



# Assessing the quality of restored images in optical long-baseline interferometry

Nuno Gomes, Paulo J. V. Garcia, Éric Thiébaud

## ► To cite this version:

Nuno Gomes, Paulo J. V. Garcia, Éric Thiébaud. Assessing the quality of restored images in optical long-baseline interferometry. *Monthly Notices of the Royal Astronomical Society*, 2017, 465, pp.3823-3839. <10.1093/mnras/stw2896>. <insu-03710628>

**HAL Id: insu-03710628**

**<https://insu.hal.science/insu-03710628v1>**

Submitted on 1 Jul 2022

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

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



HAL Authorization

# Assessing the quality of restored images in optical long-baseline interferometry

Nuno Gomes,<sup>1,2,3★</sup> Paulo J. V. Garcia<sup>1,2</sup> and Éric Thiébaud<sup>4</sup>

<sup>1</sup>Universidade do Porto - Faculdade de Engenharia, Rua Dr. Roberto Frias, P-4200-465 Porto, Portugal

<sup>2</sup>CENTRA, Instituto Superior Técnico, Av. Rovisco Pais, P-1049-001 Lisboa, Portugal

<sup>3</sup>Universidade do Porto - Faculdade de Ciências, Rua do Campo Alegre, P-4169-007 Porto, Portugal

<sup>4</sup>Centre de Recherche Astrophysique de Lyon/Observatoire de Lyon, 9 Avenue Charles André, F-69561 Saint-Genis Laval Cédex, France

Accepted 2016 November 6. Received 2016 October 27; in original form 2016 July 19

## ABSTRACT

Assessing the quality of aperture synthesis maps is relevant for benchmarking image reconstruction algorithms, for the scientific exploitation of data from optical long-baseline interferometers, and for the design/upgrade of new/existing interferometric imaging facilities. Although metrics have been proposed in these contexts, no systematic study has been conducted on the selection of a robust metric for quality assessment. This article addresses the question: what is the best metric to assess the quality of a reconstructed image? It starts by considering several metrics and selecting a few based on general properties. Then, a variety of image reconstruction cases are considered. The observational scenarios are phase closure and phase referencing at the Very Large Telescope Interferometer (VLTI), for a combination of two, three, four and six telescopes. End-to-end image reconstruction is accomplished with the MiRA software, and several merit functions are put to test. It is found that convolution by an effective point spread function is required for proper image quality assessment. The effective angular resolution of the images is superior to naive expectation based on the maximum frequency sampled by the array. This is due to the prior information used in the aperture synthesis algorithm and to the nature of the objects considered. The  $\ell_1$ -norm is the most robust of all considered metrics, because being linear it is less sensitive to image smoothing by high regularization levels. For the cases considered, this metric allows the implementation of automatic quality assessment of reconstructed images, with a performance similar to human selection.

**Key words:** instrumentation: high angular resolution – instrumentation: interferometers – methods: data analysis – techniques: high angular resolution – techniques: image processing – techniques: interferometric.

## 1 INTRODUCTION

Existing optical long-baseline interferometers provide information at angular scales a factor of 10 smaller than any existing or planned single aperture telescope. This is achieved by measuring interference fringes from pairs of telescopes. The fringes' contrast and position at the detector can be related to the spatial coherence of the incoming electromagnetic field, which in turn contains information on the object brightness distribution (cf. e.g. Buscher 2015; Glindemann 2011). This makes an imaging interferometer very different from an imaging camera. The first difference is related to the information content. A camera generates an image from a continuous sampled pupil, while an interferometer only obtains information at a much smaller number of specific locations of an effective ‘meta-pupil’ – the so-called *uv*-coverage of the data. A second difference

is that while in a camera all the information is obtained simultaneously, in an interferometer data are taken from diverse array combinations separated in time. Finally, for an interferometer an algorithm must be used to synthesize an image.

In optical long-baseline interferometry, phase information degradation by atmospheric turbulence is normally overcome by phase closure triangulation (e.g. Jennison 1958; Monnier 2007), at the expense of further reducing the information content of the measurement. It is therefore not surprising that the first optical long-baseline images were of binaries (morphological simple objects) and were first obtained with three telescopes (Baldwin et al. 1996; Benson et al. 1997). Since the publication of the first relevant results, the technique of image reconstruction of long-baseline interferometric data in the optical/infrared (O/IR; 0.4–20  $\mu\text{m}$ ) regime has evolved and it is nowadays well established. A major breakthrough in optical long-baseline interferometry was the availability of the CHARA and Very Large Telescope Interferometer (VLTI) arrays (ten Brummelaar et al. 2005; Schöller 2007) coupled to the

★E-mail: [nunogomes.pt@gmail.com](mailto:nunogomes.pt@gmail.com)

control of atmospheric effects with spatial filtering (Coudé du Foresto, Ridgway & Mariotti 1997; Tatulli et al. 2010) and adaptive optics (e.g. Arsenault et al. 2004). By combining three or more telescopes and reasonable *uv*-coverages, the information content allowed us to overcome the binary barrier and enter into more complex morphologies such as stellar surfaces and discs (e.g. Le Bouquin et al. 2009; Benisty et al. 2011; Che et al. 2011; Millour et al. 2011; Kloppenborg et al. 2015; Mourard et al. 2015; Hillen et al. 2016).

Because of the low information content of interferometric data, the generation of images is an ill-posed problem with more unknowns than available data. Therefore, images are reconstructed by minimizing a cost function that includes both the data and some prior information on the object brightness distribution (e.g. Thiébaud 2013). To overcome the effects of the turbulence, optical long-baseline interferometry data traditionally rely on the closure phase (and not on the baseline phase). The non-convex nature of the problem makes image reconstruction a difficult task, and algorithms are still a matter of active research (cf. Berger et al. 2012 for a recent review). The availability of dispersed fringes increased the information content of interferometry data, enabling spectral self-calibration (e.g. Millour et al. 2011; Schutz et al. 2014). Other developments are algorithms joining imaging and parametric descriptions of the astronomical objects (e.g. Kluska et al. 2014), or different types of regularization (Renard, Thiébaud & Malbet 2011; Baron et al. 2014).

With the advent of *GRAVITY* at European Southern Observatory, the first common instrument allowing phase referencing observations (Eisenhauer et al. 2008), most of the aperture synthesis algorithms may be simplified, because when a reference source is available, the phase closure is no longer required to remove atmospheric effects and the baseline phase becomes accessible. Standard radio interferometry approaches have proved successful with simulated data in this context (e.g. Vincent et al. 2011).

The large variety of aperture synthesis methods naturally leads to the question on which is the best approach. In 2001, the Working Group on Optical Interferometry of the International Astronomical Union (IAU) decided to compare and promote the development of different algorithms to restore O/IR interferometric images on a regular basis. Starting in 2004, an ‘Imaging Beauty Contest’ has been held by SPIE every two years (Lawson et al. 2004, 2006; Cotton et al. 2008; Malbet et al. 2010; Baron et al. 2012; Monnier et al. 2014a), where contestants present blindly restored images from synthetic or observational data provided by the organization of the contest. They are also asked to interpret the results, indicating what is believed to be real features and what are the potential artefacts of the imaging process. Subsequently, the restored images obtained from the different software are compared to their corresponding reference images by means of a best-fitting method. This method typically comprises a resampling of the restored image to the grid of the reference one, the normalization of the restored image to its peak brightness, and the comparison with the reference image convolved with the effective point spread function (PSF) of the interferometer, using a root-mean-square agreement. However, this approach is limited, because a particular metric might favour a special algorithm for a specific object morphology. This is a pertinent objection which, to our knowledge, is not addressed in the literature.

The work presented here addresses this very question: how can we equitably measure the quality of an image obtained in aperture synthesis? This is a topic of relevance not only for algorithms, but also to the scientific exploitation of aperture synthesis, and for any future infrastructure relying on aperture synthesis imaging,

such as the Planet Formation Imager (Kraus et al. 2014; Monnier et al. 2014b).

