

Implémentation du jeu Hex via la Programmation Orientée Objet et stratégies de décision en C++

1st Jair Anderson Vasquez Torres
Ingénieur Degree Programme STIC
ENSTA Paris
Paris, France
jair-anderson.vasquez@ensta.fr

2nd Santiago Florido Gomez
Ingénieur Degree Programme STIC
ENSTA Paris
Paris, France
santiago.florido@ensta-paris.fr

Abstract—Ce projet implémente en C++ un système complet pour le jeu Hex, conçu principalement comme un exercice de Programmation Orientée Objet. La solution organise le domaine du jeu au moyen de classes qui encapsulent le plateau, l'état, les coordonnées et les règles, permettant la génération de coups légaux et la vérification de victoire au moyen d'un parcours BFS sur les connexions des pions. Sur cette base s'intègrent des stratégies de décision (p. ex., Negamax avec hachage et table de transposition) et une évaluation de positions découpée du moteur, qui peut être heuristique ou s'appuyer sur un modèle neuronal de valeur intégré à l'exécutable (exportable en TorchScript), sans lier la conception à une architecture spécifique. De plus, une interface graphique (GUI) en SFML a été développée pour faciliter l'interaction, la visualisation du plateau et les tests du comportement des agents.

Index Terms—Hex game, C++, Object-Oriented Programming, Negamax, Transposition Table, TorchScript, SFML.

I. IMAGE FORMATS AND CONVOLUTIONS

A discrete two-dimensional convolution is a linear operation that produces an output image by sliding a small matrix, the kernel K of size $(2a+1) \times (2b+1)$, over every pixel of an input image I and computing a weighted sum of the neighborhood, as defined in Eq. (1):

$$(I * K)(y, x) = \sum_{j=-a}^a \sum_{i=-b}^b I(y + j, x + i) K(j, i). \quad (1)$$

This operation is the computational basis for a range of image-processing tasks such as smoothing, sharpening, edge detection, and feature extraction, all depending on the choice of kernel [1].

A. OpenCV and Matplotlib Functions

The provided code in `Convolutions.py` relies on several OpenCV and Matplotlib functions whose behavior needs to be understood to interpret the results correctly.

The image is loaded with `cv2.imread(path, 0)`, where the flag 0 forces grayscale loading as a single-channel 8-bit array, and the result is cast to `float64` via `np.float64()` so that subsequent convolution operations can produce values outside $[0, 255]$ without overflow. For copying the image, the code calls `cv2.copyMakeBorder(img, 0, 0, 0, 0, cv2.BORDER_REPLICATE)` with zero

padding on every side, which in practice just creates a deep copy; the `BORDER_REPLICATE` flag would replicate edge pixels if positive padding were specified. The optimized convolution is done via `cv2.filter2D(img, -1, kernel)`, where -1 means the output keeps the same depth as the input; internally, this delegates to OpenCV's C++ backend, which can use SIMD and multi-threading [2].

An important detail concerns the display conventions: `cv2.imshow` interprets integer arrays on $\{0, \dots, 255\}$ but floating-point arrays on $[0, 1]$, which is why the code divides by 255 when showing the `filter2D` result. On the other hand, `plt.imshow` normalizes the array range to the colormap by default, so to get a consistent grayscale display the code explicitly sets `vmin=0.0`, `vmax=255.0` and `cmap='gray'`.

B. Direct Computation vs OpenCV filter2D

The code compares two ways of computing the convolution of the grayscale image `FlowerGarden2.png` (240×360 pixels) with the 3×3 sharpening kernel. The first is a direct method that uses a nested Python loop over all interior pixels (y, x) with $1 \leq y \leq h-2$, $1 \leq x \leq w-2$, computing explicitly:

$$v = 5I(y, x) - I(y-1, x) - I(y, x-1) - I(y+1, x) - I(y, x+1), \quad (2)$$

followed by a clamp $I_{\text{out}}(y, x) = \min(\max(v, 0), 255)$. The second is a single call to `cv2.filter2D`, which does the same operation but through OpenCV's compiled C++ backend.

Given the nature of these two implementations, the hypothesis is straightforward: the direct Python loop should be dramatically slower than `filter2D`, because every pixel iteration goes through the Python interpreter with dynamic type checking, whereas the C++ path processes contiguous memory with compiled instructions and can leverage SIMD vectorization and multi-threading. We would also expect that, for a small image like this one, enabling or disabling multi-threading in OpenCV should not make a large difference, since the overhead of thread synchronization could offset the parallelism gains when the workload per thread is small.

TABLE I
MEDIAN EXECUTION TIME FOR THE CONVOLUTION OF A 240×360 GRayscale IMAGE WITH A 3×3 KERNEL.

Method	Median time	Runs
Direct (Python loop)	0.6987 s	5
filter2D (multi-thread)	0.329 ms	20
filter2D (1-thread)	0.348 ms	20

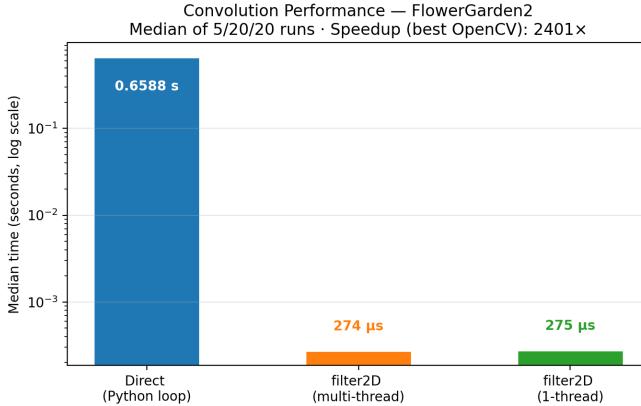


Fig. 1. Median execution time comparison (log scale) for the three convolution methods on a 240×360 image.

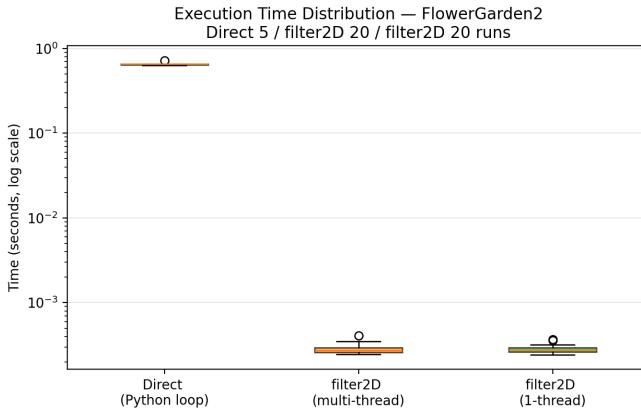


Fig. 2. Distribution of execution times across all benchmark runs for each method.

To test this, each method was benchmarked over multiple independent runs (5 for the slow direct method, 20 for the fast `filter2D` variants), and the median execution time was retained to mitigate warm-up and scheduling noise. A single-thread variant of `filter2D` was measured separately by calling `cv2.setNumThreads(1)` before execution.

