settings satisfy such general principles, the performance gap between the different settings is generally acceptable (see Appendixes F and G), implying that the DIR is not sensitive to such settings.

6 EXPERIMENTS

In this section, we will evaluate the performance of the proposed DIR, covering experiments at both the descriptor level and the forensic level.

At the descriptor level, the performance statistics for accuracy/complexity and robustness/invariance of DIR are provided. For the accuracy and complexity, several implementation strategies in Section 4 will be compared quantitatively. For the robustness and invariance, our DIR and state-of-the-art dense descriptors are tested in dense detection and matching tasks under geometric or signal perturbations.

At the forensic level, we will measure the usefulness of the proposed DIR in passive and active semantic forensic paths, where the copy-move forgery detection and perceptual hashing are considered as examples, respectively.

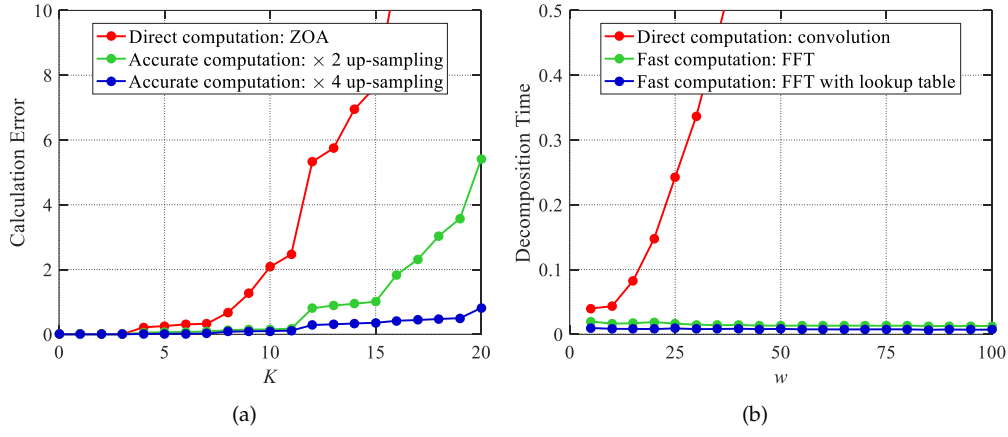


Fig. 4. Calculation error (a) and decomposition time (b) for different implementation methods.

Here, the DIR is applied directly to the algorithm as a feature extraction module. Such a direct application will be compared with state-of-the-art methods in terms of accuracy and efficiency. In Appendixes F and G, we provide some additional forensic experiments; in Appendix H, we discuss the applicability of DIR to forensics beyond exposing semantic artifacts.

For the sake of brevity, a set of cosine functions is chosen as the radial basis functions in all experiments, i.e., a DIR extended from the Polar Cosine Transform (PCT) [42]. In fact, we implement all the DIR extensions for the methods listed in Appendix A. The code is available online at <https://github.com/ShurenQi/DIR>. Note that all experiments are performed in Matlab R2021a, with 2.90-GHz CPU and 16-GB RAM, under Microsoft Windows environment.

6.1 Accuracy and Complexity

Accuracy. Firstly, we evaluate the accuracy performance of different implementation strategies in Section 4.1. Let us consider an image with unity gray-level $\{f_{\text{uni}}(x, y) = 1 : (x, y) \in D\}$, such that the moments with form (7) can be easily derived as:

$$\langle f_{\text{uni}}, V_{nm}^{uvw} \rangle = \begin{cases} 0 & m \neq 0 \\ 2\pi \int_0^1 R_n^*(r') r' dr' & m = 0 \end{cases} \quad (19)$$

This equation indicates that the theoretical value of $\langle f_{\text{uni}}, V_{nm}^{uvw} \rangle$ is zero for $m \neq 0$ case, but in practice such property usually does not hold due to various implementation errors. Hence, we can introduce a simple measure for evaluating the error, named Calculation Error (CE), as follows [32]:

$$\text{CE} = \sum_{(n,m) \in \{S(K), m \neq 0\}} |\langle f_{\text{uni}}, \widehat{V_{nm}^{uvw}} \rangle|, \quad (20)$$

where $\langle f_{\text{uni}}, \widehat{V_{nm}^{uvw}} \rangle$ is the actual calculated value by a discrete implementation.

The comparison methods include the direct computation by ZOA and the accurate computation by pseudo up-sampling (14). As we mentioned, in theory, the error is

positively/inversely proportional to the frequency/number of samples, thus we consider $w = 8$ and $K \in \{0, 1, \dots, 20\}$.

The comparison results of CE are shown in Fig. 4 (a). We observe that the up-sampling strategy yields less error than the ZOA strategy. This observation is especially true for higher up-sampling rate and higher order constraint K . Accordingly, the coefficient with small w and large (n, m) by direct computation may contains considerable errors, hence compromising the discriminability. Such experimental evidence supports our theoretical analysis on integration error, verifying the usefulness of the accurate strategy (14) for real-world scenarios.

Complexity. Secondly, we evaluate the complexity performance of different implementation strategies in Section 4.2. The experiment is designed to compute the dense coefficients for an image of size 512×512 , i.e., $(u, v) \in \{1, 2, \dots, 512\}^2$, under fixed order $(n, m) = (1, 1)$ and varying scales $w \in \{5, 10, \dots, 200\}$. The measure is the Decomposition Time (DT), i.e., the CPU elapsed time in single-thread modality. The comparison methods include the direct computation by convolution (15) and the fast computation by FFT (17), where the trick of lookup table is also considered.

The comparison results of DT are shown in Fig. 4 (b). As w increases, the DT curve for the convolution-based computation rises sharply, while the FFT-based computations exhibit almost-constant time costs. Obviously, this observation is consistent with our theoretical expectation on the complexity. In addition, the introduction of the lookup table further reduces the running time of FFT-based computation. Note that only the computation time for a given (n, m, w) is considered here. While in real-world scenarios, since a set of (n, m, w) is generally needed, the performance gap between the two strategies will be even greater.

6.2 Robustness and Invariance

Two typical vision tasks, dense pattern detection and matching, are chosen for evaluating the robustness and invariance. Note that such two tasks have different natures in the representation: the detection generally works with large window, allowing the use of multi-scale representation as the results are sparse; while the matching usually relies on

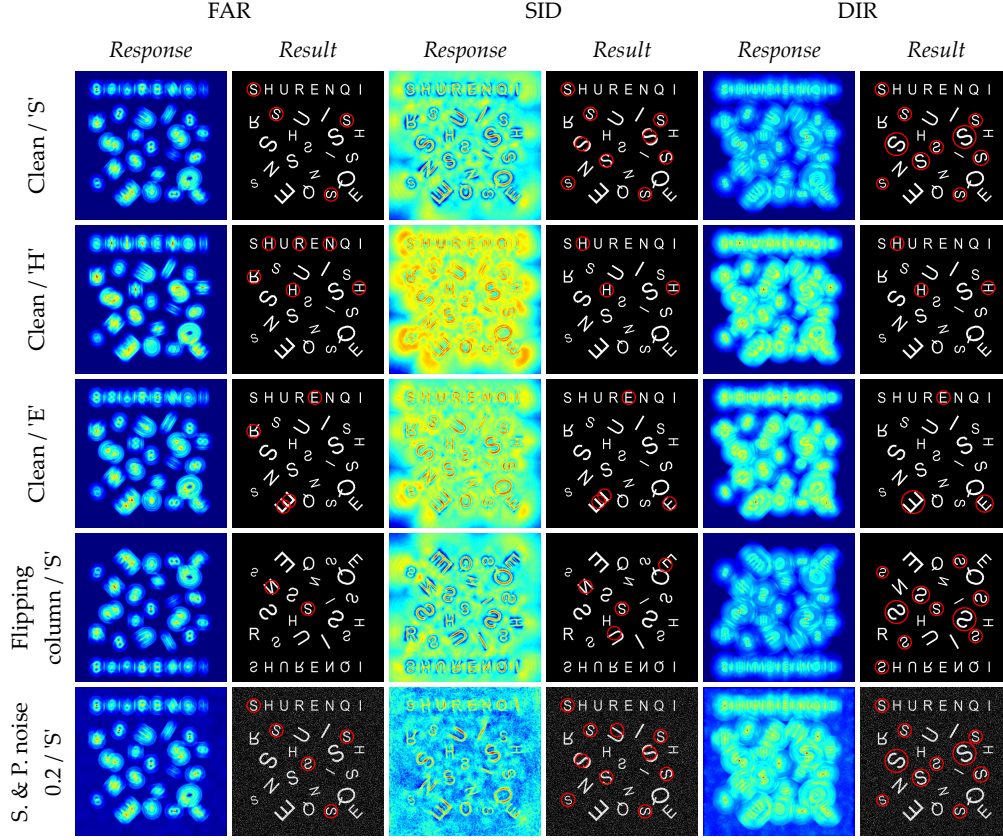


Fig. 5. Some samples of the pattern detection by different dense descriptors.

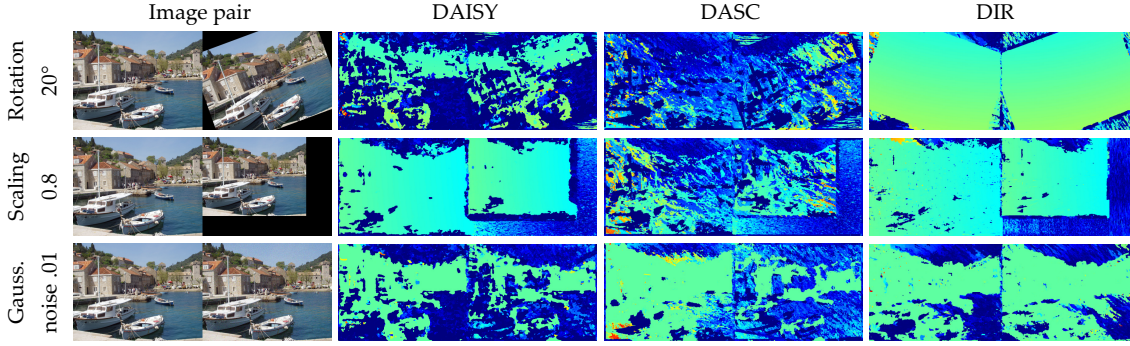


Fig. 6. Some samples of the pattern matching by different dense descriptors.

small window, forcing the use of compact representation as the results are dense.

In the experiments, the algorithm framework is designed in a simple form. This is because we aim to reflect the inherent performance of different descriptors in the framework, rather than directly achieving state-of-the-art results on such tasks.