This article is structured as follows. In Section 2, we review merit functions used for image quality assessment, and we select a few for further analysis. It is underlined that image convolution with an effective PSF is mandatory. In Section 3, we present the methods we used to recover the interferometric images, explaining how we generate the observables and respective noise, how we restored the images, and how we assess their quality. Important aspects of this approach are (a) both phase closure and phase referencing techniques are addressed, and (b) the array configurations are selected from available stations at the VLTI, particularly the case for four telescopes using phase closure, where the configurations are the ones used with the *PIONIER* instrument. Section 4 concerns about the reconstructed images and the analysis of the behaviour of the selected merit functions. We discuss the results and provide a summary of our findings. The most surprising outcome is that the metric used in the ‘Imaging Beauty Contest’ is biased, but it can be replaced by a simple metric. A side bonus of our approach is that it paves the way for image quality assessment without human intervention. In Section 5 we conclude and present directions for future developments.

## 2 IMAGE QUALITY

The quality of an image has to be assessed by an objective quantitative criterion. What is the best criterion also largely depends on the context. Here we will assume that the *metric*  $\Theta(x, y)$  is used to estimate the discrepancy between a reconstructed image  $x$  and a reference image  $y$ . To simplify the discussion, we also assume that the lower the  $\Theta(x, y)$  the better the agreement between  $x$  and  $y$ . In other words,  $\Theta(x, y)$  can be thought as a measure of the distance between  $x$  and  $y$ .

When assessing image quality, it is important that the result does not depend on irrelevant changes. This, however, depends on the type of images and on the context. For instance, for object detection or recognition, the image metric should be insensitive to the background level, to a geometrical transform (translation, rotation, magnification, etc.) or to a multiplication of the brightness by some positive factor which does not affect the shape of the object. In cases where image reconstruction has underdeterminations, these should not have any incidence on the metric. For optical interferometry and when only power-spectrum and closure phase data are available, the images to be compared may have to be shifted for best matching. In general, the metric should be minimized with respect to the undetermined parameters.

When comparing a true image  $z$  (with potentially an infinitely high resolution) to a restored image  $x$ , the effective resolution achievable by the instrument and the image restoration process must be taken into account. Otherwise and because image metrics are in general based on pixel-wise comparisons, the slightest displacement of sharp features would lead to large loss of quality (according to the metric) whereas the images may look very similar at a lower and more realistic resolution. The easiest solution is then to define the reference image  $y$  to be the true image  $z$  blurred by an effective PSF  $h_{\text{ref}}$ , whose shape corresponds to the effective resolution

$$y = h_{\text{ref}} * z, \quad (1)$$

where the symbol asterisk (\*) denotes the convolution. The choice of the effective resolution is then a parameter of the metric.

To summarize and to be specific, using the distance  $\Theta(x, y)$  between the restored image  $x$  and the reference image  $y$ , the discrepancy between  $x$  and the true image  $z$  would be given by:

$$d(x, z) = \min_{\alpha, \beta, \sigma, t} \Theta(\alpha h_{\sigma, t} * x + \beta, h_{\text{ref}} * z), \quad (2)$$

with  $\alpha$  a brightness scale,  $\beta$  a background, and  $h_{\sigma, t}$  a matching PSF of width parameter<sup>1</sup>  $\sigma > 0$  and centred at position  $t$ . Note that the merit function should be minimized with respect to the width  $\sigma$  of the effective PSF in order to estimate the effective resolution achieved by a given restored image. Our choice to assigning the translation to the matching PSF is to avoid relying on some particular method to perform sub-pixel interpolation (of  $x$ ,  $y$  or  $z$ ) for fine tuning the position. Not doing so would add another ingredient to the metric. When dealing with images with different pixel sizes, the resampling of the images at a given common resolution can be implemented by a linear operator which performs at the same time the resampling, the fine shifting and the blurring by one of the PSFs.

In the following subsections, we first review the most common metrics found in the literature and argue whether they are appropriate or not in the context of optical interferometry. We then propose a family of suitable metrics.

## 2.1 Merit functions

### 2.1.1 Quadratic metrics

Quadratic merit functions are probably the most widely used ones, for they are easy to manipulate and can be made insensitive to various effects, such as an affine change in the image levels (see Section 2.1.2). Even though it is not always obvious, they are, in fact, related to various metrics proposed for comparing images. Compared to the Kullback–Leibler divergence (see Section 2.1.7), quadratic merit functions amount to assuming a simple distribution of the differences between two images (that is to say, independent and Gaussian). The most general expression of a quadratic metric to measure the discrepancy between two images  $x$  and  $y$  takes the form of a weighted (squared)  $\ell_2$ -norm:

$$\text{WL2N}(x, y; W) = \|x - y\|_W^2,$$

where we denote by  $\|q\|_W^2 = q^T W q$  the weighted squared Euclidean norm, with  $W$  a positive (semi-)definite weighting operator. Using a diagonal weighting operator  $W = \text{diag}(w)$  yields:

$$\text{WL2N}(x, y; w) = \sum_i w_i (x_i - y_i)^2, \quad (3)$$

where the sum is carried out for all pixels of the images and where the  $w_i \geq 0$  is the weight of pixel  $i$ .

By choosing specific weights, it is possible to mimic a number of commonly used metrics. For instance, the metric of the *Interferometric Imaging Beauty Contest* (Lawson et al. 2004) is

$$\begin{aligned} \text{IBC}(x, y) &= \sqrt{\text{WL2N}(x, y; w = y / \sum_i y_i)} \\ &= \left[ \frac{\sum_i y_i (x_i - y_i)^2}{\sum_i y_i} \right]^{1/2}, \end{aligned} \quad (4)$$

which amounts to taking the weights as being proportional to the reference image:  $w = y / \sum_i y_i$ . The main drawbacks of this merit function are that it overemphasizes the brighter regions of the image and discards pixels where the reference image  $y$  is zero, which

occurs for many pixels for a compact astronomical source on a dark background. For these reasons, we anticipate that IBC may not be the best metric.

The most simple quadratic metric is the squared  $\ell_2$ -norm (also known as the *squared Euclidean norm*) of the pixel-wise differences between the images:

$$\begin{aligned} \text{L2N}(x, y) &= \|x - y\|_2^2 \\ &= \sum_i (x_i - y_i)^2, \end{aligned} \quad (5)$$

which is WL2N when  $w = 1$ . The *Mean Squared Error* (MSE) is directly derived from the Euclidean norm by taking  $w = 1/N_{\text{pix}}$ , with  $N_{\text{pix}}$  the number of pixels:

$$\text{MSE}(x, y) = \frac{1}{N_{\text{pix}}} \|x - y\|_2^2. \quad (6)$$

The MSE was used by Renard et al. (2011) to benchmark the effects of the regularization in the image reconstruction from interferometric data. For all the metrics presented so far, the smaller the merit value, the more similar are the images.

Some other commonly used metrics are also based on the Euclidean norm of the differences. For instance, the *Peak Signal to Noise Ratio* (PSNR) is

$$\text{PSNR}(x, y) = 10 \times \log_{10} \left( \frac{[\max(y) - \min(y)]^2}{\text{MSE}(x, y)} \right). \quad (7)$$

Here,  $\min(y)$  and  $\max(y)$  correspond respectively to the minimum and maximum possible pixel value of the reference image  $y$ . The PSNR is given in decibel (db) units and the higher the PSNR, the more similar are the images.

Clearly, MSE and PSNR are the squared Euclidean norm of the pixel-wise difference between the images (L2N) but expressed in different units. They can be used interchangeably and we will only consider IBC and L2N in what follows.

### 2.1.2 Minimizing the discrepancy with respect to the brightness distortion

In order to make a formal link between different metrics, it is worth investigating what happens when the minimization with respect to the brightness distortion parameters  $\alpha$  and  $\beta$  is carried on. As we will show, this minimization has a closed form solution with a quadratic metric:

$$\|\alpha x + \beta \mathbb{1} - y\|_W^2,$$

with  $x$  and  $y$  the images to compare,  $\alpha \in \mathbb{R}^+$  a positive factor,  $\beta \in \mathbb{R}$  a constant background, and  $\mathbb{1}$  an image where all pixels are equal to 1.