As reported in Table I and illustrated in Fig. 1, the direct Python loop takes about 0.70 s while `filter2D` completes in roughly 0.33 ms, giving a speedup of approximately 2100 \times . This is consistent with the hypothesis: the per-pixel overhead of the Python interpreter is enormous compared to the compiled, vectorized C++ implementation. The difference between the multi-thread and single-thread variants of `filter2D` is

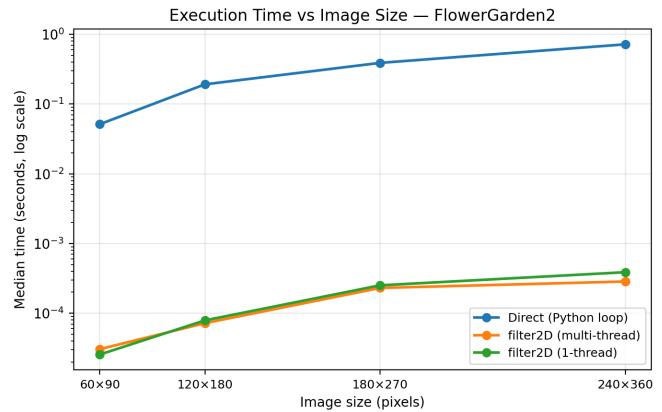


Fig. 3. Execution time as a function of image resolution (25 %, 50 %, 75 %, and 100 % of the original 240×360 image).

almost negligible (0.329 ms vs 0.348 ms), which also confirms our expectation that thread synchronization cost offsets any parallelism benefit at this image size. Moreover, Fig. 2 shows that the `filter2D` measurements have very low variance, reflecting the stability of a compiled implementation, while the direct method exhibits more dispersion—likely caused by Python’s garbage collector and OS scheduling interrupting the long-running loop.

To see how these results scale, the benchmark was repeated on resized versions of the image at 25 %, 50 %, 75 %, and 100 % of the original resolution. The hypothesis here is that the direct method should grow roughly quadratically with the number of pixels (since it visits every pixel once in a doubly-nested loop), whereas `filter2D` should remain nearly flat at these small sizes because the actual convolution time is dwarfed by the fixed overhead of the function call. As shown in Fig. 3, this is exactly what happens: the direct method curve grows steeply as the resolution increases, consistent with $\mathcal{O}(H \cdot W \cdot k^2)$ complexity, while both `filter2D` variants stay close to the bottom of the plot across all tested resolutions.

C. Visual and Numerical Comparison

Since both methods are computing the same mathematical convolution, the hypothesis is that their outputs should be identical in the interior of the image, and that any differences should appear only at the borders due to the different ways each method handles boundary pixels. The direct loop simply skips the border (it only iterates from $y = 1$ to $h-2$) and leaves the original values there via the initial copy, whereas `filter2D` applies its own default border extrapolation (`BORDER_REFLECT_101`) before computing the convolution, so the border values will differ.

The results in Fig. 4 confirm this. A pixel-wise comparison shows that 98.9 % of the pixels are numerically identical between the two outputs, with a mean absolute difference of only 0.575 gray levels. The maximum difference reaches 153 gray levels, but as visible in the amplified difference heatmap (bottom right of Fig. 4), these large differences are concen-

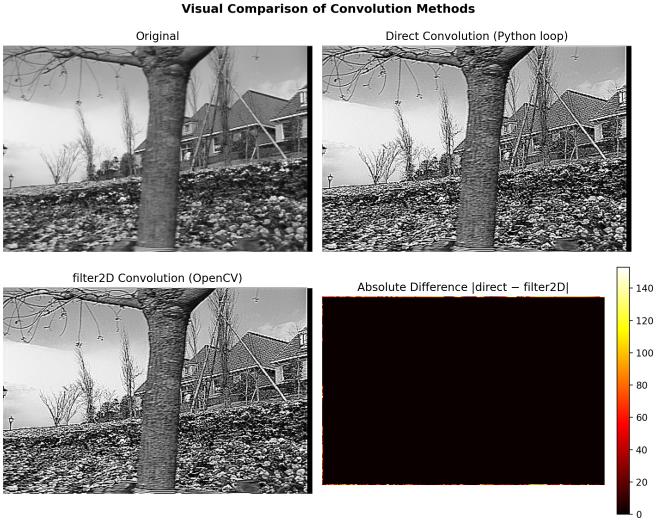


Fig. 4. Visual comparison: original image (top left), direct method result (top right), `filter2D` result (bottom left), and absolute difference amplified $\times 10$ (bottom right).

trated exclusively at the image borders. In the interior, both methods produce the same result, which makes sense because the convolution formula is the same—the only divergence comes from boundary handling, not from the computation itself. This confirms that the performance gap discussed in the previous subsection is purely an implementation-level phenomenon and does not affect the mathematical result.

D. Contrast Enhancement by Unsharp Masking

The convolution kernel provided in the TP code is:

$$K = \begin{pmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{pmatrix}. \quad (3)$$

To understand why this kernel enhances contrast, we can decompose it algebraically. Noting that the coefficients sum to $0 + (-1) + 0 + (-1) + 5 + (-1) + 0 + (-1) + 0 = 1$ and that the central value is 5 while the four cardinal neighbors are -1 , this can be rewritten as the identity kernel minus the discrete Laplacian:

$$K = \underbrace{\begin{pmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{pmatrix}}_{\delta} - \underbrace{\begin{pmatrix} 0 & 1 & 0 \\ 1 & -4 & 1 \\ 0 & 1 & 0 \end{pmatrix}}_{\nabla^2}. \quad (4)$$

This means the convolution result can be expressed as:

$$I_{\text{out}} = I * K = I - \nabla^2 I. \quad (5)$$

The Laplacian $\nabla^2 I$ measures the second-order intensity variation: it is positive where a pixel is darker than its neighbors and negative where it is brighter. Subtracting this from the original image amplifies intensity transitions at edges—the dark side of an edge becomes darker and the bright

Visual Comparison of Convolution Methods		
Absolute Difference direct - filter2D 		
 $\text{Color Scale: } 0 \text{ to } 140$		
Kernel Decomposition: $K = \delta - \nabla^2$		
K (sharpening)	=	δ (identity)
$\begin{pmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{pmatrix}$	$-$	$\begin{pmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{pmatrix}$
∇^2 (Laplacian)		$\begin{pmatrix} 0 & 1 & 0 \\ 1 & -4 & 1 \\ 0 & 1 & 0 \end{pmatrix}$

Fig. 5. Visual verification of the kernel decomposition $K = \delta - \nabla^2$.