6.2.1 Dense Detection

For this experiment, we consider detecting five letter 'S', 'H', 'R', 'E', and 'N' from a 800×800 image containing multiple letters and their rotated/scaled/flipped versions, as shown in Fig. 5. Note that other letters are not involved as their detection leads to false positives in backgrounds or lines. Furthermore, this image is globally degraded with different

affine transformations and signal corruptions, leading to more challenging conditions.

Algorithm Design. The template is a patch taken directly from the top left corner of the image, containing a clean letter 'S' with the fixed orientation and size. The calculated feature vector of this template is considered to be the ground-truth for comparing with the dense features from the host image. By calculating the Euclidean distances of such features, we treat the positions with smaller distance values as the pattern detection results. We use the F1 score to evaluate the detection performance.

The competing methods include the state-of-the-art invariant descriptors Fourier-Argand Representation (FAR) [43] and Scale Invariant Descriptor (SID) [19]. Here, the FAR is designed for rotation-invariant pattern detection, with

TABLE 3
F1 Scores (%) for Different Dense Descriptors in Pattern Detection Experiment.

Method	FAR [43]	SID [19]	DIR
Clean	61.54	92.68	97.56
Rotation 20°	60.00	73.08	100.00
Rotation 45°	56.41	56.25	97.56
Flipping column	34.48	31.58	97.56
Flipping row	34.48	26.47	97.56
Scaling 0.8	25.64	35.05	74.07
Scaling 0.5	19.35	19.18	45.61
Gaussian noise 0.01	61.54	92.31	100.00
Gaussian noise 0.02	61.54	92.68	93.33
Salt and pepper noise 0.01	60.00	90.48	97.56
Salt and pepper noise 0.02	60.00	92.68	100.00
Average filtering 7 × 7	43.24	45.57	61.76
Average filtering 9 × 9	40.00	37.11	60.61
Gaussian filtering 7 × 7	45.71	50.75	84.00
Gaussian filtering 9 × 9	43.24	46.34	59.15
Median filtering 7 × 7	60.00	92.68	100.00
Median filtering 9 × 9	61.22	92.68	97.67
JPEG compression 10	61.54	92.68	100.00
JPEG compression 5	61.54	92.68	100.00
Laplacian sharpening	57.89	92.68	100.00
Average ↑	50.47	67.28	88.20
Standard deviation ↓	13.26	27.26	17.12

impressive robustness to severe noise conditions. The SID is characterized by the inherent rotation and scaling invariance, without the need for common pooling or selection operations.

For implementation details of DIR, we set the order constraint $K = 5$ in (11) with ∞ -norm, resulting in 36-dimensional features for each position (u, v, w) in the scale space. Here, we consider a dense sampling of (u, v) over the image grid, and the scale w increases from 30 to 120 with 10 samples. As for FAR and SID, we use the common parameter settings in the original papers. Note that both FAR and SID work on the single scale, and their feature dimensions at each position are 21 and 1008, respectively. In terms of sampling, the FAR is dense, while the SID employs a 2-pixel sampling interval due to the considerable time/space cost.

Robustness and Invariance. Fig. 5 shows some samples of the response map and detection result, while the F1 scores for all above comparison methods are given in Table 3. As can be observed here, under the degradation operations, it is challenging for the descriptor to be robust while maintaining discriminability. The FAR exhibits rotation invariance and higher tolerance for severe noise than other methods. However, the FAR cannot handle scale changes, and the features are relatively unstable under the filtering-like operations. As for the SID, its log-polar sampling allows the rotation/scaling-invariant representation on the single scale. However, the response maps suggest that the features are weak in rejecting the irrelevant letters i.e., less discriminability. In addition, both FAR and SID are not flip-invariant. In contrast, the proposed DIR exhibits stronger stability for rotation, flipping and scaling, which should be attributed to the full exploitation of the covariance. Also, thanks to the orthogonality and completeness of the basis functions,

TABLE 4
Repeatability Scores (%) for Different Dense Descriptors in Pattern Matching Experiment.

Method	DAISY [13]	DASC [18]	DIR
Clean	95.58	94.48	94.26
Rotation 20°	28.67	0.16	85.01
Rotation 45°	0.30	0.35	76.61
Flipping column	0.23	0.15	92.81
Flipping row	0.48	0.13	91.03
Scaling 0.8	72.07	7.54	49.36
Scaling 1.3	74.79	3.20	32.27
Gaussian noise 0.01	23.66	31.43	35.60
Gaussian noise 0.02	15.54	20.73	21.98
Salt and pepper noise 0.01	60.17	63.68	58.96
Salt and pepper noise 0.02	45.44	47.59	45.43
Average filtering 5 × 5	52.56	68.09	60.53
Average filtering 7 × 7	27.58	29.63	18.77
Gaussian filtering 5 × 5	54.81	69.70	64.47
Gaussian filtering 7 × 7	32.87	41.91	29.95
Median filtering 5 × 5	54.09	72.93	59.88
Median filtering 7 × 7	28.84	42.30	32.21
JPEG compression 10	48.49	40.25	35.97
JPEG compression 5	28.15	17.02	13.69
Laplacian sharpening	43.75	79.58	47.39
Average ↑	39.40	36.54	52.31
Standard deviation ↓	25.10	29.58	24.84

our DIR is able to distinguish well between relevant and irrelevant patterns under signal corruptions.

Efficiency. In terms of efficiency, the DIR maintains a reasonable time cost: ~ 2 seconds, ~ 35 seconds, and ~ 6 seconds for FAR, SID, and DIR, respectively, even though our method works on multiple scales.

6.2.2 Dense Matching

For this experiment, we consider establishing dense correspondences between original image and its degraded version, as shown in Fig. 6. Here, 10 images from the INRIA Holidays dataset [44] are selected for providing average experimental results, and such images are normalized to a same size of 1000×1333 .

Algorithm Design. The experiment is performed on a common framework: the PatchMatch [45] for matching the dense features, and the RANSAC [46] for excluding false matches through data modeling. Note that the regular spatial relationship of dense features allows PatchMatch to achieve geometric-invariant matching with high efficiency [47]. We use the repeatability score [10] to evaluate the matching performance.

The competing methods include the state-of-the-art dense descriptors DAISY [13] and DASC [18]. Here, the DAISY is designed for wide-baseline stereo matching, hence considering large perspective distortions. The DASC further developed the idea of DAISY, paying special attention to photometric variations.

For implementation details of DIR, we chose a set of (n, m) with a maximum order of 3, generating 10-dimensional features for each position (u, v, w) in the scale space. Here, we consider a dense sampling of (u, v) over the image grid, and the scale w increases from 8 to 32 with 10 samples. For a compact representation with scaling tolerance, the feature vectors are average-pooled together over

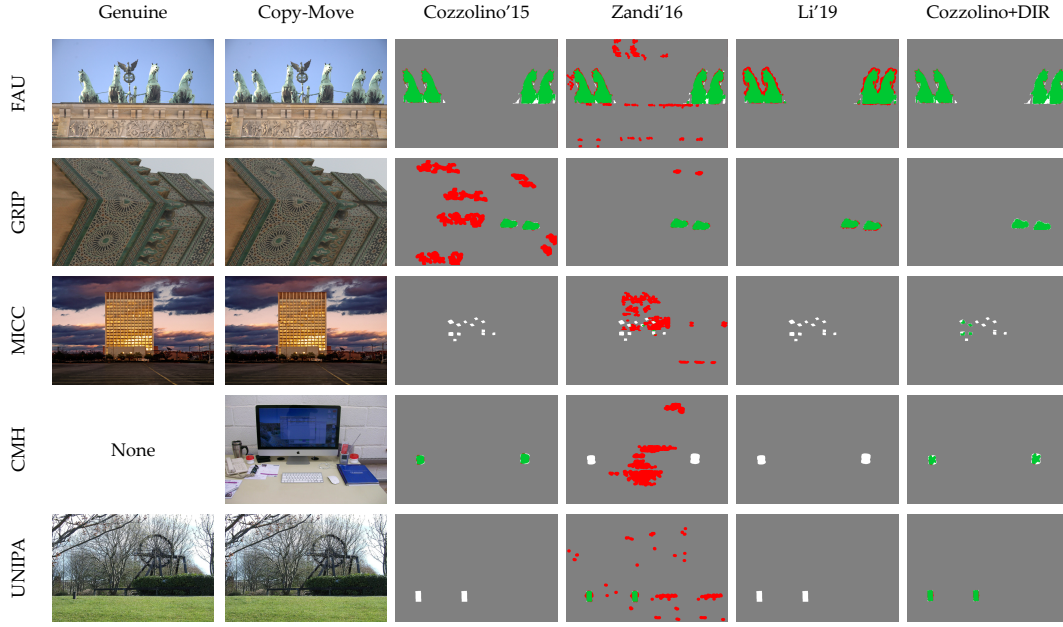


Fig. 7. Some samples of the copy-move detection by different forensic methods on the copy-move forensic benchmarks.

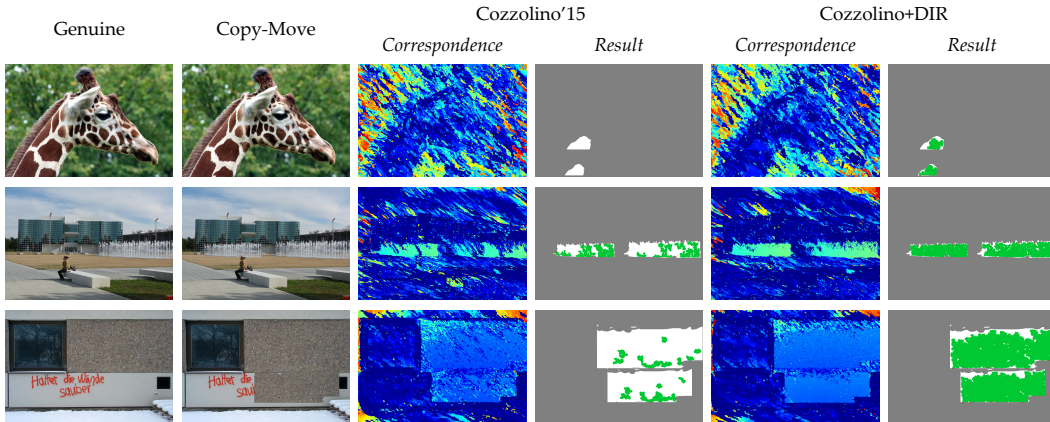


Fig. 8. Some samples of the copy-scale-move detection by different forensic methods.

the scales. As for DAISY and DASC, we use the common parameter settings in the original papers. Note that feature dimensions of DAISY and DASC at each position are 200 and 128, respectively, much more than the 10 dimensions of DIR.