Let us first consider the constant background correction. Introducing  $r = y - \alpha x$ , we want to minimize  $\|r - \beta \mathbb{1}\|_W^2$  with respect to  $\beta$ . Expanding the quadratic norm yields

$$\|r - \beta \mathbb{1}\|_W^2 = \|r\|_W^2 - 2(\mathbb{1}^T W r) \beta + \|\mathbb{1}\|_W^2 \beta^2.$$

This is a simple 2nd order polynomial in  $\beta$  and the minimum is achieved for the optimal background

$$\beta^* = \frac{\mathbb{1}^T W r}{\mathbb{1}^T W \mathbb{1}}, \quad (8)$$

which can be seen as a weighted averaging of  $r$ . Thus,

$$\min_{\beta} \|r - \beta \mathbb{1}\|_W^2 = \|r - \beta^* \mathbb{1}\|_W^2 = \|C r\|_W^2, \quad (9)$$

<sup>1</sup> In this paper we took  $\sigma$  to be the standard deviation of the PSF profile.

where the linear operator  $C$  is given by

$$C = I - \frac{\mathbb{1}^T W}{\mathbb{1}^T W \mathbb{1}}, \quad (10)$$

and  $I$  is the identity. The linear operator  $C$  has the effect of removing the weighted average of its argument. Replacing  $r$  by  $y - \alpha x$  yields:

$$\min_{\beta} \|\alpha x + \beta \mathbb{1} - y\|_W^2 = \|\alpha Cx - Cy\|_W^2, \quad (11)$$

which amounts to comparing the weighted average subtracted images.

The expansion

$$\|\alpha x - y\|_W^2 = \|y\|_W^2 - 2(y^T W x)\alpha + \|x\|_W^2 \alpha^2$$

readily shows that the optimal factor  $\alpha$  is

$$\arg \min_{\alpha} \|\alpha x - y\|_W^2 = \frac{y^T W x}{x^T W x},$$

and, after trivial simplifications, that

$$\min_{\alpha} \|\alpha x - y\|_W^2 = \|y\|_W^2 - \frac{(y^T W x)^2}{\|x\|_W^2}.$$

Putting all together we have shown that

$$\min_{\alpha, \beta} \|\alpha x + \beta \mathbb{1} - y\|_W^2 = \|Cy\|_W^2 - \frac{(y^T C^T W C x)^2}{\|Cx\|_W^2}, \quad (12)$$

where the linear operator  $C$  is given in equation (10). If no background correction is wanted, it is sufficient to take  $C = I$ . The above expression can be divided by  $\|Cx\|_W^2$  to obtain a symmetric distance between  $x$  and  $y$  which is independent of an affine transform of the brightness of any of the two images

$$d(x, y) = 1 - \text{Corr}(x, y)^2, \quad (13)$$

with

$$\text{Corr}(x, y) = \frac{y^T C^T W C x}{\|Cx\|_W \|Cy\|_W} \quad (14)$$

the (weighted) correlation between the two images  $x$  and  $y$ . If  $W \propto I$ , then the usual definition of the correlation, given in equation (16), is retrieved.

The distance  $d(x, y)$  takes values in the range  $[0, 1]$ , the smaller it is the better is the agreement. Conversely, the better the agreement the larger the absolute value of the (weighted) correlation. It is therefore clear now that comparing images by means of their (weighted) correlation coefficient is equivalent to using a quadratic norm minimized with respect to an affine transform of the image intensity.

### 2.1.3 Universal image quality index and image structural similarity

The *universal image quality index* was proposed by Wang & Bovik (2002) to overcome MSE and PSNR, which were found to be very poor estimators of the image quality for common brightness distortions and image corruptions (like salt-and-pepper noise, lossy compression artefacts, etc.). The universal image quality index is defined as

$$Q(x, y) = \frac{4 \text{Avg}(x) \text{Avg}(y) \text{Cov}(x, y)}{(\text{Avg}(x)^2 + \text{Avg}(y)^2)(\text{Var}(x) + \text{Var}(y))}, \quad (15)$$

where  $\text{Avg}(x)$ ,  $\text{Var}(x)$  and  $\text{Cov}(x, y)$  are respectively the empirical average, variance and covariance of  $x$  and  $y$ , given by:

$$\text{Avg}(x) = \frac{1}{N_{\text{pix}}} \sum_i x_i,$$

$$\text{Var}(x) = \text{Cov}(x, x),$$

$$\text{Cov}(x, y) = \frac{1}{N_{\text{pix}} - 1} \sum_i (x_i - \text{Avg}(x))(y_i - \text{Avg}(y)).$$

The universal image quality index takes values in the range  $[-1, 1]$ .  $Q(x, y)$  is maximal for the best agreement, which occurs when  $y = \alpha x + \beta$ , and minimal when  $y = -\alpha x + \beta$ , for any  $\alpha > 0$  and any  $\beta$ . Although the universal image quality index was designed to cope with brightness distortions such as mean shift or dynamic shrinkage, this indicator is not exactly insensitive to any affine transform of the intensity as is (see the demonstration in Section 2.1.2) the correlation coefficient:

$$\text{Corr}(x, y) = \frac{\text{Cov}(x, y)}{\sqrt{\text{Var}(x) \text{Var}(y)}}. \quad (16)$$

In order to improve over the universal image quality index, Wang et al. (2004) introduced the *image Structural SIMilarity* (SSIM):

$$\begin{aligned} \text{SSIM}(x, y) &= \frac{2 \text{Avg}(x) \text{Avg}(y) + \varepsilon_1}{\text{Avg}(x)^2 + \text{Avg}(y)^2 + \varepsilon_1} \\ &\times \frac{2 \text{Cov}(x, y) + \varepsilon_2}{\text{Var}(x) + \text{Var}(y) + \varepsilon_2}, \end{aligned} \quad (17)$$

where  $\varepsilon_1 > 0$  and  $\varepsilon_2 > 0$  are small values introduced to avoid divisions by zero. Note that with  $\varepsilon_1 = 0$  and  $\varepsilon_2 = 0$ , the SSIM is just the image quality index defined in equation (15). The higher the SSIM, the better the agreement. In principle SSIM and the quality index should be used *locally*, that is on small regions of the images.

### 2.1.4 Accuracy function

Similarly to the IBC metric, the *accuracy function* (ACC, Gomes 2016) is based on a normalized weighted quadratic difference between the reconstructed image  $x$  and the reference image  $y$ :

$$\text{ACC}(x, y) = \frac{\sum_i w_i (x_i - y_i)^2}{\sum_i (x_i + y_i)^2}. \quad (18)$$

Here  $w$  is a normalized weighting function, a mask that eliminates all pixels where the reference and the restored images have intensities smaller than the image's dynamic range. On all non-negligible pixels,  $w$  is equal to 1.

ACC varies between 0 and 1 and the smaller it is, the greater the resemblance between both images. Note that the accuracy function is neither quadratic in  $x$  nor in  $y$ .

### 2.1.5 Sum of absolute differences

One of the drawbacks of quadratic metrics is that they strongly emphasize the largest differences. To avoid this, an  $\ell_p$ -norm can be used with an exponent  $p < 2$ . For instance, the *sum of absolute differences* or  $\ell_1$ -norm is given by:

$$\begin{aligned} \text{L1N}(x, y) &= \|x - y\|_1 \\ &= \sum_i |x_i - y_i|. \end{aligned} \quad (19)$$



### 2.1.6 Fidelity function

The *fidelity function* was introduced by Pety, Gueth & Guilleaume (2001b) in the context of image reconstruction for ALMA. It is defined as the ratio of the total flux of the reference  $y$  to the difference between the restored image  $x$  and the reference one:

$$\text{FID}(x, y) = \frac{\sum_i y_i}{\sum_i \max\{\eta, |y_i - x_i|\}}, \quad (20)$$

where  $\eta$  is some non-negative threshold. The higher the fidelity value, the better the agreement.

Choosing  $\eta > 0$  avoids divisions by zero, and Pety et al. (2001b) took  $\eta = 0.7 \text{RMS}(x - y)$ , where  $\text{RMS}(\dots)$  yields the root mean squared value of its argument. We note that with  $\eta > 0$ , all differences smaller than  $\eta$  have the same incidence on the total cost and are therefore irrelevant. To avoid this, one has to take  $\eta = 0$ , in which case the reciprocal of the fidelity function is then just the  $\ell_1$ -norm defined in equation (19) times some constant factor which only depends on the reference  $y$ . As the fidelity function would then yield the same results as the  $\ell_1$ -norm, we only consider the latter in our study.