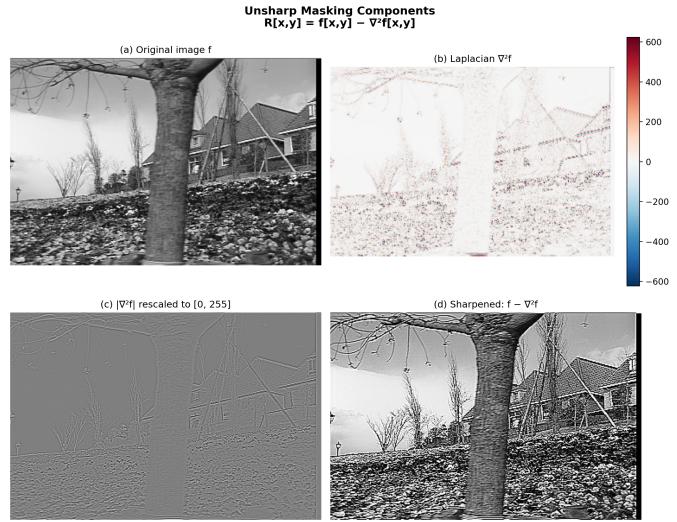


Fig. 6. Unsharp masking components: (a) original image f , (b) Laplacian $\nabla^2 f$ (red = positive, blue = negative), (c) $|\nabla^2 f|$ rescaled to $[0, 255]$, (d) sharpened result $f - \nabla^2 f$ after clipping.

side becomes brighter—while leaving flat regions essentially unchanged (since $\nabla^2 I \approx 0$ there). This is the classical *unsharp masking* mechanism [1].

Based on this decomposition, the hypothesis is as follows: (1) the sharpened image should show visually crisper edges compared to the original, with the effect concentrated around contours and fine structures; (2) the global histogram should become wider (higher standard deviation) because edge pixels are pushed to more extreme values; and (3) the mean intensity should be approximately preserved, since the Laplacian kernel has zero DC response (its coefficients sum to zero) and therefore does not add or remove energy from the image on average.

Fig. 5 shows the decomposition verified numerically: the element-by-element subtraction $\delta - \nabla^2$ reconstructs K exactly.

When applying the decomposition to the actual image, the components are shown in Fig. 6. The Laplacian map (panel b) reveals the edge structure: transitions between regions of different intensity appear as pairs of positive and negative values, while flat areas stay near zero. As expected from the hypothesis, the sharpened image (panel d) shows noticeably crisper edges, especially visible in the flower contours and the fence structures, while smooth regions remain essentially identical to the original.

To visualize the mechanism more clearly, Fig. 7 shows a

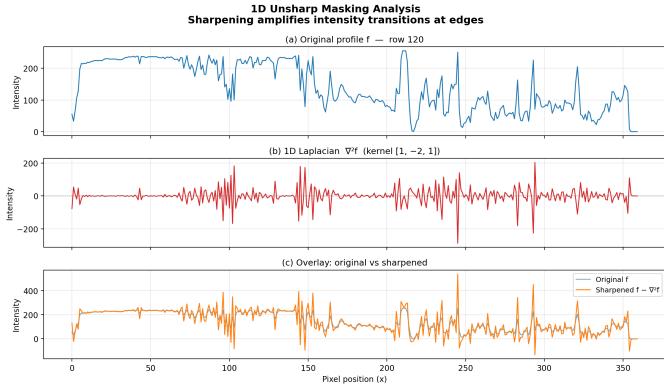


Fig. 7. 1D profile at row 120: (a) original signal, (b) 1D Laplacian, (c) overlay showing how sharpening amplifies transitions.

TABLE II
PIXEL INTENSITY STATISTICS BEFORE AND AFTER SHARPENING.

	Original	Sharp (float)	Sharp (clip)
Min	0.0	-609.0	0.0
Max	255.0	877.0	255.0
Mean	127.14	126.99	129.68
Std	65.90	136.08	89.58

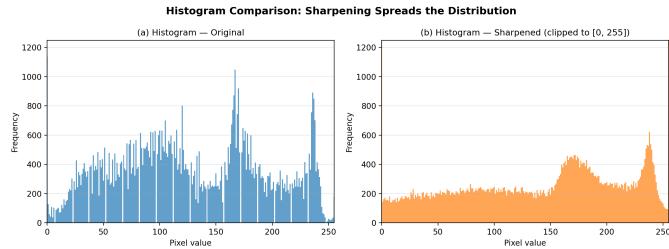


Fig. 8. Histogram comparison: (a) original (linear), (b) sharpened clipped (linear) with saturation markers, (c) sharpened clipped (log scale), (d) sharpened float before clipping.

horizontal 1D intensity profile extracted at the middle row ($y = 120$), with the 1D Laplacian approximated by $[1, -2, 1]$. At every intensity transition, the Laplacian produces a zero-crossing pattern with positive and negative lobes on either side. Subtracting this from the original causes the signal to “overshoot” at each edge: the bright side becomes brighter and the dark side becomes darker, while flat portions of the signal are left untouched. This is precisely the pointwise mechanism behind the contrast enhancement.

To evaluate the hypotheses quantitatively, Table II reports the intensity statistics and Fig. 8 shows the histograms. Looking at the results, the standard deviation increases from 65.90 to 89.58 after clipping (a 35.9 % increase), which confirms that the distribution is indeed wider—the pixel values are more spread out, exactly as predicted. The mean is nearly preserved (127.14 vs 129.68), consistent with the zero-DC-response property of the Laplacian; the small shift of about 2.5 gray levels is explained by the asymmetric effect of clipping, which removes more negative outliers (16.36 % of pixels saturated at 0) than positive ones (9.81 % at 255), slightly biasing the mean

upward.

Before clipping, the sharpened values range from -609 to 877 (Table II), which shows that the Laplacian subtraction can push pixel values far outside the representable $[0, 255]$ interval, especially near strong edges. Panel (d) of Fig. 8 makes this visible: the full float-range histogram extends well beyond the dashed red lines that mark the $[0, 255]$ boundaries. When these values are clamped, they accumulate at the boundaries, producing the saturation peaks visible in panel (b). The log-scale view in panel (c) helps see past these peaks and reveals that the body of the distribution follows a shape similar to the original but stretched, which is consistent with the increased standard deviation.

In summary, the kernel K achieves contrast enhancement because its algebraic structure is that of unsharp masking: the identity component preserves the base brightness while the subtracted Laplacian amplifies local intensity differences around edges. The experimental results confirm the three hypotheses: edges are visually sharper, the histogram is wider, and the mean is approximately preserved. The side effect is that pixels near strong edges can exceed the dynamic range, causing saturation at 0 and 255 after clipping.

II. DESCRIPTORS AND PAIRING

In computer vision, a descriptor is understood as a numerical representation, generally vectorial, that is used for the summarized representation in features of the neighborhood of an entity in an image, which is specifically conceived in order to be able to perform matching or comparison operations of those same entities with other images even when there are events such as changes of scale, rotations, illuminations, or noise that can interfere with the normal matching process [3]; a descriptor can then be understood, therefore, as the result of mapping into a vector a local neighborhood in an image, guaranteeing robustness understood as stability under image changes, viewpoints, geometric distortion and occlusions, the discrimination of different objects given the vector, and efficiency, in such a way that it will be easy to compute and compact in memory, that is, it is a feature vector corresponding to the key points.