Robustness and Invariance. Fig. 6 shows some samples of the offset length map exported from PatchMatch, while the repeatability scores for all above comparison methods are given in Table 4. Regarding geometric variations, the DAISY and DASC are difficult to maintain a stable matching, especially for the flipping and large-angle rotation. The DAISY exhibits a good repeatability for scaling, while the DASC does not have such a property. As for signal corruptions, the overall repeatability of DASC is better than DAISY, which is consistent with the design goal of DIR. In general, compared with such well-designed descriptors for the matching task, our DIR is typically more robust to geometric variations and provides comparable repeatability against signal corruptions.

TABLE 5
Precision, Recall, and F1 Scores (%) for Different Methods on Various Copy-Move Forensic Benchmarks.

Method		Cozzolino'15 [48]	Zandi'16 [49]	Li'19 [50]	Cozzolino +DIR
Precision	FAU	93.08	79.59	89.44	96.75
	GRIP	93.01	81.9	92.7	96.45
	MICC	92.02	66.06	88.77	94.3
	CMH	82.98	54.96	85.26	88.62
	UNIPA	85.37	73.07	86.18	96.5
Recall	FAU	91.94	95.09	88.94	92.44
	GRIP	96.41	98.64	97.43	95.89
	MICC	89.07	74.98	86.53	88.76
	CMH	78.99	65.27	71.98	75.91
	UNIPA	91.76	99.05	93.95	95.03
F1	FAU	91.89	83.39	88.63	93.62
	GRIP	93.85	86.67	94.75	95.74
	MICC	89.22	66.41	86.73	89.82
	CMH	80.12	58.25	76.35	79.5
	UNIPA	88.31	81.01	89.3	95.66

TABLE 6
Precision, Recall, and F1 Scores (%) for Different Methods on the FAU Copy-Move Forensic Benchmark.

Method	Ryu'13 [51]	Li'13 [52]	Silva'15 [53]	Emam'16 [54]	Pun'18 [55]	Bi'18 [56]	Wu'18 [57]	Zhong'20 [58]	Cozzolino+DIR
Precision	95.02	58.05	88.02	-	91.07	90.55	44.59	75.61	96.75
Recall	88.15	92.26	89.72	-	90.21	91.66	31.60	74.13	92.44
F1	91.50	71.21	88.95	84.91	90.33	91.07	37.11	74.82	93.62

The results for comparison methods in this table are cited directly from [58].

TABLE 7
Precision, Recall, F1 Scores (%), and Matching Performance (Number of Matches per Image) Gain Rate in Copy-Scale-Move Robustness Experiment.

Method	Cozzolino'15 [48]	Cozzolino+DIR	Gain rate
Precision	84.47	90.86	7.56
Recall	46.06	63.32	37.47
F1	56.15	71.18	26.77
#matches/image	144519.75	270730.23	87.33

Efficiency. Turning to efficiency, the feature extraction time for DAISY, DASC, and DIR is ~ 6 seconds, ~ 106 seconds, and ~ 6 seconds, respectively; the feature matching time for DAISY, DASC, and DIR is ~ 79 seconds, ~ 77 seconds, and ~ 20 seconds, respectively. Owing to the proposed constant-time implementation, the feature extraction of DIR is quite fast and comparable to DAISY, which is well known for its efficiency. In addition, the feature vector of DIR is more compact, due to the orthogonality of basis functions, allowing a significant time saving in the matching process.

6.3 Passive Forensics: Copy-Move Forgery Detection

Copy-move is one of the most basic operations in image forgery. It involves copying and pasting specific patches from and to an image. Typical detection algorithm extracts sparse or dense features from the image and reveals potential copy-move regions by matching such features. Obviously, the quality of the extracted features has a significant impact on the performance. In general, sparse methods are more efficient and geometrically invariant but exhibit lower accuracy; conversely, dense methods are more accurate but relatively less efficient and geometrically invariant. This phenomenon is consistent with our analysis in Section 1.1.

Algorithm Design. A representative work of the dense approach is proposed by Cozzolino et al. [48], relying on rotation-invariant orthogonal moments. Theoretically, its feature extraction module can be considered as a special case of DIR, with a fixed scale. As can be expected, in practice, this method is stable under rotation, flipping and

signal corruptions, but sensitive to scale changes. We hence introduce the DIR framework into this algorithm mainly for improving its scaling robustness. Specifically, the original PCT feature extraction module is directly replaced by the DIR extension of PCT. For implementation details of DIR, we chose a set of (n, m) with a maximum order of 3, generating 10-dimensional features for each position (u, v, w) in the scale space. Here, we consider a dense sampling of (u, v) over the image grid, and the scale w increases from 8 to 32 with 10 samples. For a compact representation with scaling tolerance, the feature vectors are average-pooled together over the scales. As a common trick, the low-resolution image will be upsampled (long edge = 2000 pixels) to suppress the parameter sensitivity. Note that we exclude the border pixels between forgery and background in the score computation [23], [48].

Copy-Move Benchmark. Firstly, we perform a quantitative comparison on five copy-move forensic benchmarks: FAU [23], GRIP [48], MICC [59], CMH [53], and UNIPA [60]. Here, FAU, GRIP, and UNIPA contain only rigid copy-move manipulation; while MICC and CMH are with further attacks (e.g., scaling and rotation) for a convincing visual effect. The experiment involves the basic Cozzolino'15 [48] and its DIR version, as well as state-of-the-art sparse algorithms: Silva'15 [53], Zandi'16 [49] and Li'19 [50]; state-of-the-art dense algorithms: Ryu'13 [51], Li'13 [52], Emam'16 [54], Pun'18 [55], Bi'18 [56], Wu'18 [57], Zhong'20 [58]. Here, Wu'18 and Zhong'20 are based on deep neural networks. Note that the feature extraction modules in Ryu'13, Li'13, Emam'16, Pun'18, and Bi'18 can also be considered as special cases of DIR.

Fig. 7 shows some samples of copy-move detection on benchmarks. The precision, recall, and F1 scores for popular open-source algorithms [48], [49], [50] over all above benchmarks are given in Table 5. For the rest of the comparison methods, the scores are summarized in Table 6, only on the FAU due to the lack of the code. It can be observed that, in general, the dense approach, especially Cozzolino'15 and its DIR extension, provides higher detection accuracy than the sparse approach. As far as the sparse approach is concerned, the Li'19 exhibits the performance advantage over the

TABLE 8
Precision, Recall, and F1 Scores (%) for Different Methods in Copy-Scale-Move Robustness Experiment.

Method	Ryu'13 [51]	Li'13 [52]	Emam'16 [54]	Pun'18 [55]	Bi'18 [56]	Wu'18 [57]	Zhong'20 [58]	Cozzolino+DIR
Precision	< 25	< 25	-	41.48	62.09	34.84	68.05	90.86
Recall	< 20	24.95	-	39.69	< 20	20.12	64.59	63.32
F1	< 20	20.97	< 20	40.51	22.91	25.22	64.67	71.18

The results for comparison methods in this table are cited directly from [58].

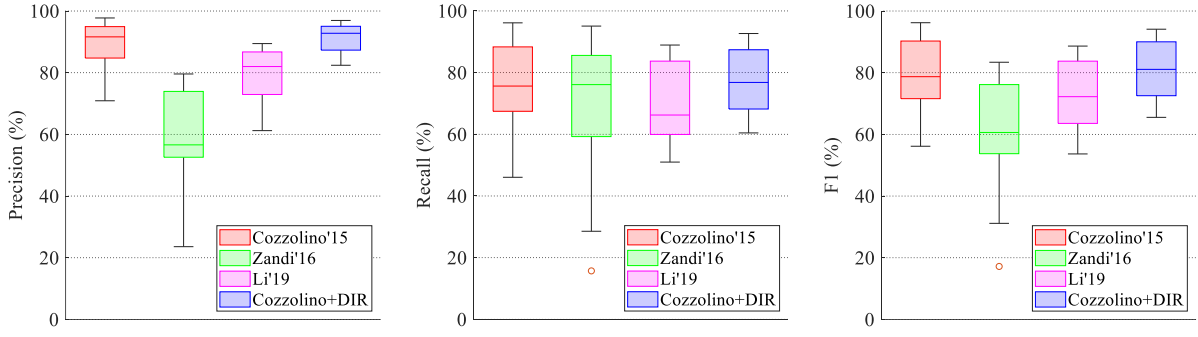


Fig. 9. Precision, recall, and F1 box-plots by different copy-move forgery detection methods in comprehensive robustness experiment.

Zandi'16, which can be mainly attributed to the more dense keypoints. Such phenomenon confirms that the coverage is a vital factor in the image representation for forensics. With the introduction of DIR, the accuracy of Cozzolino'15 on most benchmarks is further improved (basically the same on CMH), meaning that DIR has practical usefulness rather than just as a mathematical extension. Such accuracy gain is largely due to the multi-scale framework in our DIR, which allows features to be more informative also with scaling tolerance.

Copy-Scale-Move Robustness. Secondly, we evaluate the scaling tolerance of the basic Cozzolino'15 and its DIR version. The experiment is designed on FAU benchmark with copy-scale-move manipulation, where the factor is set to 0.8. The experiment also covers some state-of-the-art dense algorithms as comparison baselines: Ryu'13 [51], Li'13 [52], Emam'16 [54], Pun'18 [55], Bi'18 [56], Wu'18 [57], and Zhong'20 [58].

Fig. 8 shows some samples of the correspondence map and detection result for copy-scale-move. The performance statistics are summarized in Tables 7 and 8. In the scenario with scale variations, the accuracy gap between Cozzolino'15 and its DIR version widens significantly. Also, strong performance degradation is generally observed for other dense algorithms. The introduction of DIR improves the recall ($\sim 37\%$ gain) and the precision ($\sim 8\%$ gain), thus exhibiting a higher F1 ($\sim 27\%$ gain). The more noticeable fact is that the average number of matches per image has nearly $\times 2$ in the DIR version. In addition, the proposed algorithm also demonstrates benefits compared to the similar dense strategies.

Comprehensive Robustness. Finally, we evaluate the robustness under comprehensive attacks, involving six signal corruptions of whole images and four geometric transformations of manipulated regions, on the FAU benchmark. The comparison methods include Cozzolino'15 [48], Zandi'16 [49], and Li'19 [50]. To reveal the distribution properties of the forensic scores (especially the average and deviation nature), the corresponding box-plots are presented in Fig. 9. Note that we provide the detailed scores for each attack in Appendix F.

It is clear from Fig. 9 that, compared with Cozzolino'15 in precision, recall, and F1 metrics, the proposed DIR version exhibits lower score deviations while maintaining higher or similar average scores. As for Zandi'16 and Li'19, greater

performance gain is achieved by our method from both average and deviation perspectives. These common phenomena confirm that, with the introduction of our DIR, the stability of the copy-move forensic algorithm is significantly improved (especially for scaling) while not compromising the average-level of accuracy.