### 2.1.7 Kullback–Leibler divergence

Being non-negative everywhere and normalized, the images can be thought as distributions (over the pixels). The *Kullback–Leibler divergence* measures the similarity between two distributions. When applied to our (normalized) images it writes

$$\text{KL}(x, y) = \sum_i y_i \log(x_i/y_i). \quad (21)$$

A restriction for the Kullback–Leibler divergence is that  $x$  and  $y$  must be strictly positive everywhere. It is however possible to account for non-negative distributions by modifying the definition of the Kullback–Leibler divergence as follows:

$$\text{KL}(x, y) = \sum_i c_{\text{KL}}(x_i, y_i),$$

where  $c_{\text{KL}}(q, r)$  extends  $r \log(q/r)$  by continuity:

$$c_{\text{KL}}(q, r) = \begin{cases} 0 & \text{if } q = r, \text{ or } q > 0 \text{ and } r = 0, \\ -\infty & \text{if } q = 0 \text{ and } r > 0, \\ r \log(q/r) & \text{otherwise.} \end{cases}$$

Note that the Kullback–Leibler divergence is not symmetric, i.e.  $\text{KL}(x, y) \neq \text{KL}(y, x)$ . The Kullback–Leibler divergence is less or equal to zero. The lower the Kullback–Leibler divergence the worse is the agreement between  $x$  and  $y$ . The maximal value of the Kullback–Leibler divergence is equal to zero and is achieved when  $x = y$ .

Like the IBC metric, the Kullback–Leibler divergence disregards  $x_i$  where  $y_i = 0$ . In addition, any image  $x$  with at least one pixel, say  $i_0$ , such that  $x_{i_0} = 0$  while  $y_{i_0} > 0$  yields  $\text{KL}(x, y) = -\infty$ , which corresponds to the maximum possible discrepancy. These are serious drawbacks for using the Kullback–Leibler divergence as an image metric, because it could not make a distinction between restored images such that  $x_{i_0} = 0$ , whatever the values of the other pixels.

### 2.1.8 Designing the metric

We want to derive an image metric that is adapted to our particular case: we consider images of compact objects (i.e. with finite size

support) over a constant background, and which may be shifted by an arbitrary translation.

We assume that  $d(x, y, t)$  yields the discrepancy between the image  $x$  and the image  $y$  shifted by a translation  $t$ . Quite naturally, we require that the following properties hold:

- (i) The metric does not change if the images are extended with pixels set with the background level; likewise, the metric does not change if the images are truncated, provided that the values of the removed pixels equal the background level;
- (ii) The metric is non-negative and equal to zero if the two images are the same (for a given relative translation); in particular  $d(x, x, 0) = 0$ , whatever the image  $x$ ;
- (iii) The metric is *stationary* in the sense that whatever the images  $x$  and  $y$  and the translations  $t, t'$  and  $t''$ ,

$$d(s(x, t), s(y, t'), t'') = d(x, y, t + t'' - t'), \quad (22)$$

where  $s(x, t)$  yields image  $x$  shifted by translation  $t$ :

$$s(x, t)_i = x_{i-t}.$$

A last requirement, although optional, could be:

- (i) the metric is *symmetric* in the sense that

$$d(y, x, -t) = d(x, y, t), \quad (23)$$

whatever the images  $x$  and  $y$  and the translation  $t$ .

To limit the number of possibilities, we consider that the metric is the sum of a pixel-wise cost. Then, accounting for property (i),

$$d(x, y, t) = \sum_{i \in \mathbb{Z}^n} c(\tilde{x}_i, \tilde{y}_{i-t}), \quad (24)$$

where  $n$  is the number of dimensions of the images  $x$  and  $y$  (in our case,  $n = 2$ ),  $\mathbb{Z}$  is the set of integers,  $t \in \mathbb{Z}^n$  is the considered translation,  $c(q, r)$  is the pixel-wise cost, and  $\tilde{x}$  (resp.  $\tilde{y}$ ) is the image  $x$  (resp.  $y$ ) infinitely extended with the background level  $\beta$ :

$$\tilde{x}_i = \begin{cases} x_i & \text{if } i \in \mathbb{X}; \\ \beta & \text{else,} \end{cases} \quad (25)$$

with  $\mathbb{X} \subset \mathbb{Z}^n$  (resp.  $\mathbb{Y} \subset \mathbb{Z}^n$ ) the support of the image  $x$  (resp.  $y$ ). We note that property (ii) implies that  $c(q, q) = 0$  whatever  $q \in \mathbb{R}$ , and also that the background level must be the same for the two images. We also note that property (iv) implies that the pixel-wise cost be a symmetric function, i.e.  $c(q, r) = c(r, q)$  whatever  $(q, r) \in \mathbb{R}^2$ . Finally, property (iii) holds because the same pixel-wise cost is used whatever the index  $i$ .

As  $c(\beta, \beta) = 0$ , the sum over the infinite set  $\mathbb{Z}^n$  in equation (24) simplifies to sums over three finite (and possibly empty) subsets:

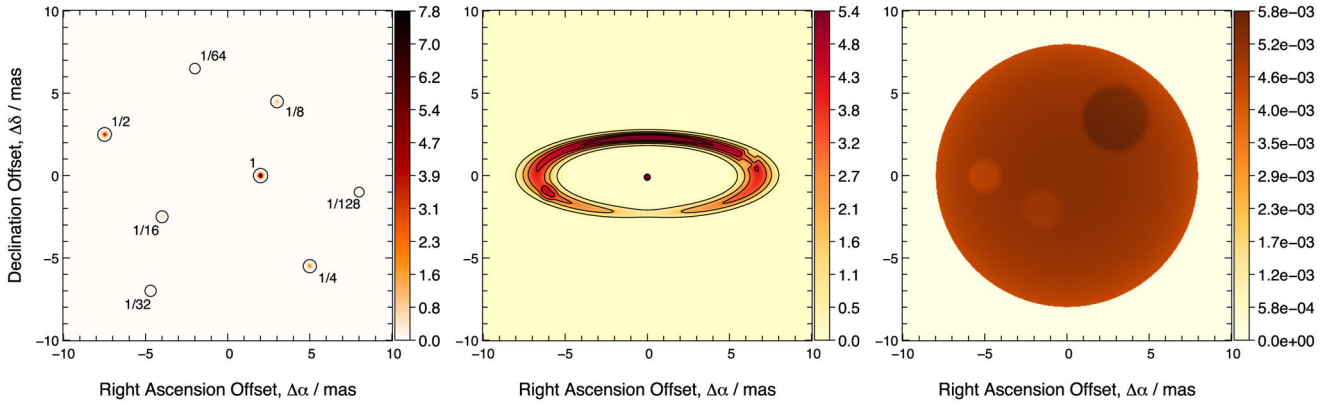
$$d(x, y, t) = \sum_{i \in \mathbb{X} \cap \mathbb{Y}_t} c(x_i, y_{i-t}) + \sum_{i \in \mathbb{X} \setminus \mathbb{Y}_t} c(x_i, \beta) + \sum_{i \in \mathbb{Y} \setminus \mathbb{X}_t} c(y_i, \beta), \quad (26)$$

where  $\mathbb{A} \setminus \mathbb{B}$  denotes the set of elements of  $\mathbb{A}$  which do not belong to  $\mathbb{B}$ , and

$$\mathbb{X}_t = \{i \in \mathbb{Z}^n \mid i - t \in \mathbb{X}\}$$

is the set of indices  $i$  such that  $i - t$  belongs to the support of  $x$ . An efficient implementation of the metric may be achieved with:

$$d(x, y, t) = \gamma + \sum_{i \in \mathbb{X} \cap \mathbb{Y}_t} [c(x_i, y_{i-t}) - c(x_i, \beta) - c(y_{i-t}, \beta)], \quad (27)$$