On the other hand, a detector is an algorithm whose function is to find in the image a set of points, or keypoints, that are repeatable and well localizable [4]; generally, detectors operate from the definition of a measure $r(x, y)$ or $R(x, y, \omega)$ that is used for the detection of corners, blobs, or stable regions in the domain mainly of local multi-scale characterization, seeking local maxima or minima of these functions and filtering unstable candidates, focusing on repeatability, localization precision of these keypoints in the figure, and stability in detection [5].

A. Oriented FAST and Rotated BRIEF

ORB is a local feature method that outputs keypoints and a binary descriptor; it can be conceptually understood as a pipeline that combines a FAST-family detector and a BRIEF-family descriptor. Considering that since FAST is not scale-

invariant, ORB detects FAST keypoints on multiple rescaled versions of the image, on an image pyramid [6].

The FAST method works by considering a circle of generally 16 pixels in the neighborhood of a candidate pixel p and classifies the pixel p as a corner in the case where there exists a group of pixels on that circle (the neighborhood of p) that have sufficient intensity in comparison with p using a threshold value t , as defined in Eq. (6):

$$S_{p \rightarrow x} = \begin{cases} d, & I_{p \rightarrow x} \leq I_p - t \\ s, & I_p - t < I_{p \rightarrow x} < I_p + t \\ b, & I_p + t \leq I_{p \rightarrow x} \end{cases} \quad (6)$$

where $S_{p \rightarrow x}$ is the “label” of pixel x on the circle.

- d (*dark*): the circle pixel is at least t darker than the center.
- b (*bright*): the circle pixel is at least t brighter than the center.
- s (*similar*): it is inside the band $(I_p - t, I_p + t)$, that is, it does not differ enough.

FAST *declares* that p is a corner if there exists a set of contiguous pixels on the circle such that all of them are *bright* or all of them are *dark*. It is a very fast method because it can be implemented via process optimizations such as the use of a learned decision tree to select which positions to evaluate first [7].

By itself, the FAST segment test constitutes only a binary classification; however, FAST introduces a score value so that non-maximum suppression can be performed and only local maxima are kept [7]. One (efficient) definition given is in Eq. (7):

$$V = \max \left(\sum_{x \in S_{\text{bright}}} (|I_{p \rightarrow x} - I_p| - t), \sum_{x \in S_{\text{dark}}} (|I_p - I_{p \rightarrow x}| - t) \right). \quad (7)$$

In ORB, the FAST threshold is set sufficiently low so as to obtain more than N candidates that can then be evaluated using a Harris corner measure, and to keep the top N values [6].

A standard Harris measure uses the second-moment (structure tensor) matrix H and response in Eq. (8):

$$C = |H| - k(\text{trace}(H))^2. \quad (8)$$

The selection between a score type, more stable with Harris, or faster but slightly less stable with FAST_SCORE, is made available by OpenCV [8].

Given the nature of FAST, it does not naturally produce an orientation; this is why ORB needs to add an orientation step, and for that it implements the centroid idea. ORB improves stability by computing moments only within a circular region of radius r , with raw moments defined in Eq. (9).

$$m_{pq} = \sum_{x,y} x^p y^q I(x, y) \quad (9)$$

$$C = \left(\frac{m_{10}}{m_{00}}, \frac{m_{01}}{m_{00}} \right) \quad (10)$$

Then the orientation is the angle of the vector from patch center O to centroid, as in Eq. (11):

$$\theta = \text{atan2}(m_{01}, m_{10}) \quad (11)$$

On the other hand, BRIEF is responsible for representing a patch p by means of a bitstring that is produced from simple intensity comparisons, as shown in Eq. (12).

$$\tau(p; x, y) = \begin{cases} 1, & p(x) < p(y) \\ 0, & p(x) \geq p(y) \end{cases} \quad (12)$$

An n -bit descriptor can then be built as in Eq. (13).

$$f_n(p) = \sum_{i=1}^n 2^{i-1} \tau(p; x_i, y_i) \quad (13)$$

It focuses on binary strings given the inherent ease of comparing them using Hamming distance rather than L_2 distances in vectors [6]. Plain BRIEF is very sensitive to rotation; in response to this it suggests the concept of steered BRIEF, which rotates the sampling pattern as a function of a discretized angle θ . Let the test locations be encoded in a matrix as in Eq. (14).

$$S = \begin{pmatrix} x_1 & \dots & x_n \\ y_1 & \dots & y_n \end{pmatrix} \quad (14)$$

These locations are rotated according to Eq. (15).

$$S_\theta = R_\theta S \quad (15)$$

Finally, the descriptor is computed using the rotated test coordinates, as expressed in Eq. (16).

$$g_n(p, \theta) := f_n(p)|_{(x_i, y_i) \in S_\theta} \quad (16)$$

ORB then analyzes a subtle but critical issue: when you orient BRIEF consistently, the bit statistics change—the means move away from 0.5 and the tests become less discriminative and more correlated.

The final descriptor used in ORB is an rBRIEF constructed by generating tests with a mean close to 0.5 and high variance, which also preserve low correlation with the tests already selected. The procedure can be summarized as follows: (i) enumerate all pairs of subwindows (then remove overlapping tests), yielding candidate tests; (ii) run each test over all training patches; (iii) sort tests by the distance of their mean from 0.5 (best first); and (iv) apply a greedy selection, keeping a test only if its absolute correlation with all selected tests is below a threshold. This produces a final descriptor that remains binary and fast (Hamming), but with bits that are more informative and less redundant than a naive “steered BRIEF”.

For binary descriptors $d_1, d_2 \in \{0, 1\}^{256}$, matching uses the Hamming distance, as defined in Eq. (17):

$$\text{Ham}(d_1, d_2) = \sum_{i=1}^{256} \mathbf{1}[d_{1,i} \neq d_{2,i}]. \quad (17)$$

This can be computed efficiently as in Eq. (18):

$$\text{Ham}(d_1, d_2) = \text{popcount}(d_1 \oplus d_2). \quad (18)$$

BRIEF emphasizes this efficiency, and ORB notes SSE popcount optimizations in their matching implementation [9].

Finally, in ORB an image pyramid of scales is constructed and, for each scale, FAST corners are detected using a test threshold and are assigned a score using FAST_SCORE or Harris; the strongest features are retained, generally using Harris; the orientation θ is computed via the intensity centroid moments; and finally a descriptor is computed using a smoothed patch, a rotated sampling pattern, and an rBRIEF learned test set, in order to be able to perform descriptor matching using Hamming distance computed via popcount. It is worth highlighting that the way ORB is constructed allows it to handle in-plane rotation, which is addressed through the centroid-based orientation and the steered sampling pattern. Additionally, it can handle image scale by implementing pyramidal detection; however, it remains not fully affine-invariant to viewpoint changes, because projective changes can still break patch appearance.