Such experimental evidence supports our theoretical expectation on geometric transformation robustness, verifying the usefulness of DIR for passive forensic scenarios. We would like to make a comment that the above application of DIR is naive; more careful design for scaling invariance can be achieved by incorporating the dense scale selection.

6.4 Active Forensics: Perceptual Hashing

Perceptual hashing is a well-known active forensic framework for multimedia content authentication. The main idea relies on deriving a compact hash sequence as digital abstract of the image content. With this idea in mind, the image representation should be robust to content-preserving operations and discriminative to visually distinct contents; as for efficiency, the hash sequence is expected to be compact enough. Note that such requirements are also consistent with the discussion in Section 1.1.

Algorithm Design. Currently, the state-of-the-art perceptual hashing algorithms generally follow a general framework that combines sparse and dense features, where sparse features (e.g., SIFT) for geometric correction and dense features (e.g., dense DCT) for forgery localization [61], [62]. Despite their popularity, the DCT-like dense features still suffer from flaws in both representation capability and hash compactness. The forensic results are quite unstable under certain signal corruptions, such as blurring and compression. Motivated by this, we replace such dense DCT by our DIR for improving its signal corruption robustness while shortening the hash sequence. Here, we consider an 8-pixel sampling interval for (u, v) over the image grid, which is a common setting for compactness. The order constraint is set as $K = 3$ in (11) with ∞ -norm, resulting in 16-dimensional features for each position. The scale w only increases from 8 to 12 with 3 samples due to the presence of geometric correction, and such feature vectors are then average-pooled together over the scales. In the implementation, the threshold for hash distance comparison is important. For a fair experiment, this threshold is determined by the adaptive Otsu's method [63] for all comparison methods.

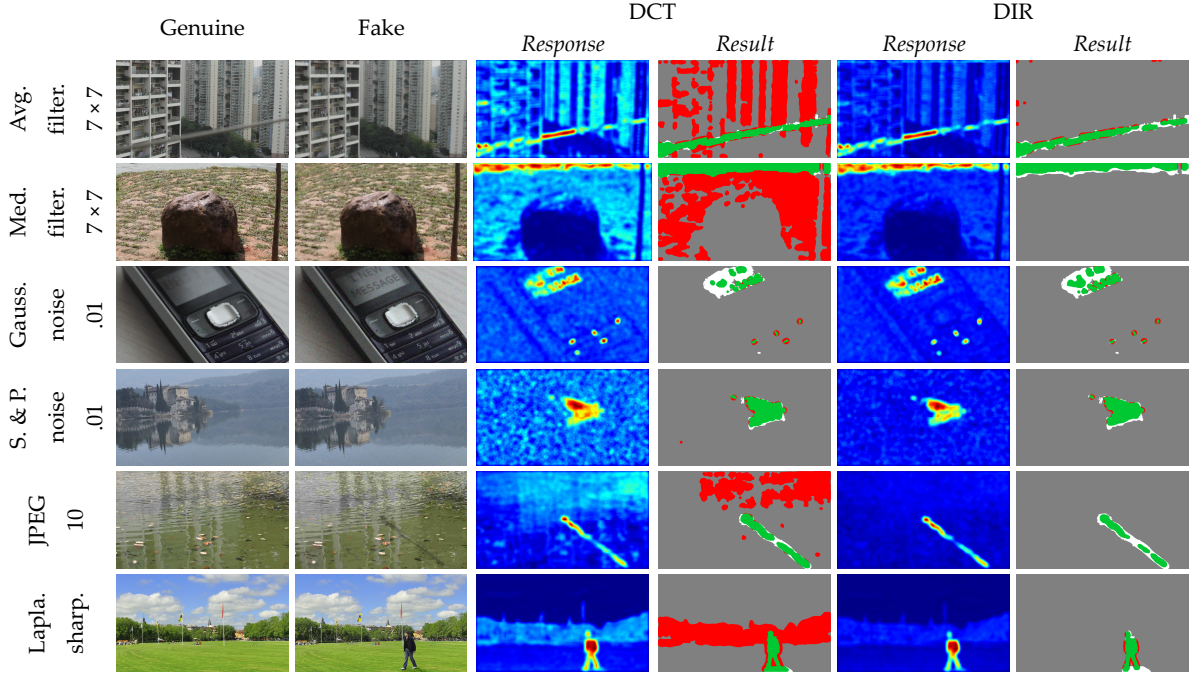


Fig. 10. Some samples of the content authentication by different perceptual hashing methods.

TABLE 9
Precision, Recall, and F1 Scores (%) for Different Perceptual Hashing Algorithms in Hash Robustness Experiment.

Method	DCT [61], [62]			DIR		
	Precision	Recall	F1	Precision	Recall	F1
Clean	87.24	77.7	79.85	89.03	71.97	76.95
Gaussian noise 0.01	87.02	74.81	77.8	88.63	70.13	75.46
Gaussian noise 0.02	86.74	74.2	77.21	87.8	70.05	74.68
Salt and pepper noise 0.01	86.38	77.83	79.52	88.61	71.88	76.68
Salt and pepper noise 0.02	85.73	77.66	78.97	88.29	71.65	76.35
Average filtering 7×7	55.43	75.33	53.43	78.8	68.09	68.06
Average filtering 14×14	37.87	74.02	38.45	52.89	65.55	50.18
Gaussian filtering 7×7	59.49	75.48	56.72	80.43	68.73	69.41
Gaussian filtering 14×14	49.86	74.97	49.03	71.47	68.17	63.35
Median filtering 7×7	60.61	74.91	57.84	82.43	68.51	70.17
Median filtering 14×14	44.72	69.49	43.13	66.83	61.51	56.72
JPEG compression 10	84.22	76.08	76.28	88.61	71.02	76.15
JPEG compression 5	77.12	75.57	69.86	87.55	70.88	75.29
Laplacian sharpening	55.36	80.17	54.55	78.76	72.79	69.64
Average \uparrow	68.41	75.59	63.76	80.72	69.35	69.94
Standard deviation \downarrow	17.53	2.36	14.32	10.26	2.89	7.88

Hash Robustness. A quantitative comparison is performed on forensic benchmark RTD [64], which involves not only the inpainting, splicing, and copy-move manipulations, but also other subtle changes by hand. Due to the inherent nature of perceptual hashing, the forensics on such clean images is effortless. Therefore, we mainly consider the more challenging tasks under global signal corruptions. The comparison methods include the basic SIFT-DCT framework and our SIFT-DIR version. Note that such basic framework is involved with several state-of-the-art perceptual hashing algorithms, e.g., Zhang’20 [61], Wang’15 [62], Hao’21 [65], and Biswas’20 [66].

Fig. 10 shows some samples of the response map and detection result, while the precision, recall, and F1 scores are summarized in Table 9. One can note that, even in the

forensic scenarios employing geometric correction, the DIR still exhibits certain advantages over common descriptors. As shown in the response map, under some operations, it is challenging for DCT to be robust while maintaining discriminability. The DCT typically exhibits false positives in regions with content-preserving operations. In general, compared with such state-of-the-art framework, our DIR version provides comparable detection accuracy for clean and noisy images, while generally being more robust to blurring, compression, and sharpening. The average and standard deviation of the scores also demonstrate the advantage of DIR, especially in terms of precision (w.r.t. false positives). The DIR exhibits better overall localization accuracy (for average) and is more stable under different attacks (for standard deviation).

Hash Compactness. Turning to compactness, the dimensions of DCT and DIR features on each position are 32 and 16, respectively. This means a nearly 50% saving in terms of storage and transmission costs.

Such experimental evidence supports our theoretical expectation on signal perturbation robustness and representation compactness, verifying the usefulness of DIR for active forensic scenarios.

7 CONCLUSION

The main goal of this paper is to provide a principled design of image representation for forensic tasks. We name this complete pipeline “Dense Invariant Representation”, meaning a local representation with covariance for position, orientation and scale variations.

The key ingredients of our work are as follows. i) At the theoretical level, the global definition of classical orthogonal moments is extended to the local with scale space. With such generic definition, the deformation stability backed with “invariance-equivariance-covariance” framework is explicitly analyzed (Section 3). ii) At the implementation level, the fine-tuned discrete computation strategy is designed, which is characterized by low integral/numerical error, constant complexity, and generality for arbitrary basis functions (Section 4). iii) At the application level, the above ideas are fully validated in two vision tasks (dense pattern detection and matching) and two forensic tasks (copy-move forgery detection and perceptual hashing). In general, our method gives state-of-the-art accuracy and efficiency performance, proving its promise in small-scale robust vision problems (Section 6).

The limitations for current approach includes scaling invariance construction and space complexity. We note that the common pooling operations only provide a limited tolerance for scaling. A more elegant construction is in our plan, relying on a suitable fusion of DIR and some recent advances of dense scale selection. Additionally, our fast implementation based on the convolution theorem may increase the space complexity. As a long-term research path, such problem is expected to be mitigated by hardware-oriented optimization and more careful adaptive time-space trade-offs.

ACKNOWLEDGMENTS

This work was supported in part by the National Key R&D Program of China under Grant 2019YFB1406500, in part by the Postgraduate Research & Practice Innovation Program of Jiangsu Province under Grant KYCX22_0383, in part by the National Natural Science Foundation of China under Grants 62072237, 61971476, and U2001202, in part by Guangxi Key Laboratory of Trusted Software under Grant KX202027, in part by Basic Research Program of Jiangsu Province under Grant BK20201290, in part by Macau Science and Technology Development Fund under Grants SKLIOTSC-2021-2023 and 0072/2020/AMJ, and in part by Research Committee at University of Macau under Grant MYRG2020-00101-FST.

APPENDIX A

COMMON DEFINITIONS OF RADIAL BASIS FUNCTIONS IN (4)

In Table A1, we include some common definitions of radial basis functions for (4), with orthogonality and completeness [32]. The main discussion in this paper is generic to all such definitions. Also, we implement all the DIR extensions for the methods listed in Table A1. The code is available online at <https://github.com/ShurenQi/DIR>.