**Figure 1.** True images ( $z$ ) used for the image reconstruction study: stellar cluster (left), YSO (centre), and stellar photosphere (right). The images are normalized by their total flux. The colour bars indicate surface flux. The stars of the cluster have relative intensities as indicated in the figure. The circles point the position of the stars. The colour maps have been chosen in order to maximize the contrast of the features in each image.

where  $c(x_i, \beta)$  (resp.  $c(y_i, \beta)$ ) can be pre-computed for all  $i \in \mathbb{X}$  (resp. for all  $i \in \mathbb{Y}$ ) and

$$\gamma = \sum_{i \in \mathbb{X}} c(x_i, \beta) + \sum_{i \in \mathbb{Y}} c(y_i, \beta).$$

Finally, it remains to choose the pixel-wise cost  $c(q, r)$ . A whole family of merit functions can be derived with the following pixel-wise cost

$$c(q, r) = \left| \Gamma(q) - \Gamma(r) \right|^p \quad (28)$$

where  $p > 0$  is a chosen exponent and  $\Gamma$  is a function used to emphasize the discrepancy in the low/high range of the brightness distribution. For example, taking

$$\Gamma(q) = \text{sign}(q) |q|^\gamma, \quad (29)$$

with  $\gamma \in [0, 1]$ , it amounts to paying more attention to the least bright part of the images. Taking  $p = 2$  and  $\gamma = 1$  yields the  $\ell_2$ -norm (L2N), while taking the quadratic merit  $p = 1$  and  $\gamma = 1$  yields the  $\ell_1$ -norm (L1N). Incidentally, this shows that the required aforementioned properties (including the symmetry) do hold for these norms.

### 2.1.9 Choice of the candidates

We already mentioned that not all merit functions reviewed in this paper are appropriate for comparing synthetic aperture images. For example, we disregarded the Kullback–Leibler divergence (see Section 2.1.7) because of its inability to distinguish between very different images which have pixels equal to zero while they are non-zero in the reference image. In our context, the background level is known (i.e.  $\beta = 0$  which corresponds to the positivity constraint) and should not have to be adjusted when comparing images. The Universal Quality Index and Image Structural Similarity described in Section 2.1.3 are therefore not appropriate for our needs. However, these metrics can be of value in image patches with non-zero backgrounds.<sup>2</sup> The brightness scale  $\alpha$  may have to be tuned so as to minimize the discrepancy between the images because, on the one hand, they may have different normalization constraints and, on the other hand, they may have been interpolated to cope with

different pixel sizes. As we have shown in Section 2.1.2, minimizing a quadratic cost function in  $\alpha$  would be equivalent to use the correlation of the images as a metric.

To summarize, we will compare images using the  $\ell_2$ -norm (L2N), the  $\ell_1$ -norm (L1N), the metric used in the past IBC and the accuracy function (ACC).

## 3 METHODS

### 3.1 Synthetic image library

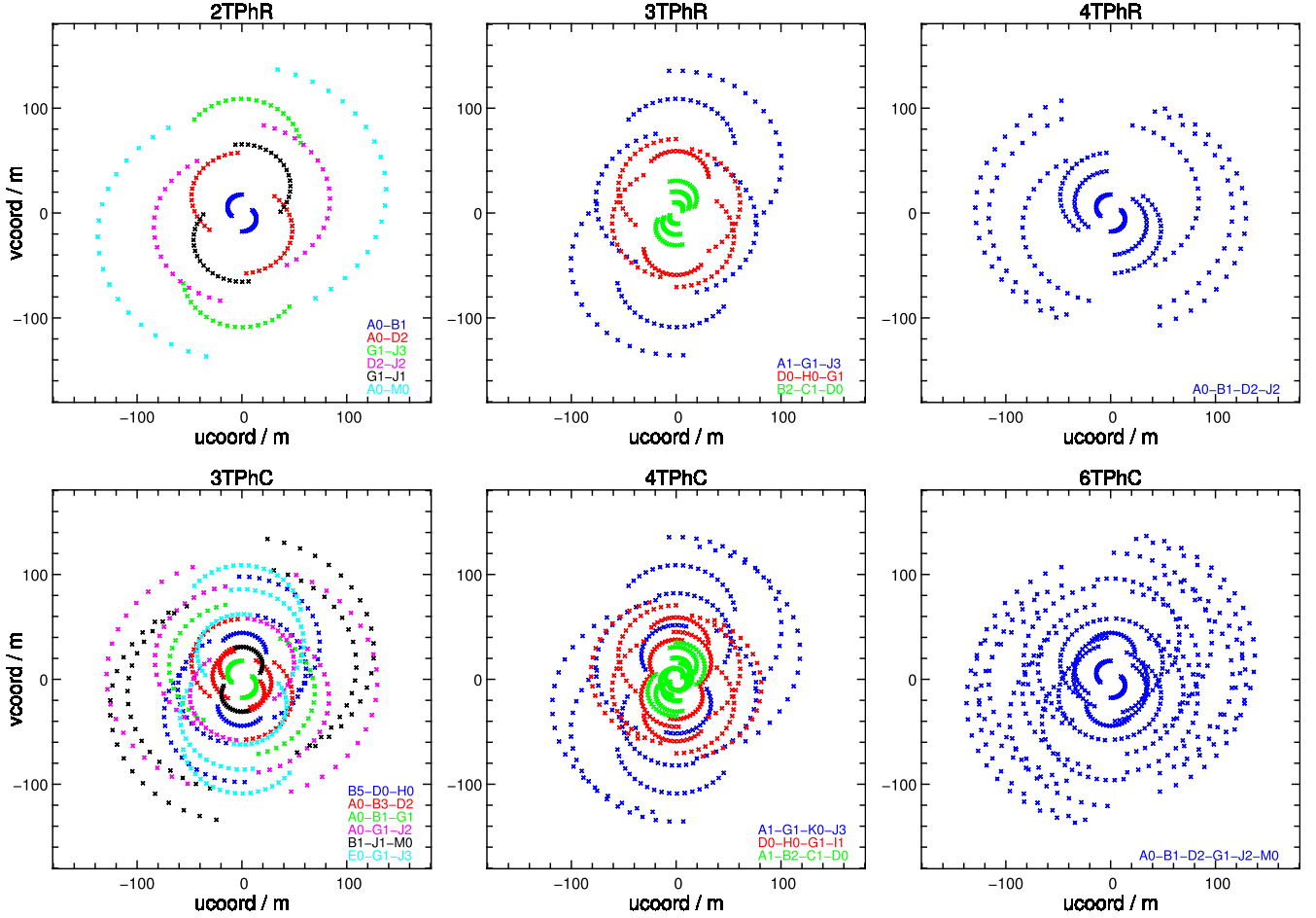
The true images ( $z$ ) used in the study are presented in Fig. 1. They span representative science cases of interferometric imaging (cf. e.g. Berger et al. 2012): compact clusters/multiple stellar systems, young stellar objects (YSOs) and stellar surfaces. We fixed the size of the images to ease the interpretation of the results. The width of the pixel is 0.04 mas. The images cover a wide range of visibilities, from the very sharp cluster to the over-resolved stellar photosphere. The cluster consists of eight stars ‘randomly’ spread in the FOV, with a Gaussian profile of standard deviation 0.1 mas, whose intensities decrease in factors of 2. The typical separation between neighbouring stars is 5 mas. The YSO consists of a central star and a circumstellar disc, with a total flux ratio of 10 to 1. The disc has two features: a dark spot on the first quadrant and a bright spot in the third quadrant. The stellar surface has two bright spots in the third quadrant, and a dark spot on the first quadrant.

### 3.2 UV-space generation

We used realistic  $uv$ -coverages for the VLTI station positions.<sup>3</sup> Six observational configurations are considered, corresponding to one, three and six nights of observation, and to phase referencing (PhR) and phase closure (PhC) data. The station configurations are inspired in previous imaging studies (Filho et al. 2008a,b), and are representative of several instruments: *PRIMA* (2TPhR; Delplanck 2008), *AMBER* (3TPhC; Petrov et al. 2007), *GRAVITY* (3T-4TPhR; Eisenhauer et al. 2011), *PIONIER* (4TPhC; Eisenhauer et al. 2011; Le Bouquin et al. 2011), and *VSI* (6TPhC, Malbet, Kern & Berger 2006). To compute the  $uv$ -tracks, which depend on

<sup>2</sup> Using these metrics would also imply the definition of a patch size, which would open other questions outside the scope of this article.