B. KAZE

As in the case of ORB, KAZE is an algorithm for the detection and description of keypoints, but unlike ORB it constructs the scale space with nonlinear, edge-preserving diffusion; it detects points with a Hessian-type detector and describes them with a (M-)SURF-type descriptor over that nonlinear scale space [10].

KAZE starts with the construction of the nonlinear scale space; for that, a family of images $L(x, y, t)$ is defined, where t plays the role of scale or diffusion time, and it is modeled using the PDE given in Eq. (19):

$$\frac{\partial L}{\partial t} = \text{div}(c(x, y, t) \nabla L). \quad (19)$$

Here, ∇L points toward where the image changes the fastest (edges = large gradient). The diffusion “flow” can be seen as Eq. (20):

$$\mathbf{J} = -c \nabla L, \quad (20)$$

then $\text{div}(\mathbf{J})$ measures how much flow “accumulates” or “leaves” a point, resulting in a definition of brightness propagation analogous to heat propagation, but with a conductivity c that is a controllable parameter [10]. This is precisely the key, because c depends on the gradient, as expressed in Eq. (21):

$$c(x, y, t) = g(\|\nabla L(x, y, t)\|). \quad (21)$$

It should be noted that this gradient is not the raw gradient of the image, but rather the gradient computed on a Gaussian-smoothed version.

If $|\nabla L_\sigma|$ is small, it corresponds to a flat region in the image, and what is sought is to increase diffusion; however, in the opposite case, if $|\nabla L_\sigma|$ is large, it corresponds to a strong edge in the image where the main intention is not to

cross that edge [11]. This is ensured by the two typical forms of g given in Eq. (22) and Eq. (23):

$$g_1(s) = \exp\left(-\frac{s^2}{k^2}\right), \quad (22)$$

$$g_2(s) = \frac{1}{1 + \frac{s^2}{k^2}}. \quad (23)$$

In both definitions, a fundamental element that emerges is k , which is a contrast threshold used as a separation between variations that are considered small for the image, such as noise or textures, and those that are considered large, such as an edge. When k is small, many gradients are considered as edges, which implies that diffusion is stopped in many parts of the image, that is, it is smoothed less; whereas if k is large, only very large gradients are considered as edges, so there is more diffusion and the image is smoothed more.

Since this problem does not have a closed-form analytic solution, KAZE uses numerical schemes that are semi-implicit and that use additive operator splitting in order to construct the scale space with stability [10].

At each level or scale of KAZE, the Hessian response is computed through its determinant and maxima are searched both in position and in scale, as given in Eq. (24):

$$L_{\text{Hessian}} = \sigma^2 (L_{xx} L_{yy} - L_{xy}^2). \quad (24)$$

Here, L_{xx} , L_{yy} , and L_{xy} are the second derivatives (local curvature) measured on L , and the factor σ^2 is a normalization so that the response is comparable across scales (because derivatives “shrink” as scale increases).

Then, KAZE searches for extrema in spatial neighborhoods and across scales, and estimates with precision the position of the maximum that was found with the Hessian response in Eq. (24). KAZE also computes a SURF-type orientation: it takes first-order derivatives in a circular neighborhood with a radius proportional to ω , weights them with a Gaussian, and searches—through a sliding angular window—for a dominant angle.

For the descriptor, KAZE makes use of an M-SURF descriptor, as already mentioned, adapted to the nonlinear scale space, integrating gradient-type responses over sub-patches; the objective of this type of descriptor is to capture how intensity changes around the point in a way that is robust to noise [12].

For a keypoint at scale σ_i , it computes derivatives L_x and L_y at that scale. It builds a 4×4 grid of subregions around the keypoint [10]. In each subregion it sums a vector of the form shown in Eq. (25):

$$d_v = \left(\sum L_x, \sum L_y, \sum |L_x|, \sum |L_y| \right). \quad (25)$$

It then concatenates all subregion vectors to obtain a typical 64-dimensional descriptor, and finally normalizes it. If orientation is used, the sampling is rotated and the derivatives are also computed in that orientation.



Fig. 9. Cross-check matching results using ORB.

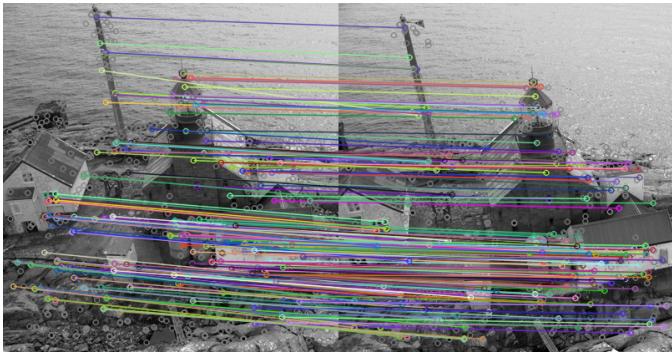


Fig. 10. Cross-check matching results using KAZE.

This allows, finally, each KAZE keypoint to be described as in Eq. (26), together with a descriptor (64D or extended depending on the implementation).

$$(x, y, \sigma, \theta) \quad (26)$$

Given the implementation of extrema detection in nonlinear scale spaces constructed by diffusion, the detector is scale-invariant. The computation of the dominant orientation of the keypoint in order to build the descriptor makes it rotation-invariant. If the upright mode of OpenCV is used, although the algorithm becomes faster, it loses that invariance property by not computing the dominant orientation. It can be said that, due to the normalization inherent to the M-SURF descriptor, this algorithm also obtains descriptors that are approximately contrast-invariant. Finally, KAZE typically behaves well for moderate viewpoints, but it is not “affine-invariant” in the strong sense.

C. Qualitative Performance of descriptors-pairing

As a preliminary stage, it is proposed to qualitatively analyze the results obtained by the descriptor and pairing for matching, FLANN and CROSSCHECK, with and without ratio test, for the ORB and KAZE descriptors on two images for which the second is the result of a transformation of the first.

In Figs. 9–10 the results of the implementation of the detectors and matching using ORB and KAZE, respectively, are

plotted with cross-check matching for a unique value, which consists of taking the descriptors with the smallest distance in the two images and verifying that the correspondence is mutual; if it is, it is identified as a match. In both cases, the best 200 matches between the two images are plotted with lines. From the comparison of the use of both pipelines, results emerge that are immediately striking for the analysis.

It was mentioned in the description of the ORB algorithm that, by combining FAST with steered BRIEF, it is very fast and rotation-invariant, and if it uses a pyramid as is the case, it can handle scale reasonably well; however, including a change in the point of view of the image, which is the case between the compared images, causes anisotropy that makes the pointwise intensity comparisons of the bits performed by BRIEF reduce the rate of correct matches by comparing different entities.

The effect of viewpoint on pairing using cross-check is visible for ORB, since under a slight variation in viewpoint one would expect vectors between the identified pairs with similar directions and lengths; however, in this case, a considerable number of vectors with directions that are not very coherent with the viewpoint change becomes evident. There is also clear evidence of some matches connecting structures that are semantically distinct in the sense of the image, and there is also an accumulation of matches on repeated similar structures that generate ambiguity.