APPENDIX B

PROOF OF (8): EQUIVARIANCE TO TRANSLATION

By plugging f_T into (7), one can check that image translation operation only affects the representation parameters (u, v) , as follows:

$$\begin{aligned} & \langle f_T(x, y), V_{nm}^{uvw}(x, y) \rangle \\ &= \iint_D R_n^* \left(\frac{1}{w} \sqrt{(x-u)^2 + (y-v)^2} \right) A_m^* \left(\arctan \left(\frac{y-v}{x-u} \right) \right) \\ & \quad \times f(x + \Delta x, y + \Delta y) dx dy \\ &= \iint_D R_n^* \left(\frac{1}{w} \sqrt{(x - \Delta x - u)^2 + (y - \Delta y - v)^2} \right) \\ & \quad \times A_m^* \left(\arctan \left(\frac{y - \Delta y - v}{x - \Delta x - u} \right) \right) f(x, y) dx dy \\ &= \iint_D R_n^* \left(\frac{1}{w} \sqrt{(x - (u + \Delta x))^2 + (y - (v + \Delta y))^2} \right) \\ & \quad \times A_m^* \left(\arctan \left(\frac{y - (v + \Delta y)}{x - (u + \Delta x)} \right) \right) f(x, y) dx dy \\ &= \langle f(x, y), V_{nm}^{(u+\Delta x)(v+\Delta y)w}(x, y) \rangle, \end{aligned}$$

where the same offset $(\Delta x, \Delta y)$ also appears in representation $\langle f(x, y), V_{nm}^{(u+\Delta x)(v+\Delta y)w}(x, y) \rangle$, implying the equivariance w.r.t. translation.

APPENDIX C

PROOF OF (9): COVARIANCE AND INVARIANCE TO ROTATION

By plugging f_R into (7) with $(u, v) = (0, 0)$ and writing down in polar form, one can check that image rotation operation only affects the phase of representation, as follows:

$$\begin{aligned} & \langle f_R(r, \theta), V_{nm}^{uvw}(r', \theta') \rangle \\ &= \iint_D R_n^*(r') A_m^*(\theta') f(r, \theta + \phi) r' dr' d\theta' \\ &= \frac{1}{w^2} \int_0^{2\pi} \int_0^w R_n^* \left(\frac{r}{w} \right) A_m^*(\theta) f(r, \theta + \phi) r dr d\theta \\ &= \frac{1}{w^2} \int_0^{2\pi} \int_0^w R_n^* \left(\frac{r}{w} \right) A_m^*(\theta - \phi) f(r, \theta) r dr d\theta \\ &= \frac{1}{w^2} \int_0^{2\pi} \int_0^w R_n^* \left(\frac{r}{w} \right) A_m^*(\theta) A_m^*(-\phi) f(r, \theta) r dr d\theta \\ &= \frac{1}{w^2} \langle f(r, \theta), V_{nm}^{uvw}(r', \theta') \rangle w^2 A_m^*(-\phi) \\ &= \langle f(r, \theta), V_{nm}^{uvw}(r', \theta') \rangle A_m^*(-\phi), \end{aligned}$$

where $A_m^*(\theta - \phi) = A_m^*(\theta)A_m^*(-\phi)$ is true as $A_m(\theta) = \exp(jm\theta)$; the same angle ϕ also appears in phase of the representation $\langle f(r, \theta), V_{nm}^{uvw}(r', \theta') \rangle A_m^*(-\phi)$, implying the covariance w.r.t. rotation. Therefore, the magnitude-only strategy is able to derive the rotation-invariant features: $|\langle f_R(r, \theta), V_{nm}^{uvw}(r', \theta') \rangle| = |\langle f(r, \theta), V_{nm}^{uvw}(r', \theta') \rangle|$.

APPENDIX D

PROOF OF (10): COVARIANCE TO SCALING

By plugging f_S into (7) with $(u, v) = (0, 0)$, one can check that image scaling operation only affects the representation parameters w , as follows:

$$\begin{aligned}
& \langle f_S(x, y), V_{nm}^{uvw}(x, y) \rangle \\
&= \iint_{x^2+y^2 \leq w^2} R_n^*\left(\frac{1}{w}\sqrt{x^2+y^2}\right) A_m^*\left(\arctan\left(\frac{y}{x}\right)\right) f(sx, sy) dx dy \\
&= \iint_{x^2+y^2 \leq w^2 s^2} R_n^*\left(\frac{1}{w}\sqrt{\left(\frac{x}{s}\right)^2 + \left(\frac{y}{s}\right)^2}\right) A_m^*\left(\arctan\left(\frac{y}{x}\right)\right) \\
&\quad \times f(x, y) d\frac{x}{s} d\frac{y}{s} \\
&= \frac{1}{s^2} \iint_{x^2+y^2 \leq w^2 s^2} R_n^*\left(\frac{1}{ws}\sqrt{x^2+y^2}\right) A_m^*\left(\arctan\left(\frac{y}{x}\right)\right) \\
&\quad \times f(x, y) dx dy \\
&= \frac{1}{s^2} \langle f(x, y), V_{nm}^{uv(ws)}(x, y) \rangle s^2 \\
&= \langle f(x, y), V_{nm}^{uv(ws)}(x, y) \rangle,
\end{aligned}$$

where it should be reminded that the area for integration changes with the factor of s^2 ; the same factor s also appears in the representation $\langle f(x, y), V_{nm}^{uv(ws)}(x, y) \rangle$, implying the covariance w.r.t. scaling.

APPENDIX E

AN ANALYSIS OF THE ROBUSTNESS TO SIGNAL CORRUPTION

In fact, our DIR exhibits a certain robustness to some signal corruptions, e.g., additive noise and low-pass filtering/blurring. This is mainly inherited from the classical orthogonal moment theory, i.e., the low-order moments are less affected by the signal corruptions that act on the high-frequency components of the image. Here, we give a formal analysis of this property from a frequency-domain perspective.

Without loss of generality, the signal corruptions of the modern imaging process can be modeled as:

$$f = p \otimes s + n,$$

where f is the captured image, s is the ideal scene (i.e., ground-truth image), p denotes the Point-Spread Function (PSF) w.r.t. image blurring, n denotes the additive noise, and \otimes denotes the spatial convolution operation. Here, we assume that PSF p is a common low-pass filter, e.g., a Gaussian kernel function, and n is random white noise with limited intensity.

Considering the DIR of f with dense sampling of (u, v) over the discrete grid, it can be written as follows:

$$\langle f, V_{nm}^{uvw} \rangle = f(i, j) \otimes (H_{nm}^w(i, j))^T.$$

By plugging $f = p \otimes s + n$ into this equation, we can derive an expanded version as follows:

$$\begin{aligned}
\langle f, V_{nm}^{uvw} \rangle &= f \otimes (H_{nm}^w)^T \\
&= (p \otimes s + n) \otimes (H_{nm}^w)^T \\
&= p \otimes s \otimes (H_{nm}^w)^T + n \otimes (H_{nm}^w)^T.
\end{aligned}$$

Let us analyze this result from a frequency-domain perspective:

$$\begin{aligned}
& \mathcal{F}(p \otimes s \otimes (H_{nm}^w)^T + n \otimes (H_{nm}^w)^T) \\
&= \mathcal{F}(p)\mathcal{F}(s)\mathcal{F}((H_{nm}^w)^T) + \mathcal{F}(n)\mathcal{F}((H_{nm}^w)^T),
\end{aligned}$$

where the equation holds due to the convolution theorem.

Mathematically, when the order (n, m) is small, e.g., under a l -norm-based constraint condition $\{(n, m) : \|(n, m)\|_l \leq K\}$ with small positive constant K , the frequency distribution of $\mathcal{F}((H_{nm}^w)^T)$ is sparse [22], and $\mathcal{F}((H_{nm}^w)^T)$ is basically composed of low-frequency coefficients [22].

With this property and the low-pass nature of p , we have $\mathcal{F}(p)\mathcal{F}((H_{nm}^w)^T) \simeq \mathcal{F}((H_{nm}^w)^T)$.

As for $\mathcal{F}(n)$, it is random, specifically the phase is completely random and the amplitude is random but constrained by the noise intensity. Hence, for the frequency (ξ_x, ξ_y) such that $\mathcal{F}((H_{nm}^w)^T)(\xi_x, \xi_y) \simeq 0$, we have $\mathcal{F}(n)(\xi_x, \xi_y)\mathcal{F}((H_{nm}^w)^T)(\xi_x, \xi_y) \simeq 0$. Note that due to the sparsity of $\mathcal{F}((H_{nm}^w)^T)$, such (ξ_x, ξ_y) satisfying $\mathcal{F}((H_{nm}^w)^T)(\xi_x, \xi_y) \simeq 0$ is dominant in the frequency domain. As for other frequency (ξ_x, ξ_y) , the effect of the noise for the frequency coefficients is also mitigated, since the amplitude of $\mathcal{F}((H_{nm}^w)^T)$ is quite small relative to the amplitude of $\mathcal{F}(n)$.

Based on the above assumptions, it is reasonable to derive the following approximation:

$$\begin{aligned}
& \mathcal{F}(\langle f, V_{nm}^{uvw} \rangle) \\
&= \mathcal{F}(p)\mathcal{F}(s)\mathcal{F}((H_{nm}^w)^T) + \mathcal{F}(n)\mathcal{F}((H_{nm}^w)^T) \\
&\simeq \mathcal{F}(s)\mathcal{F}((H_{nm}^w)^T) \\
&= \mathcal{F}(\langle s, V_{nm}^{uvw} \rangle),
\end{aligned}$$

i.e., the DIR-based low-order moments for f and s , $\langle f, V_{nm}^{uvw} \rangle$ and $\langle s, V_{nm}^{uvw} \rangle$, are approximation in frequency domain, meaning certain robustness in many practical tasks for the signal corruptions like blurring and noise.

APPENDIX F

COMPLEMENTARY EXPERIMENTS FOR COPY-MOVE FORGERY DETECTION

We have conducted several complementary experiments for copy-move forgery detection w.r.t. different parameter settings, different basis function definitions, and systematic image attacks.