<sup>3</sup> Available at <https://www.eso.org/observing/etc/doc/viscalc/vltstations.html>.



**Figure 2.** *UV*-coverages of the observational configurations used in the study. PhR stands for phase referencing and PhC for phase closure. The observing nights are fixed for each column and are as follows: six nights left column, three nights central column and one night right column. The stations used in each configuration are indicated.

the object position, observatory location, station positions and hour-angle of the observations (Thompson, Moran & Swenson 2001), the following assumptions were made: (i) object declination of  $-60^\circ$ , (ii) a full *uv*-track corresponding to 19 instantaneous and evenly sampled data points, during a 9 h transit, and (iii) fixed station configurations during each night. The corresponding *uv*-coverages are presented in Fig. 2.

### 3.3 Noise model

The observables used in this study are the visibility amplitude  $V$ , the baseline visibility phase  $\phi$ , the squared visibility  $V^2$ , the bi-spectrum  $\mathcal{B}$ , and the closure phase  $\phi_c$ . A synthetic observable  $o_s$  is generated by

$$o_s \sim \mathcal{N}(E\{o\}, \text{Var}\{o\}),$$

where the expected value of the observable ( $E\{o\}$ ) is computed by interpolating the reference image at the angular frequencies of the observations,<sup>4</sup> using the MiRA package.<sup>5</sup> We adopted the *Simple Noise Model* (Gomes 2016), which is Gaussian and described by

one free parameter, the signal-to-noise ratio (SNR). It is assumed to be  $\text{SNR} = 20$ , a value typical of good quality interferometric observations. The variance of the noise for the  $n$ th visibility amplitude is defined as

$$\text{Var}\{V_n\} = \left( \frac{\langle V \rangle}{\text{SNR}} \right)^2, \quad (30)$$

where  $\langle V \rangle$  is the average of all visibility amplitudes for a given *uv*-coverage (cf. Table 1).

In order to derive the noise for the baseline phase, we assume that the complex visibility has independent real and imaginary parts, with the same Gaussian noise (Goodman approximation; Goodman 1985). The variance of the noise for the  $n$ th baseline phase becomes

$$\text{Var}\{\phi_n\} = \frac{\text{Var}\{V_n\}}{V_n^2}. \quad (31)$$

The noise for the remaining observables can be determined by error propagation.

The simple noise model is in contrast with the one used by Renard, Thiébaud & Malbet (2011), since it initially sets the noise in the visibility amplitude instead of the phase, making the noise in the phase increase with decreasing visibility amplitude. It also qualitatively agrees with Tatulli & Chelli (2005), where the visibility SNR increases with the visibility amplitude.

<sup>4</sup> The observing wavelength is taken at the centre of the  $K$  band:  $2.179 \mu\text{m}$ .

<sup>5</sup> Available for download at <http://cral.univ-lyon1.fr/labo/perso/eric.thiebaut/?Software/MiRA>.



**Table 1.** Mean values of the distribution of the visibility amplitudes for the objects in each  $uv$ -configuration. The errors correspond to the standard deviation.

Object	2TPhR	3TPhC	3TPhR	4TPhC	4TPhR	6TPhC
Stellar cluster	$0.6 \pm 0.2$	$0.6 \pm 0.2$	$0.6 \pm 0.2$	$0.6 \pm 0.2$	$0.6 \pm 0.2$	$0.6 \pm 0.2$
YSO	$0.4 \pm 0.2$	$0.3 \pm 0.2$	$0.5 \pm 0.3$	$0.5 \pm 0.3$	$0.4 \pm 0.2$	$0.3 \pm 0.2$
Stellar photosphere	$0.3 \pm 0.2$	$0.2 \pm 0.2$	$0.4 \pm 0.3$	$0.4 \pm 0.3$	$0.3 \pm 0.2$	$0.2 \pm 0.2$

### 3.4 Image reconstruction with MiRA

The noisy data generated are saved in an OIFITS file (Pauls et al. 2005) and used as input for the MiRA image reconstruction software, assuming monochromatic data. As the goal of the study is to find the best metric for image reconstruction, the actual algorithm is not relevant, as long as it remains the same for all metrics. The MiRA software and its principles are described in detail by Thiébaud (2008, 2013). To summarize, MiRA searches for the image  $x^+$  which minimizes the two-term penalty criterion:

$$x^+ = \arg \min_x \left\{ f(x) = f_{\text{data}}(x|d) + \mu f_{\text{prior}}(x) \right\}. \quad (32)$$

The term  $f_{\text{data}}(x|d)$ , usually known as the *likelihood term*, measures the discrepancy between the actual data  $d$  (e.g. squared visibilities  $V^2$ , visibility amplitudes  $V$ , baseline phases  $\phi$ , and closure phases  $\phi_c$ ) and their model, given the image,  $x$ . The term  $f_{\text{prior}}(x)$ , commonly designated as the *regularization term*, is a penalty which enforces additional priors, and it is required to avoid artefacts. It is needed because the data alone cannot unambiguously yield a unique image. The so-called *level of regularization* or *hyper-parameter*  $\mu > 0$  is adjusted to set the relative weight of the priors. In addition to minimizing the cost  $f(x)$ , the sought image  $x^+$  is strictly constrained to be non-negative and normalized (the sum of the pixels being equal to 1).

For the regularization term, we chose a relaxed version of the total variation criterion (Rudin, Osher & Fatemi 1992; Strong & Chan 2003), which enforces edge-preserving smoothness (Charbonnier et al. 1997), and that was found by Renard et al. (2011) to be the most effective for a large variety of astronomical objects:

$$f_{\text{prior}}(x) = \sum_{i,j} \sqrt{(x_{i+1,j} - x_{i,j})^2 + (x_{i,j+1} - x_{i,j})^2 + \varepsilon^2}, \quad (33)$$

with  $x$  the image and  $i, j$  the pixel indexes ( $\varepsilon > 0$  is a small value to have a differentiable prior term).

#### 3.4.1 Practical implementation

Once the regularization is defined, MiRA takes as input (i) the data, (ii) an optional initial estimate for the image – assumed a square of  $N \times N$  pixels – (iii) the pixel size  $\delta\theta$ , (iv) the hyper-parameter  $\mu$ , and (v) the maximum number of iterations. MiRA stops once the convergence criterion is fulfilled or the maximum number of iterations is reached. It then outputs a reconstructed image.

The image lateral size is  $\Omega = N \delta\theta$ . It provides a strict constraint which limits the support of the restored object and strongly impacts on the reconstruction process. As we want to have as few constraints as possible for the reconstruction, we chose an image size significantly larger than that of the object. In the present work,  $\Omega$  was set to be 40 mas, roughly 2.5 times the object size. The pixel size should sample the maximum angular resolution in the Nyquist–Shannon sense, i.e.  $\delta\theta < \lambda/(2B_{\text{max}})$ , with  $B_{\text{max}}$  the maximum projected baseline length. However, it was found that to make

image comparison of point-like structures reliable, a much smaller value had to be used:  $\delta\theta \simeq \lambda/(12B_{\text{max}})$ . By combining the above constraints and taking into account that the maximum baseline of the configurations in Fig. 2 is  $B_{\text{max}} = 144$  m, we adopted  $N = 160$  and  $\delta\theta = 0.25$  mas.

The only remaining parameters in the reconstruction are the (optional) initial image estimate, the number of iterations, and the value of the hyper-parameter. Their joint management is described in the following subsection.

#### 3.4.2 Tuning the hyper-parameter $\mu$

For the phase referencing reconstructions, MiRA is called without an initial image estimate, which amounts to starting with a random guessing image whose pixels are drawn following an independent uniform law. For the phase closure restorations, the initial image was a quick reconstruction from the corresponding phase referencing observation<sup>6</sup> with a large value of  $\mu$ . Because of the strong level of regularization, this image is a highly blurred version of the true image  $z$ . Other procedures could be devised to obtain the starting image for phase closure, such as a short image recover without any phase information, but this aspect is not important for the goal of this study, to wit, devise a method to assess the quality of final reconstructed images. The first restoration step (with or without initial guess) is performed for 300 iterations.