Additionally, in the case of the ORB algorithm, selecting a number n of features—in this case 500—causes regions with strong textures, such as the central area of the image, to occupy most of the keypoints identified in the figure, mainly due to the very fast corner response that is accentuated in structures that are mostly repetitive.

On the other hand, in the case of the KAZE implementation with cross-checking, a much more uniform behavior of the vectors between the paired keypoints identified in the images is observed, with few outliers among the 200 best matches. It is important to highlight that this is favored because, although KAZE is also not constructed to be invariant to point-of-view modifications, the fact that KAZE operates on a nonlinear scale space, achieved through a diffusion process, in some sense reduces localization precision and distinctiveness, so that under slight viewpoint variations it can still deliver repeatable results in the presence of small geometric deformations. Additionally, the form of the descriptor, which uses the notion of gradients in KAZE, is by definition more tolerant to small geometric deformations (such as those in the images under study) than a binary identifier based on pointwise comparisons, because it is integrated over areas rather than exact points, which favors this improved response.

It is important to highlight that the KAZE detector is based on extrema of the normalized Hessian determinant in its nonlinear scale space, so in textures/structures it can produce quite a few valid extrema, even in outer regions. Thus, areas that in the ORB case appeared as a concentration zone due to the presence of repetitive textures, in KAZE tend to yield keypoints that are more spatially distributed. In the image, it can be observed how the keypoints identified with the KAZE

TABLE III
MEASURED RUNTIME FOR DETECTION/DESCRIPTION AND MATCHING CROSS-CHECK.

Detector	Detect+Describe (s)	Matching (s)
ORB	0.011849542	0.001364103
KAZE	0.195190992	0.003179951

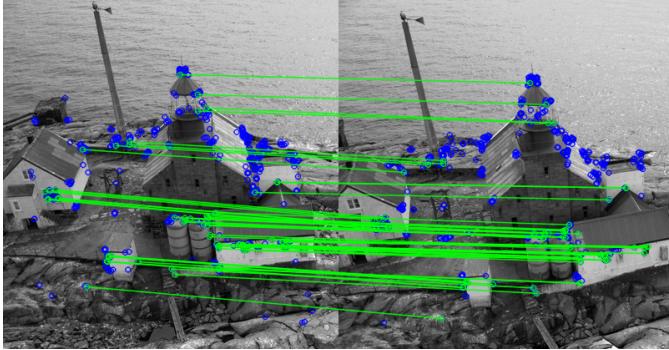


Fig. 11. Ratio-test matching results using ORB.

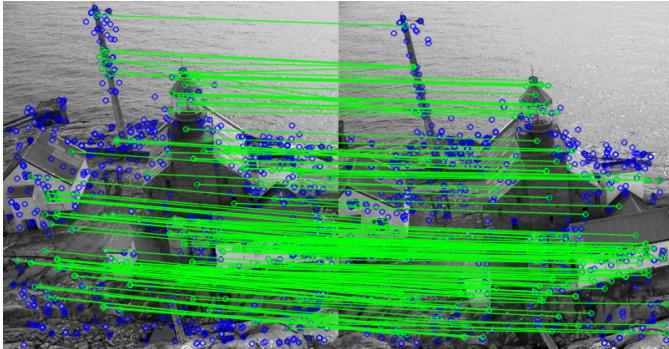


Fig. 12. Ratio-test matching results using KAZE.

algorithm are more spread across the scene, even leading some of them to form matches between figures.

We also report the ratio-test matching results for the same pair. Fig. 11 shows ORB with ratio test, and Fig. 12 shows KAZE with ratio test.

In terms of keypoint detection and description time, the result is consistent with the implemented algorithms. In the case of ORB, since it is based on FAST with simple comparisons and rBRIEF with binary operations, it is computationally cheaper and takes a much shorter time for detection and description. In contrast, KAZE, which constructs scale spaces by nonlinear diffusion, implies solving a diffusion equation at different levels, multiple iterations, and Hessian-type derivative operations; thus it is much heavier and therefore takes significantly more time to execute detection and description on the image. KAZE matching is also slower, mainly because ORB uses binary operations based on Hamming distance, whereas KAZE uses floating-point operations with L_2 distances, implying higher computational complexity, as is evident in Table III.

We also include the FLANN matching results for visual

TABLE IV
MEASURED RUNTIME FOR RATIO-TEST MATCHING.

Detector	Detect+Describe (s)	Matching (s)
ORB	0.019036635	0.003361369
KAZE	0.181437621	0.001895746

TABLE V
MEASURED RUNTIME FOR FLANN MATCHING.

Detector	Detect+Describe (s)	Matching (s)
ORB	0.018674926	0.002685916
KAZE	0.205246737	0.015027942

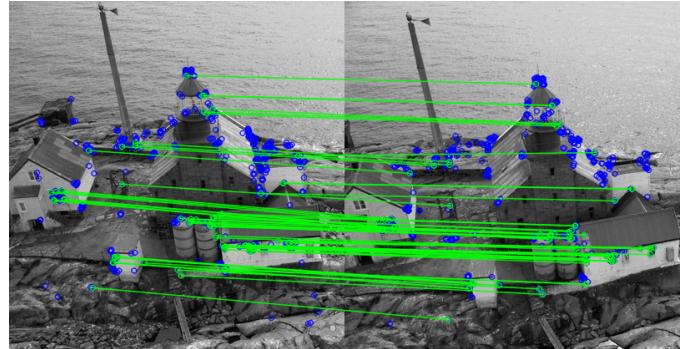


Fig. 13. FLANN matching results using ORB.

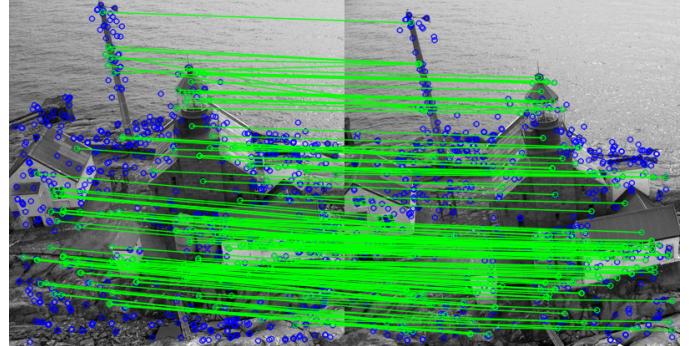


Fig. 14. FLANN matching results using KAZE.

comparison. Fig. 13 shows ORB with FLANN, and Fig. 14 shows KAZE with FLANN.