In Table A2, the copy-move forensic accuracy on the FAU benchmark under different DIR settings is provided, involving three different values of K/w and five different radial basis function definitions. As can be seen from this

TABLE A1
Definitions of Radial Basis Functions of Unit Disk-based Orthogonal Moments

Method	Radial Basis Function
ZM [67]	$R_{nm}^{(ZM)}(r) = \sqrt{\frac{n+1}{\pi}} \sum_{k=0}^{\lfloor \frac{n- m }{2} \rfloor} \frac{(-1)^k (n-k)! r^{n-2k}}{k! (\frac{n+ m }{2} - k)! (\frac{n- m }{2} - k)!}$
PZM [68]	$R_{nm}^{(PZM)}(r) = \sqrt{\frac{n+1}{\pi}} \sum_{k=0}^{n- m } \frac{(-1)^k (2n+1-k)! r^{n-k}}{k! (n+ m +1-k)! (n- m -k)!}$
OFMM [69]	$R_n^{(OFMM)}(r) = \sqrt{\frac{n+1}{\pi}} \sum_{k=0}^n \frac{(-1)^{n+k} (n+k+1)! r^k}{k! (n-k)! (k+1)!}$
CHFM [70]	$R_n^{(CHFM)}(r) = \frac{2}{\pi} \left(\frac{1-r}{r}\right)^{\frac{1}{4}} \sum_{k=0}^{\lfloor \frac{n}{2} \rfloor} \frac{(-1)^k (n-k)! (4r-2)^{n-2k}}{k! (n-2k)!}$
PJFM [71]	$R_n^{(PJFM)}(r) = \sqrt{\frac{(n+2)(r-r^2)}{\pi(n+3)(n+1)}} \sum_{k=0}^n \frac{(-1)^{n+k} (n+k+3)! r^k}{k! (n-k)! (k+2)!}$
JFM [72]	$R_n^{(JFM)}(p, q, r) = \sqrt{\frac{r^{q-2} (1-r)^{p-q} (p+2n) \cdot \Gamma(q+n) \cdot n!}{2\pi \Gamma(p+n) \cdot \Gamma(p-q+n+1)}} \sum_{k=0}^n \frac{(-1)^k \Gamma(p+n+k) r^k}{k! (n-k)! \Gamma(q+k)}$
RHFM [73]	$R_n^{(RHFM)}(r) = \begin{cases} \frac{1}{\sqrt{2\pi r}} & n = 0 \\ \sqrt{\frac{1}{\pi r}} \sin(\pi(n+1)r) & n > 0 \text{ \& } n \text{ odd} \\ \sqrt{\frac{1}{\pi r}} \cos(\pi n r) & n > 0 \text{ \& } n \text{ even} \end{cases}$
EFM [74]	$R_n^{(EFM)}(r) = \frac{1}{\sqrt{2\pi r}} \exp(j2n\pi r)$
PCET [42]	$R_n^{(PCET)}(r) = \frac{1}{\sqrt{\pi}} \exp(j2n\pi r^2)$
PCT [42]	$R_n^{(PCT)}(r) = \begin{cases} \frac{1}{\sqrt{\pi}} & n = 0 \\ \sqrt{\frac{2}{\pi}} \cos(n\pi r^2) & n > 0 \end{cases}$
PST [42]	$R_n^{(PST)}(r) = \sqrt{\frac{2}{\pi}} \sin(n\pi r^2)$
BFM [75]	$R_n^{(BFM)}(r) = \frac{1}{\sqrt{\pi} J_{v+1}(\lambda_n)} J_v(\lambda_n r), J_v(x) = \sum_{k=0}^{\infty} \frac{(-1)^k}{k! \Gamma(v+k+1)} \left(\frac{x}{2}\right)^{v+2k}$
GRHFM [76]	$R_n^{(GRHFM)}(\alpha, r) = \begin{cases} \sqrt{\frac{\alpha r^{\alpha-2}}{2\pi}} & n = 0 \\ \sqrt{\frac{\alpha r^{\alpha-2}}{\pi}} \sin(\pi(n+1)r^\alpha) & n > 0 \text{ \& } n \text{ odd} \\ \sqrt{\frac{\alpha r^{\alpha-2}}{\pi}} \cos(\pi n r^\alpha) & n > 0 \text{ \& } n \text{ even} \end{cases}$
GPCET [76]	$R_n^{(GPCET)}(\alpha, r) = \sqrt{\frac{\alpha r^{\alpha-2}}{2\pi}} \exp(j2n\pi r^\alpha)$
GPCT [76]	$R_n^{(GPCT)}(\alpha, r) = \begin{cases} \sqrt{\frac{\alpha r^{\alpha-2}}{2\pi}} & n = 0 \\ \sqrt{\frac{\alpha r^{\alpha-2}}{\pi}} \cos(\pi n r^\alpha) & n > 0 \end{cases}$
GPST [76]	$R_n^{(GPST)}(\alpha, r) = \sqrt{\frac{\alpha r^{\alpha-2}}{\pi}} \sin(\pi n r^\alpha)$
FJFM [40]	$R_n^{(FJFM)}(\alpha, p, q, r) = \sqrt{\frac{\alpha r^{\alpha q-2} (1-r^\alpha)^{p-q} (p+2n) \cdot \Gamma(q+n) \cdot n!}{2\pi \Gamma(p+n) \cdot \Gamma(p-q+n+1)}} \sum_{k=0}^n \frac{(-1)^k \Gamma(p+n+k) r^{\alpha k}}{k! (n-k)! \Gamma(q+k)}$

TABLE A2
Precision, Recall, and F1 Scores (%) for Different DIR Settings on the FAU Copy-Move Forensic Benchmark.

Parameter	Setting	#1	#2	#3	#4	#5	#6	#7	#8	#9
K	3	✓			✓	✓	✓	✓	✓	✓
	1		✓							
	5			✓						
w	8 ~ 32	✓	✓	✓			✓	✓	✓	✓
	10 ~ 20				✓					
	6 ~ 36					✓				
R	PCT	✓	✓	✓	✓	✓				
	PCET						✓			
	ZM							✓		
	OFMM								✓	
	BFM									✓
Scores	Precision	96.75	96.75	96.75	97.95	96.94	96.50	96.65	96.68	92.80
	Recall	92.44	89.38	92.67	93.69	91.25	91.00	90.43	92.18	89.58
	F1	93.62	92.41	93.67	95.18	93.09	92.42	92.14	93.41	90.39

table, the forensic localization scores fluctuate slightly (in a reasonable interval) as the changes of settings, and the

setting used in the paper (i.e., setting #1) is not a special setting with highest scores. These phenomena imply that

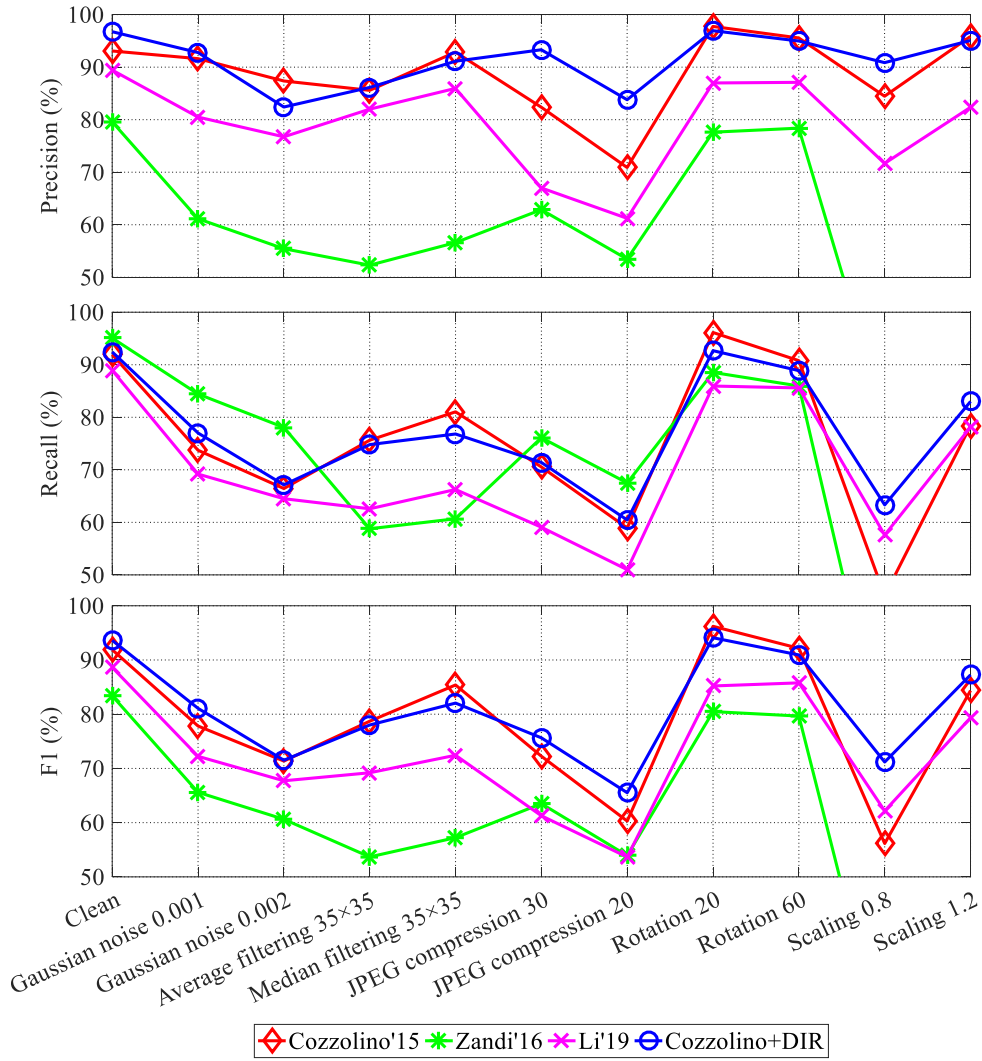


Fig. A1. Precision, recall, and F1 curves by different copy-move forgery detection methods in comprehensive robustness experiment.

TABLE A3
Precision, Recall, and F1 Scores (%) for Different DIR Settings in Hash Robustness Experiment.

Parameter	Setting	#1	#2	#3	#4	#5	#6	#7	#8	#9
K	3	✓			✓	✓	✓	✓	✓	✓
	1		✓							
	5			✓						
w	8 ~ 12	✓	✓	✓			✓	✓	✓	✓
	10				✓					
	6 ~ 16					✓				
R	PCT	✓	✓	✓	✓	✓				
	PCET						✓			
	ZM							✓		
	OFMM								✓	
	BFM									✓
Scores	Precision	80.72	84.11	79.19	79.34	81.29	82.13	85.68	80.63	80.44
	Recall	69.35	68.31	69.53	69.80	69.30	68.78	67.15	69.40	69.12
	F1	69.94	71.75	68.83	69.14	70.54	70.51	71.97	69.93	69.47

the DIR-based copy-move algorithm is not sensitive to such settings and the relevant experimental evaluations with

setting #1 in the paper are convincing (i.e., not an isolated case).

In Fig. A1, the forensic accuracy curves on the FAU

benchmark under different attacks are provided, involving six signal corruptions of whole images and four geometric transformations of manipulated regions. As can be seen from Fig. A1, for signal corruptions, the proposed method exhibits very similar scores to Cozzolino'15, which is in line with expectations since Cozzolino'15 is actually based on a single-scale DIR, while the other competing methods (i.e., Zandi'16 and Li'19) exhibit lower scores under most attacks. For geometric transformations, Cozzolino'15 and its DIR version are also systematically more robust than other competing methods. Rotation attack does not lead to significant differences between the Cozzolino'15 and its DIR version, implying that the rotation invariance is maintained in our multi-scale DIR framework. Under the scaling, especially with scale factor 0.8, the proposed DIR version significantly outperforms Cozzolino'15 (also Zandi'16 and Li'19), which is consistent with our original intention of introducing multi-scale DIR.