The image reconstruction process then follows a cascade of calls<sup>7</sup> to MiRA, where  $\mu$  is reduced by a constant factor in each call. The intermediary restored image output in each step is used as the image estimate for the next call. The total number of calls in the cascade is five and seven, respectively, for PhR and PhC.<sup>8</sup> MiRA normally achieves convergence before the maximum number of iterations is reached. In the PhC case, the initial image for the next MiRA call was obtained by soft-thresholding the output of the previous call at 5 per cent of its maximum.<sup>9</sup>

A limitation of the previous method is that convergence can be achieved for different values of  $\mu$ . Furthermore, no objective criterion for setting  $\mu$  is available. In this work two approaches were followed to identify the best  $\mu$ . Initially, reconstructions were conducted for different values of  $\mu$ , spanning logarithmically from  $10^4$  to  $10^{-3}$ . In the first approach, a human panel was asked to select the reconstructed image that most resembled the true image  $z$ , therefore determining the value of  $\mu$ . In the second approach, the metrics selected in Section 2.1.9 were used. In our approach, the number of free parameters is kept to a minimum. In particular, we assume  $\alpha = 1$ ,  $\beta = 0$ , a matching PSF  $h = \delta$  (a

<sup>6</sup> 2TPhR for 3TPhC, 3TPhR for 4TPhC, and 4TPhR for 6TPhC.

<sup>7</sup> Each using 1000 iterations.

<sup>8</sup> The two extra steps in the PhC case are necessary for better convergence and to properly centre the image in the FOV.

<sup>9</sup>  $x_{k+1} = \max(0, x_k - 0.05 \cdot \max(x_k))$ , with  $x_k$  the recovered image in step  $k$ . This approach was required because of the non-convex nature of PhC image reconstruction. The algorithm frequently converges to local minima.

Dirac function), and an effective PSF  $h_{\text{ref}} = G_{\sigma, t}$ . The only free parameters are then the Gaussian  $G$  standard deviation  $\sigma$  and the translation  $t$ . The translation is only relevant for the closure phase case, where the object position cannot be determined from the data. Furthermore, the translation can be implemented in either  $h$  or  $h_{\text{ref}}$ . For simplicity it was implemented in  $h_{\text{ref}}$ . The translation is not relevant for this work and will not be discussed further. The  $\sigma$  is the only parameter of the metric expressing the effective resolution. Other functions (e.g. Moffat functions) could be used, but as long as they reflect the shape of a PSF (characterized by a given width) the effect is not significant, because the metrics are summing over the convolved pixels of the images. By reducing the number of free parameters, this practical implementation has the further advantage of not defining a priori a given resolution for the reference image  $y$ , which could bias the results. Instead, it is a free parameter of the metric, that can be analysed later. The restored image  $x$  is resampled to the grid of the reference image  $y$ . Then, each metric was evaluated in the 2D parameter space  $(\mu, \sigma)$ , with  $\sigma$  spanning from 0 to 0.5 mas.<sup>10</sup> The minimum of the metric would then determine  $\mu$ .

## 4 RESULTS AND DISCUSSION

### 4.1 Reconstructed images

We produced 18 mock observations of the three reference images of Fig. 1 in all aforementioned array and phase scenarios. Images were restored from the corresponding interferometric data, stopping at 15 different levels of regularization, logarithmically ranging between  $10^4$  and  $10^{-3}$ . The procedure was repeated twice, in order to create three sets of simulations and image reconstructions. Some examples of the 810 restored images are illustrated in Figs 3–5. The full sets of recovered images are available at the JMMC website.<sup>11</sup>

Fig. 3 corresponds to restored images of the stellar cluster, Fig. 4 to the YSO, and Fig. 5 to the stellar photosphere. For the former, the first column lists images obtained when  $\mu = 10^4$ , the second column to  $\mu = 10$ , and the third column to  $\mu = 10^{-3}$ ; for the YSO, the first column corresponds to  $\mu = 10^4$ , the second column to  $\mu = 3$ , and the last column to  $\mu = 10^{-3}$ ; finally, for the stellar photosphere,  $\mu = 10^4$  in the first column,  $\mu = 300$  in the middle column, and  $\mu = 10^{-3}$  in the last column. The rows are organized as follows: the phase cases alternate between PhR and PhC, and the number of telescopes increases from top to bottom – two, three, four, and six telescopes (respectively 2T, 3T, 4T, and 6T) – so as to get the scenarios 2TPhR, 3TPhC, 3TPhR, 4TPhC, 4TPhR, and 6TPhC.

### 4.2 Observational scenarios

The quality of the images changes according to the observational scenarios considered (2T, 3T, 4T and 6T, and PhR or PhC) and their respective  $uv$ -coverages. This is essentially related to the  $uv$ -coverage of the data and the amount of phase information. It is not the goal of the present study to compare phase referencing with phase closure (and the data presented do not allow us to draw conclusions), but to present a wide variety of situations in image reconstruction to successfully test merit functions.

<sup>10</sup> For  $\sigma = 0$ , the image is only shifted as expected from the analytic convolution. Because PhC does not keep the absolute position of the objects (Monnier 2007),  $h_{\text{ref}}$  included a positional displacement  $t = (t_1, t_2)$ . This displacement was found by an iterative process that minimized the metric as a function of the displacement.

<sup>11</sup> Available at <http://oidb.jmmc.fr/collection.html?id=gomes2016>.

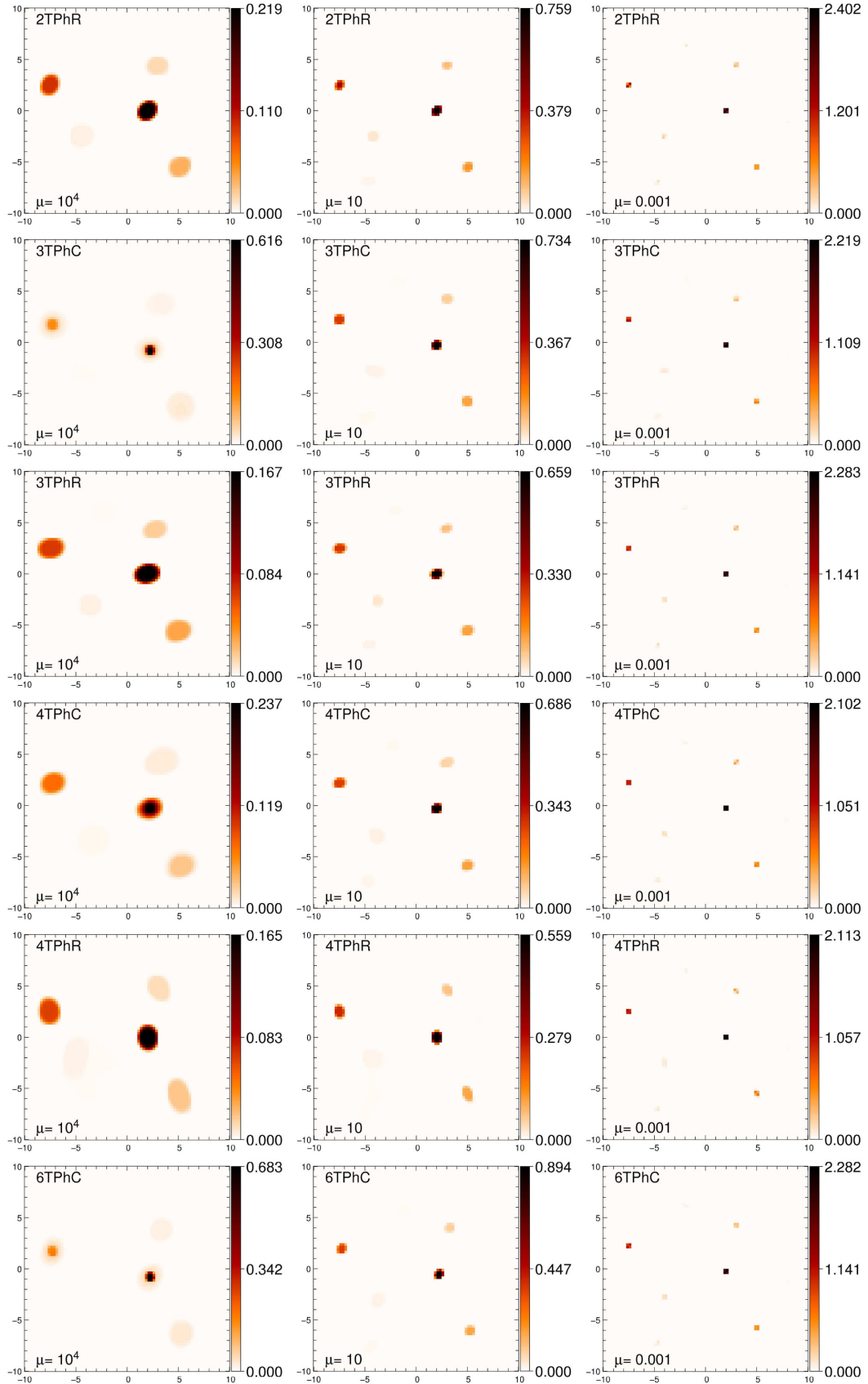
### 4.3 Effect of the level of regularization on the image reconstruction