When analyzing the results, mainly the matching between the two images using ORB and cross-check matching, it is proposed to modify the way the match results are selected between the images, with the main objective of obtaining final values that are more reliable in terms of the quality of the identified matches. For this reason, the ratio test logic is employed to validate the relationship between the paired candidates. In this way, after the detection and description stages described previously, a modification is proposed in the matching algorithm: for each keypoint it identifies the two nearest neighbors and then filters those neighbors, keeping only those for which the best match is at most 70% of the distance of the second-best match, mainly to avoid ambiguities.

As shown in Figs. 11 and 12, the results for ORB when implementing this modification in match selection become evident, since the trajectories of the vectors between pairs in the two images are more uniform than with the simple cross-checking implementation; however, since the implemented detection and description algorithm is the same, the observations regarding the effect of the viewpoint change on pairing quality remain evident, as does the difference with the KAZE results for the same images under the same pair-selection algorithm, given that KAZE can identify approximately three times more correspondences than ORB that satisfy the ratio-test criterion with a threshold of 0.7.

When analyzing the time taken by the algorithms in pairing, a striking result is found: the KAZE time decreases with respect to the cross-checking evaluation. This is explained because, in cross-checking, the algorithm must evaluate L_2 distances in both directions, whereas when cross-checking is set to false and only the two nearest neighbors are kept, the cost decreases because the L_2 computation is done in only one direction; moreover, the increase in execution time due to validating the threshold condition is smaller than the time required by the matching algorithm to perform the second mutual validation, which is why the total time decreases. In contrast, in the case of ORB, since the cost of validating in both directions is minimal because it is a binary operation, the additional computational cost of comparing the two nearest neighbors and applying the threshold causes the time spent by this algorithm in matching to be higher than in cross-checking.

Finally, the implementation of a FLANN algorithm is proposed for checking, enabling fast nearest-neighbor search through an indexed data structure and an approximate k -NN search. The results in terms of pair identification between the images are very similar to the match-ratio case without FLANN; however, the execution times of both algorithms are lower due to the optimized search for the nearest neighbors, which can make it attractive when compared with the simple ratio test. It is important to highlight that Hamming-based cross-checking for ORB remains the fastest in execution, even with FLANN optimizations. It is also important to note that, in this case, the approximate relation of about three times more pairs satisfying the ratio test for KAZE than for ORB is maintained, mainly due to the effects already discussed of nonlinear diffusion and how it favors invariance under small geometric changes such as a small change of point of view.

Although the computation of the distance in each one of the implemented algorithms has already been detailed in the introduction to their functioning, given the effects on execution time mainly in the matching stage, it is convenient to revisit this aspect in greater depth. As was already described, ORB is a binary descriptor that uses an rBRIEF-type descriptor where each component is the result of an intensity comparison, as shown in Eq. (27):

$$b_i = \mathbf{1}[I(x_i) < I(y_i)] \in \{0, 1\}. \quad (27)$$

Then the complete descriptor is given by Eq. (28):

$$\mathbf{b} = (b_1, \dots, b_n) \in \{0, 1\}^n. \quad (28)$$

Thus, the natural way to compare is to identify how many bits change between one descriptor and another, and this measure is precisely the Hamming distance, which is an operation that, as was evidenced computationally, is very efficient because it can be expressed via an XOR operation and a counting of ones. In contrast, in the case of the KAZE descriptor, it is a real vector of 64 components built from sums of derivatives; it is computed over a 4×4 grid and concatenated, yielding a vector as in Eq. (29), which is then normalized:

$$\mathbf{d} \in \mathbb{R}^{64}. \quad (29)$$

Therefore, in its case the distance employed for comparison that makes the most sense is the L_2 distance. Here the components are real (floats) and represent integrated gradient magnitudes. For this type of descriptors, the natural notion of similarity is “how close they are as vectors” in a Euclidean space, as defined in Eq. (30):

$$d_2(\mathbf{d}, \mathbf{d}') = \|\mathbf{d} - \mathbf{d}'\|_2 = \sqrt{\sum_{i=1}^n (d_i - d'_i)^2}. \quad (30)$$

D. Quantitative Performance of descriptors-pairing

It is proposed to evaluate the response of the descriptor and pairing algorithms by analyzing their performance under known geometric transformations applied to the images. For this, it is proposed to analyze the geometric transformations for which the methods are invariant, in order to make evident the pairing dynamics under modifications in rotation and scale. Additionally, it is also proposed to include a transformation that resembles a change of point of view in the image, in order to observe how the pairing precision behaves as the change becomes more significant. For this purpose, the definition of a 2D transformation matrix in OpenCV is used, and then, with that matrix and `warpAffine`, the transformation is applied to the image.

On the side of the implemented algorithm for the evaluation, it is proposed to use the ORB and KAZE detectors for the identification of a set of keypoints in each image, and afterwards, with these, to apply the transformation matrix to the positions of the keypoints of the original image and use the already defined transformation matrix to map the expected position of those keypoints in the transformed image. This is then compared with the position of the nearest keypoint found with `ratioTest` and FLANN; the L_2 distance is measured between the expected position and the detected keypoint position, and a threshold value is proposed to identify the pairing as an error, delivering as a result of the estimation the percentage of correct pairings over the reviewed pairs.

To validate the functioning of the implemented descriptors, mainly in terms of their invariance with respect to the scale

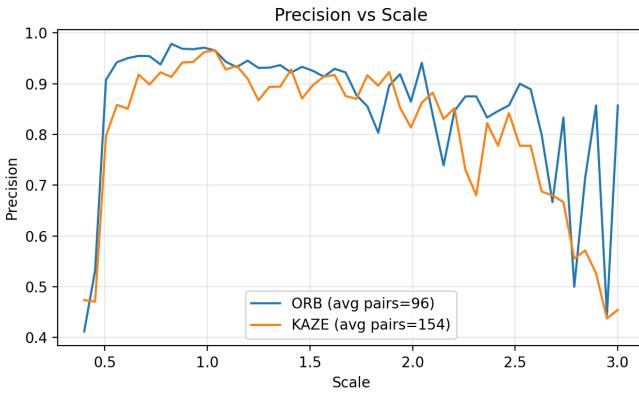


Fig. 15. Precision as a function of scale for ORB and KAZE.

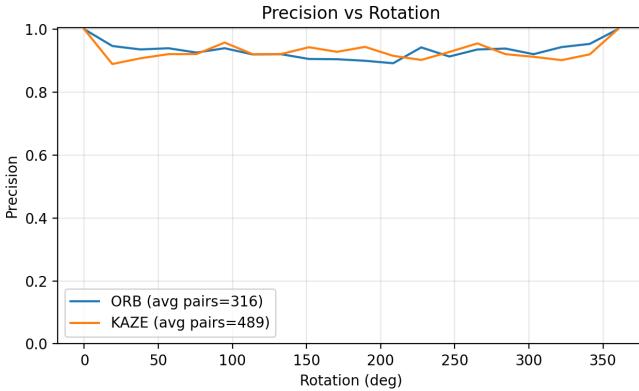


Fig. 16. Precision as a function of rotation for ORB and KAZE.

factor, and following the proposed methodology, it was proposed to analyze the behavior of pairing precision under scale variations from $\times 0.3$ to $\times 3.0$. It is observed that, under scale-only changes, the descriptor behaves as expected, being able to handle the effect of the geometric modification in feature extraction through the specific mechanisms of each algorithm; however, we note that when scale changes are below 50% of the original figure or above 200%, the quality of the results and the pairing becomes less stable and less precise, as shown in Fig. 15.