APPENDIX G COMPLEMENTARY EXPERIMENTS FOR PERCEPTUAL HASHING

We have also conducted several complementary experiments for perceptual hashing w.r.t. different parameter settings and different basis function definitions.

In Table A3, the perceptual-hashing based forensic accuracy on the RTD benchmark under different DIR settings is provided, involving three different values of K/w and five different radial basis function definitions. Note that the scores in Table A4 are averaged over all the attacks listed in Table 9. Here, a very similar conclusion can be observed, i.e., the DIR-based perceptual hashing algorithm is not sensitive to such settings and the relevant experimental evaluations with setting #1 in the paper are convincing.

APPENDIX H A REMARK ON THE FORENSIC APPLICABILITY

Following the taxonomy of DARPA, digital media forensic should look for digital integrity, physical integrity, and semantic integrity [1]. From the representation perspective, such digital, physical, and semantic integrities mainly correspond to the high-order, mid-order, and low-order statistics of the image, respectively [1]. As an image orthogonal analysis tool, our DIR has the beneficial property of information preservation (the whole image information is preserved in orthogonal moment space) [22], i.e., the full-order statistics provided by DIR allow a uniform discriminability for broad forensic scenarios.

We would like to clearly state that the forensic tasks listed in this paper, i.e., copy-move detection and perceptual hashing by semantic comparison, both belong to semantic-level algorithm. However, the discriminability provided by DIR is in fact not limited to such forensic tasks. We have initially explored the application of DIR to digital-level forensic tasks, i.e., splicing and inpainting detection by exposing digital artifacts, based on the coefficient distribution (kurtosis) modeling.

REFERENCES

- [1] L. Verdoliva, "Media forensics and deepfakes: An overview," *IEEE J. Sel. Top. Signal Process.*, vol. 14, no. 5, pp. 910–932, 2020.
- [2] H. A. Simon, *The sciences of the artificial*. MIT press, 2019.
- [3] J. Gu, Z. Wang, J. Kuen, L. Ma, A. Shahroudy, B. Shuai, T. Liu, X. Wang, G. Wang, J. Cai *et al.*, "Recent advances in convolutional neural networks," *Pattern Recognit.*, vol. 77, pp. 354–377, 2018.
- [4] X. Zhang, C. Liu, and C. Y. Suen, "Towards robust pattern recognition: A review," *Proc. IEEE*, vol. 108, no. 6, pp. 894–922, 2020.
- [5] Y. Pei, Y. Huang, Q. Zou, X. Zhang, and S. Wang, "Effects of image degradation and degradation removal to CNN-based image classification," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 43, no. 4, pp. 1239–1253, 2021.
- [6] X. Yuan, P. He, Q. Zhu, and X. Li, "Adversarial examples: Attacks and defenses for deep learning," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 30, no. 9, pp. 2805–2824, 2019.
- [7] Y. Zhang, P. Tiño, A. Leonardis, and K. Tang, "A survey on neural network interpretability," *IEEE Trans. Emerg. Top. Comput. Intell.*, vol. 5, no. 5, pp. 726–742, 2021.
- [8] T. Wiatowski and H. Bölcskei, "A mathematical theory of deep convolutional neural networks for feature extraction," *IEEE Trans. Inf. Theory*, vol. 64, no. 3, pp. 1845–1866, 2017.
- [9] M. M. Bronstein, J. Bruna, T. Cohen, and P. Veličković, "Geometric deep learning: Grids, groups, graphs, geodesics, and gauges," *arXiv preprint arXiv:2104.13478*, 2021.
- [10] V. Balntas, K. Lenc, A. Vedaldi, T. Tuytelaars, J. Matas, and K. Mikolajczyk, "H-patches: A benchmark and evaluation of hand-crafted and learned local descriptors," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 42, no. 11, pp. 2825–2841, 2019.
- [11] A. Iscen, G. Toliás, P.-H. Gosselin, and H. Jégou, "A comparison of dense region detectors for image search and fine-grained classification," *IEEE Trans. Image Process.*, vol. 24, no. 8, pp. 2369–2381, 2015.
- [12] Y. Li, J. Zhou, and A. Cheng, "SIFT keypoint removal via directed graph construction for color images," *IEEE Trans. Inf. Forensic Secur.*, vol. 12, no. 12, pp. 2971–2985, 2017.
- [13] E. Tola, V. Lepetit, and P. Fua, "DAISY: An efficient dense descriptor applied to wide-baseline stereo," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 32, no. 5, pp. 815–830, 2009.
- [14] T. Tuytelaars, "Dense interest points," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2010, pp. 2281–2288.
- [15] W. Ouyang and W.-K. Cham, "Fast algorithm for Walsh Hadamard transform on sliding windows," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 32, no. 1, pp. 165–171, 2009.
- [16] W. Ouyang, T. Zhao, W.-K. Cham, and L. Wei, "Fast full-search-equivalent pattern matching using asymmetric Haar wavelet packets," *IEEE Trans. Circuits Syst. Video Technol.*, vol. 28, no. 4, pp. 819–833, 2016.
- [17] W. Ouyang, F. Tombari, S. Mattoccia, L. Di Stefano, and W.-K. Cham, "Performance evaluation of full search equivalent pattern matching algorithms," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 34, no. 1, pp. 127–143, 2011.
- [18] S. Kim, D. Min, B. Ham, M. N. Do, and K. Sohn, "DASC: Robust dense descriptor for multi-modal and multi-spectral correspondence estimation," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 39, no. 9, pp. 1712–1729, 2016.
- [19] I. Kokkinos and A. Yuille, "Scale invariance without scale selection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2008, pp. 1–8.
- [20] B. Chen, G. Coatrieux, J. Wu, Z. Dong, J. L. Coatrieux, and H. Shu, "Fast computation of sliding discrete Tchebichef moments and its application in duplicated regions detection," *IEEE Trans. Signal Process.*, vol. 63, no. 20, pp. 5424–5436, 2015.
- [21] A. Bera, P. Klesk, and D. Sychel, "Constant-time calculation of Zernike moments for detection with rotational invariance," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 41, no. 3, pp. 537–551, 2019.
- [22] J. Flusser, B. Zitova, and T. Suk, *Moments and moment invariants in pattern recognition*. John Wiley & Sons, 2009.
- [23] V. Christlein, C. Riess, J. Jordan, C. Riess, and E. Angelopoulou, "An evaluation of popular copy-move forgery detection approaches," *IEEE Trans. Inf. Forensic Secur.*, vol. 7, no. 6, pp. 1841–1854, 2012.
- [24] Y. Liu, X. Zhu, X. Zhao, and Y. Cao, "Adversarial learning for constrained image splicing detection and localization based on atrous convolution," *IEEE Trans. Inf. Forensic Secur.*, vol. 14, no. 10, pp. 2551–2566, 2019.

- [25] H. Li, W. Luo, and J. Huang, "Localization of diffusion-based inpainting in digital images," *IEEE Trans. Inf. Forensic Secur.*, vol. 12, no. 12, pp. 3050–3064, 2017.
- [26] Y. Mirsky and W. Lee, "The creation and detection of deepfakes: A survey," *ACM Comput. Surv.*, vol. 54, no. 1, pp. 1–41, 2021.
- [27] I. Cox, M. Miller, J. Bloom, J. Fridrich, and T. Kalker, *Digital watermarking and steganography*. Morgan kaufmann, 2007.
- [28] L. Du, A. T. Ho, and R. Cong, "Perceptual hashing for image authentication: A survey," *Signal Process.-Image Commun.*, vol. 81, p. 115713, 2020.
- [29] D. Cozzolino and L. Verdoliva, "Noiseprint: A CNN-based camera model fingerprint," *IEEE Trans. Inf. Forensic Secur.*, vol. 15, pp. 144–159, 2019.
- [30] F. Matern, C. Riess, and M. Stamminger, "Gradient-based illumination description for image forgery detection," *IEEE Trans. Inf. Forensic Secur.*, vol. 15, pp. 1303–1317, 2019.
- [31] H. Qi, Q. Guo, F. Juefei-Xu, X. Xie, L. Ma, W. Feng, Y. Liu, and J. Zhao, "DeepRhythm: Exposing deepfakes with attentional visual heartbeat rhythms," in *Proc. ACM Int. Conf. Multimed.*, 2020, pp. 4318–4327.
- [32] S. Qi, Y. Zhang, C. Wang, J. Zhou, and X. Cao, "A survey of orthogonal moments for image representation: Theory, implementation, and evaluation," *ACM Comput. Surv.*, vol. 55, no. 1, 2021.
- [33] A. Bhatia and E. Wolf, "On the circle polynomials of Zernike and related orthogonal sets," in *Math. Proc. Camb. Philos. Soc.*, vol. 50, no. 1. Cambridge University Press, 1954, pp. 40–48.
- [34] K. Lenc and A. Vedaldi, "Understanding image representations by measuring their equivariance and equivalence," *Int. J. Comput. Vis.*, vol. 127, no. 5, pp. 456–476, 2019.
- [35] D. Marcos, M. Volpi, N. Komodakis, and D. Tuia, "Rotation equivariant vector field networks," in *Proc. IEEE Int. Conf. Comput. Vis.*, 2017, pp. 5048–5057.
- [36] J. Flusser, "On the independence of rotation moment invariants," *Pattern Recognit.*, vol. 33, no. 9, pp. 1405–1410, 2000.
- [37] T. Lindeberg, "Dense scale selection over space, time, and space-time," *SIAM J. Imaging Sci.*, vol. 11, no. 1, pp. 407–441, 2018.
- [38] M. Nwali and S. Liao, "A new fast algorithm to compute continuous moments defined in a rectangular region," *Pattern Recognit.*, vol. 89, pp. 151–160, 2019.
- [39] C. Singh and R. Upneja, "Accurate calculation of high order pseudo-Zernike moments and their numerical stability," *Digit. Signal Prog.*, vol. 27, pp. 95–106, 2014.
- [40] H. Yang, S. Qi, J. Tian, P. Niu, and X. Wang, "Robust and discriminative image representation: Fractional-order Jacobi-Fourier moments," *Pattern Recognit.*, vol. 115, p. 107898, 2021.
- [41] V. Sze, Y.-H. Chen, T.-J. Yang, and J. S. Emer, "Efficient processing of deep neural networks: A tutorial and survey," *Proc. IEEE*, vol. 105, no. 12, pp. 2295–2329, 2017.
- [42] P.-T. Yap, X. Jiang, and A. C. Kot, "Two-dimensional polar harmonic transforms for invariant image representation," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 32, no. 7, pp. 1259–1270, 2010.
- [43] T. Zhao and T. Blu, "The Fourier-Argand representation: An optimal basis of steerable patterns," *IEEE Trans. Image Process.*, vol. 29, pp. 6357–6371, 2020.
- [44] H. Jegou, M. Douze, and C. Schmid, "Hamming embedding and weak geometric consistency for large scale image search," in *Proc. Eur. Conf. Comput. Vision*, 2008, pp. 304–317.
- [45] C. Barnes, E. Shechtman, A. Finkelstein, and D. B. Goldman, "PatchMatch: A randomized correspondence algorithm for structural image editing," *ACM Trans. Graph.*, vol. 28, no. 3, p. 24, 2009.
- [46] M. A. Fischler and R. C. Bolles, "Random sample consensus: A paradigm for model fitting with applications to image analysis and automated cartography," *Commun. ACM*, vol. 24, no. 6, pp. 381–395, 1981.
- [47] C. Barnes, E. Shechtman, D. B. Goldman, and A. Finkelstein, "The generalized PatchMatch correspondence algorithm," in *Proc. Eur. Conf. Comput. Vision*, 2010, pp. 29–43.
- [48] D. Cozzolino, G. Poggi, and L. Verdoliva, "Efficient dense-field copy-move forgery detection," *IEEE Trans. Inf. Forensic Secur.*, vol. 10, no. 11, pp. 2284–2297, 2015.
- [49] M. Zandi, A. Mahmoudi-Aznavah, and A. Talebpour, "Iterative copy-move forgery detection based on a new interest point detector," *IEEE Trans. Inf. Forensic Secur.*, vol. 11, no. 11, pp. 2499–2512, 2016.
- [50] Y. Li and J. Zhou, "Fast and effective image copy-move forgery detection via hierarchical feature point matching," *IEEE Trans. Inf. Forensic Secur.*, vol. 14, no. 5, pp. 1307–1322, 2019.
- [51] S.-J. Ryu, M. Kirchner, M.-J. Lee, and H.-K. Lee, "Rotation invariant localization of duplicated image regions based on Zernike moments," *IEEE Trans. Inf. Forensic Secur.*, vol. 8, no. 8, pp. 1355–1370, 2013.
- [52] Y. Li, "Image copy-move forgery detection based on polar cosine transform and approximate nearest neighbor searching," *Forensic Sci. Int.*, vol. 224, no. 1–3, pp. 59–67, 2013.
- [53] E. Silva, T. Carvalho, A. Ferreira, and A. Rocha, "Going deeper into copy-move forgery detection: Exploring image telltales via multi-scale analysis and voting processes," *J. Vis. Commun. Image Represent.*, vol. 29, pp. 16–32, 2015.
- [54] M. Emam, Q. Han, and X. Niu, "PCET based copy-move forgery detection in images under geometric transforms," *Multimedia Tools Appl.*, vol. 75, no. 18, pp. 11513–11527, 2016.
- [55] C.-M. Pun and J.-L. Chung, "A two-stage localization for copy-move forgery detection," *Inf. Sci.*, vol. 463, pp. 33–55, 2018.
- [56] X. Bi and C.-M. Pun, "Fast copy-move forgery detection using local bidirectional coherency error refinement," *Pattern Recognit.*, vol. 81, pp. 161–175, 2018.
- [57] Y. Wu, W. Abd-Almageed, and P. Natarajan, "Busternet: Detecting copy-move image forgery with source/target localization," in *Proc. Eur. Conf. Comput. Vision*, 2018, pp. 168–184.
- [58] J.-L. Zhong and C.-M. Pun, "An end-to-end dense-inceptionnet for image copy-move forgery detection," *IEEE Trans. Inf. Forensic Secur.*, vol. 15, pp. 2134–2146, 2020.
- [59] I. Amerini, L. Ballan, R. Caldelli, A. Del Bimbo, L. Del Tongo, and G. Serra, "Copy-move forgery detection and localization by means of robust clustering with J-Linkage," *Signal Process.-Image Commun.*, vol. 28, no. 6, pp. 659–669, 2013.
- [60] E. Ardizzone, A. Bruno, and G. Mazzola, "Copy-move forgery detection by matching triangles of keypoints," *IEEE Trans. Inf. Forensic Secur.*, vol. 10, no. 10, pp. 2084–2094, 2015.
- [61] Y. Zheng, Y. Cao, and C.-H. Chang, "A PUF-based data-device hash for tampered image detection and source camera identification," *IEEE Trans. Inf. Forensic Secur.*, vol. 15, pp. 620–634, 2020.
- [62] X. Wang, K. Pang, X. Zhou, Y. Zhou, L. Li, and J. Xue, "A visual model-based perceptual image hash for content authentication," *IEEE Trans. Inf. Forensic Secur.*, vol. 10, no. 7, pp. 1336–1349, 2015.
- [63] N. Otsu, "A threshold selection method from gray-level histograms," *IEEE Trans. Sys. Man Cyber.*, vol. 9, no. 1, pp. 62–66, 1979.
- [64] P. Korus and J. Huang, "Evaluation of random field models in multi-modal unsupervised tampering localization," in *Proc. IEEE Int. Workshop Inf. Forensic Secur.*, 2016, pp. 1–6.
- [65] Q. Hao, L. Luo, S. T. Jan, and G. Wang, "It's not what it looks like: Manipulating perceptual image hash based applications," in *Proc. ACM Conf. Comput. Commun. Secur.*, 2021, pp. 69–85.
- [66] R. Biswas, V. Gonzalez-Castro, E. Fidalgo, and E. Alegre, "Perceptual image hashing based on frequency dominant neighborhood structure applied to tor domains recognition," *Neurocomput.*, vol. 383, pp. 24–38, 2020.
- [67] M. R. Teague, "Image analysis via the general theory of moments," *J. Opt. Soc. Am. A*, vol. 70, no. 8, pp. 920–930, 1980.
- [68] C.-H. Teh and R. T. Chin, "On image analysis by the methods of moments," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 10, no. 4, pp. 496–513, 1988.
- [69] Y. Sheng and L. Shen, "Orthogonal Fourier-Mellin moments for invariant pattern recognition," *J. Opt. Soc. Am. A*, vol. 11, no. 6, pp. 1748–1757, 1994.
- [70] Z. Ping, R. Wu, and Y. Sheng, "Image description with Chebyshev-Fourier moments," *J. Opt. Soc. Am. A*, vol. 19, no. 9, pp. 1748–1754, 2002.
- [71] G. Amu, S. Hasi, X. Yang, and Z. Ping, "Image analysis by pseudo-Jacobi ($p=4$, $q=3$)-Fourier moments," *Applied Optics*, vol. 43, no. 10, pp. 2093–2101, 2004.
- [72] Z. Ping, H. Ren, J. Zou, Y. Sheng, and W. Bo, "Generic orthogonal moments: Jacobi-Fourier moments for invariant image description," *Pattern Recognit.*, vol. 40, no. 4, pp. 1245–1254, 2007.
- [73] H. Ren, Z. Ping, W. Bo, W. Wu, and Y. Sheng, "Multidistortion-invariant image recognition with radial harmonic Fourier moments," *J. Opt. Soc. Am. A*, vol. 20, no. 4, pp. 631–637, 2003.
- [74] H. Hu, Y. Zhang, C. Shao, and Q. Ju, "Orthogonal moments based on exponent functions: Exponent-Fourier moments," *Pattern Recognit.*, vol. 47, no. 8, pp. 2596–2606, 2014.
- [75] B. Xiao, J. Ma, and X. Wang, "Image analysis by Bessel-Fourier moments," *Pattern Recognit.*, vol. 43, no. 8, pp. 2620–2629, 2010.

- [76] T. V. Hoang and S. Tabbone, "Generic polar harmonic transforms for invariant image representation," *Image Vis. Comput.*, vol. 32, no. 8, pp. 497–509, 2014.



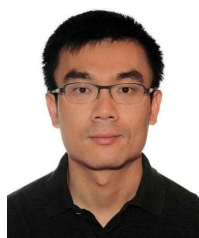
Shuren Qi received the B.A. and M.E. degrees from Liaoning Normal University, Dalian, China, in 2017 and 2020 respectively. He is currently pursuing the Ph.D. degree in computer science at Nanjing University of Aeronautics and Astronautics, Nanjing, China. His research interests include invariant feature extraction and visual signal representation with applications in robust pattern recognition and multimedia forensics/security.



security, artificial intelligence, and blockchain. Dr. Zhang is an Associate Editor of *Signal Processing and Information Sciences*.



Chao Wang received the B.S. and M.S. degrees from Liaoning Normal University, Dalian, China, in 2017 and 2020 respectively. She is currently pursuing the Ph.D. degree in computer science at Nanjing University of Aeronautics and Astronautics, Nanjing, China. Her research interests include trustworthy artificial intelligence, adversarial learning, and media forensics.



Jiantao Zhou (Senior Member, IEEE) received the B.E. degree from the Department of Electronic Engineering, Dalian University of Technology, in 2002, the M.E. degree from the Department of Radio Engineering, Southeast University, in 2005, and the Ph.D. degree from the Department of Electronic and Computer Engineering, The Hong Kong University of Science and Technology, in 2009. He held various research positions with the University of Illinois at Urbana-Champaign, The Hong Kong University of Science and Technology, and McMaster University. He is currently an Associate Professor with the Department of Computer and Information Science, Faculty of Science and Technology, University of Macau. He holds four granted U.S. patents and two granted Chinese patents. His research interests include multimedia security and forensics, and multimedia signal processing. He has coauthored three articles that received the Best Paper Award from the *IEEE Pacific-Rim Conference on Multimedia* in 2007, the Best Student Paper Award, and the Best Paper Runner Up Award from the *IEEE International Conference on Multimedia and Expo* in 2016 and 2020. He has been an Associate Editor of *IEEE Transactions on Image Processing* since 2019.



Xiaochun Cao (Senior Member, IEEE) received the B.E. and M.E. degrees in computer science from Beihang University, Beijing, China, and the Ph.D. degree in computer science from the University of Central Florida, Orlando, FL, USA. After graduation, he spent about three years at ObjectVideo Inc. as a Research Scientist. From 2008 to 2012, he was a Professor with Tianjin University, Tianjin, China. He is currently a Professor with School of Cyber Science and Technology, Shenzhen Campus, Sun Yat-sen University, Shenzhen, China. He is a Fellow of the IET. His dissertation was nominated for the University of Central Florida's university-level Outstanding Dissertation Award. In 2004 and 2010, he was the recipient of the Piero Zamperoni Best Student Paper Award at the *International Conference on Pattern Recognition*. He is on the Editorial Board of the *IEEE Transactions on Image Processing*, *IEEE Transactions on Multimedia*, and *IEEE Transactions on Circuits and Systems for Video Technology*.