In accordance with the mathematical conception of the implemented descriptors, it is observed that when applying rotation changes to the second figure and performing pairing with the ORB and KAZE algorithms, the result is stable and additionally of high precision, as shown in Fig. 16. This result was expected, mainly due to the inclusion of a rotation-angle sampling method in the discrete domain with rBRIEF in the case of ORB, and in the case of KAZE through the computation and inclusion of the dominant orientation in the image.

Finally, it is proposed to evaluate the precision of the models under the modification for which, qualitatively, the largest difference in performance between the employed algorithms

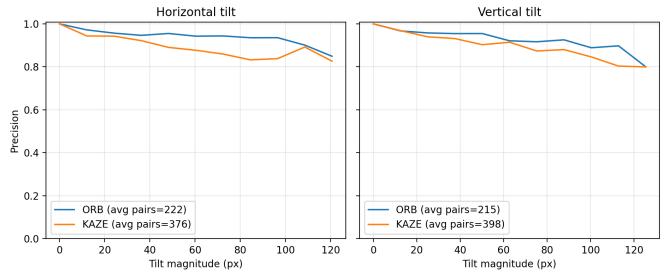


Fig. 17. Precision under viewpoint changes for ORB and KAZE (horizontal and vertical tilt).

was observed; therefore, the implementation of variable points of view is proposed in order to evidence the capacity of the descriptors, especially the one based on nonlinear diffusion, to handle these geometric modifications for which it is not invariant. We evaluated descriptor matching robustness under perspective changes by generating synthetic viewpoint tilts via homographies. Two families of transforms were considered: a horizontal tilt (trapezoidal warping with a shorter top edge and longer bottom edge) and a vertical tilt (asymmetric compression/expansion of the left edge relative to the right). For each tilt magnitude, the original image was warped using the corresponding homography, and matching precision was computed for ORB and KAZE. Precision was measured by projecting keypoints from the original image with the same homography and counting matches whose reprojection error fell below a fixed pixel threshold. The resulting curves show how matching accuracy degrades as the viewpoint distortion increases, enabling a direct comparison of ORB and KAZE under perspective changes.

As evidenced in the results in Fig. 17, although the dynamics are similar in terms of model precision—in both cases, as the tilt magnitude in pixels in the image increases from the simulated edge and point of view, precision decreases—the phenomenon observed for these geometric modifications and previously evidenced qualitatively is maintained: given the nature of the KAZE descriptor, the number of paired keypoints identified is, on average, noticeably higher than the number of keypoints obtained under ORB. Taking into account that, although both models are capable of handling to some extent the geometric change implied by the viewpoint change, KAZE adapts better to it due to the way diffusion is performed and how the descriptor is constructed; this does not necessarily appear as a different dynamic in terms of precision as tilt increases, but it does appear as a higher average number of pairs found in the image.

REFERENCES

- [1] R. C. Gonzalez and R. E. Woods, *Digital Image Processing*, 4th ed. Pearson, 2018.
- [2] OpenCV, “cv::filter2D,” *OpenCV Documentation* (OpenCV 4.x), accessed Feb. 17, 2026. [Online]. Available: https://docs.opencv.org/4.x/d4/d86/group_imgproc_filter.html
- [3] A. Huamán, “Feature Description,” *OpenCV Documentation* (OpenCV 4.14.0-pre), accessed Feb. 6, 2026. [Online]. Available: https://docs.opencv.org/4.x/d5/dde/tutorial_feature_description.html

- [4] A. Huamán, “Feature Detection,” *OpenCV Documentation* (OpenCV 4.14.0-pre), accessed Feb. 6, 2026. [Online]. Available: https://docs.opencv.org/4.x/d7/d66/tutorial_feature_detection.html
- [5] T. Tuytelaars and K. Mikolajczyk, “Local Invariant Feature Detectors: A Survey,” *Foundations and Trends in Computer Graphics and Vision*, vol. 3, no. 3, pp. 177–280, 2007, accessed Feb. 6, 2026. [Online]. Available: <https://lvelho.imp.br/ip08/reading/features.pdf>
- [6] E. Rublee, V. Rabaud, K. Konolige, and G. Bradski, “ORB: an efficient alternative to SIFT or SURF,” in *Proc. IEEE International Conference on Computer Vision (ICCV)*, Nov. 2011, doi: 10.1109/ICCV.2011.6126544, accessed Feb. 6, 2026. [Online]. Available: https://sites.cc.gatech.edu/classes/AY2024/cs4475_summer/images/ORB_an_efficient_alternative_to_SIFT_or_SURF.pdf
- [7] E. Rosten and T. Drummond, “Machine learning for high-speed corner detection,” in *Computer Vision – ECCV 2006* (Lecture Notes in Computer Science), vol. 3951, pp. 430–443, 2006, accessed Feb. 6, 2026. [Online]. Available: https://www.edwardrosten.com/work/rosten_2006-machine.pdf
- [8] OpenCV, “cv::ORB Class Reference,” *OpenCV Documentation* (OpenCV 3.4), accessed Feb. 6, 2026. [Online]. Available: https://docs.opencv.org/3.4/db/d95/classcv_1_1ORB.html
- [9] M. Calonder, V. Lepetit, C. Strecha, and P. Fua, “BRIEF: Binary Robust Independent Elementary Features,” in *Proc. European Conference on Computer Vision (ECCV)*, 2010, accessed Feb. 6, 2026. [Online]. Available: https://www.cs.ubc.ca/~lowe/525/papers/calonder_eccv10.pdf
- [10] P. F. Alcantarilla, A. Bartoli, and A. J. Davison, “KAZE Features,” in *Proc. European Conference on Computer Vision (ECCV)*, 2012, accessed Feb. 6, 2026. [Online]. Available: https://www.doc.ic.ac.uk/~ajd/Publications/alcantarilla_etal_eccv2012.pdf
- [11] P. Perona and J. Malik, “Scale-space and edge detection using anisotropic diffusion,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 12, no. 7, pp. 629–639, 1990, accessed Feb. 6, 2026. [Online]. Available: <https://www.sci.utah.edu/gerig/CS7960-S2010/materials/Perona-Malik/PeronaMalik-PAMI-1990.pdf>
- [12] H. Bay, A. Ess, T. Tuytelaars, and L. Van Gool, “SURF: Speeded Up Robust Features,” *Computer Vision and Image Understanding*, vol. 110, no. 3, pp. 346–359, 2008, accessed Feb. 6, 2026. [Online]. Available: <https://www.cs.jhu.edu/misha/ReadingSeminar/Papers/Bay08.pdf>