Concerning the reconstructions and levels of regularization (Figs 3–5), it is noticeable that all restored images become sharper as the level of regularization is decreased, that is, as more weight is given to the data. However, below a certain level of  $\mu$  – which depends on the object and telescopes+phase configuration – no visible effect on the shape and surface flux of the stellar cluster is seen, because the stars (point-like unresolved source objects) become confined to one pixel. This is not the case for objects with extended/resolved structures, such as the YSO and the stellar photosphere, where reducing the regularization below a certain level introduces reconstruction artefacts and noticeably degrades the quality of the image. For instance, in the YSO, for the highest tested level of regularization ( $\mu = 10^4$ ) all images are blurred, with the central star attached to the disc. When  $\mu = 3$ , the disc is nicely restored in all configurations, with the central star separated from it. For  $\mu = 10^{-3}$ , only the 3TPhC configuration yields a well-restored image. The configurations 2TPhR, 3TPhR, 4TPhC and 4TPhR exhibit disrupted discs, full of artefacts coming out of the reconstruction process, and the 6TPhC scenario produces an image where the disc, although intact, is very irregular. In the stellar photosphere, when  $\mu = 10^4$ , only the phase closure cases produce well enough restored images, with the most prominent spot visible. When  $\mu = 300$ , the 3TPhC and the 6TPhC cases yield images where the three spots are identifiable, but all other configurations produce discs full of restoration artefacts. For  $\mu = 10^{-3}$ , the 3TPhC and 6TPhC produce well enough restored images, with two and three spots identifiable, respectively, in the former and the latter configurations. In the remainder of the scenarios, the image is not properly restored – the disc is not produced, and the algorithm gives rise solely to restoration artefacts distributed in a circular configuration.

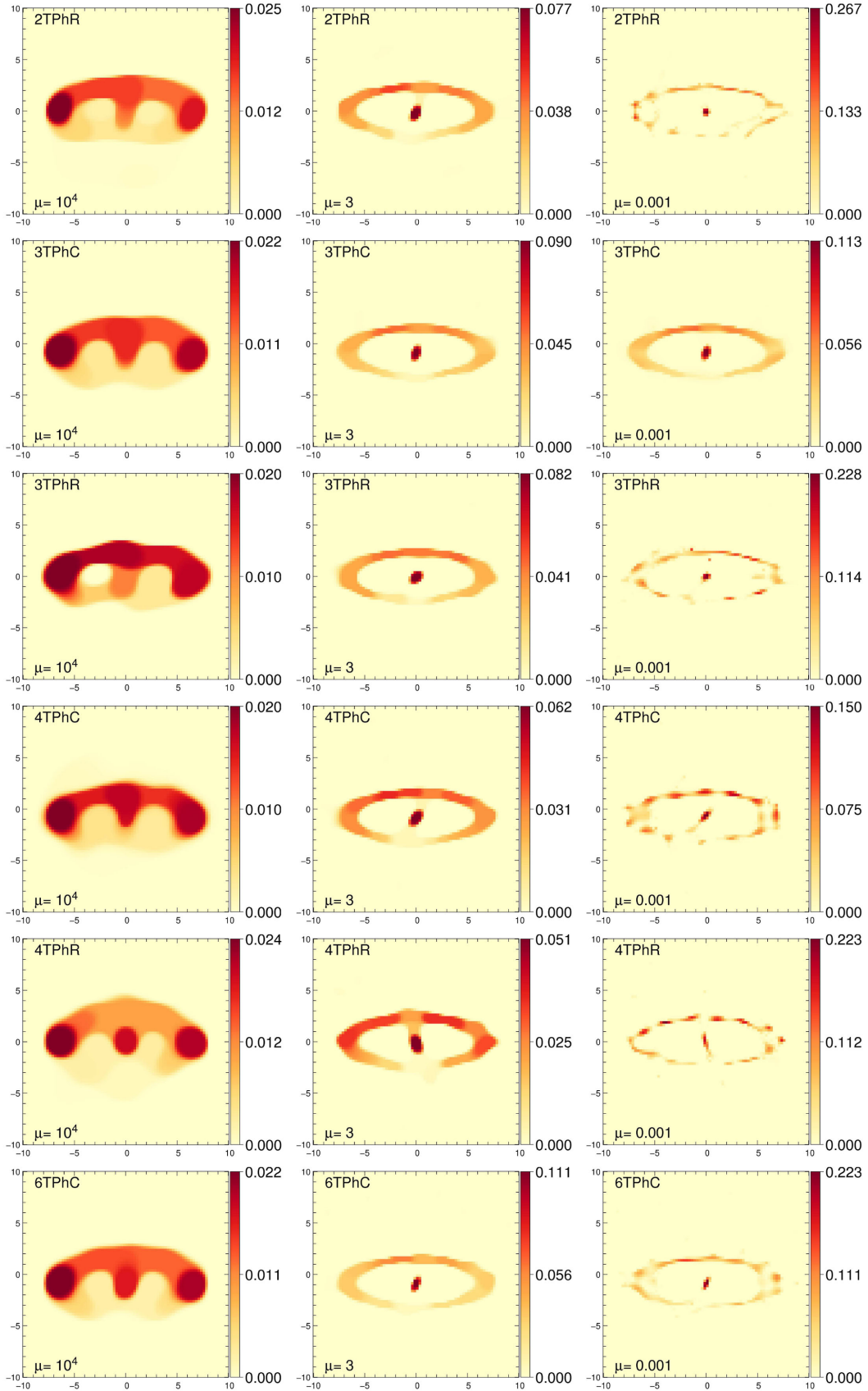
### 4.4 Human determination of the hyper-parameter

Table 2 presents the average and standard deviation of the regularization hyper-parameter  $\mu$  determined by the human panel, for each object and configuration. The value of  $\mu$  for the stellar photosphere is much larger than for the stellar cluster, which in turn is larger than that for the YSO. For a given object,  $\mu$  varies across configurations, without any specific pattern.

The values of  $\mu$  determined by human selection correspond to images that were fed to selected merit functions (see Section 2.1). The  $h_{\text{ref}}$  width is a remaining free parameter. We present in Table 3 the values of the Gaussian  $\sigma$  that minimize the metric for the human determined  $\mu$ . These values were obtained by computing the statistics for 12 realizations in each object and observational scenario. The  $\sigma$  values are of the order of 0.2 mas, which corresponds to a full width at half-maximum of about 0.5 mas. This should be compared to the angular resolution of the interferometer, which is around 3 mas, and to the reference images pixel size of 0.25 mas. Clearly the image reconstruction achieves a significant level of super-resolution, which is limited by the pixel size of the reconstructed images. This result might appear puzzling at first sight, but angular resolution is a sophisticated concept that cannot be fully enclosed in a simple Rayleigh-like criterion (e.g. den Dekker & van den Bos 1997). Because we have prior information (enforced by the regularization and positivity of the solution), a reasonable SNR and relatively smooth objects, it is expected that the image reconstruction achieves significant super-resolution.

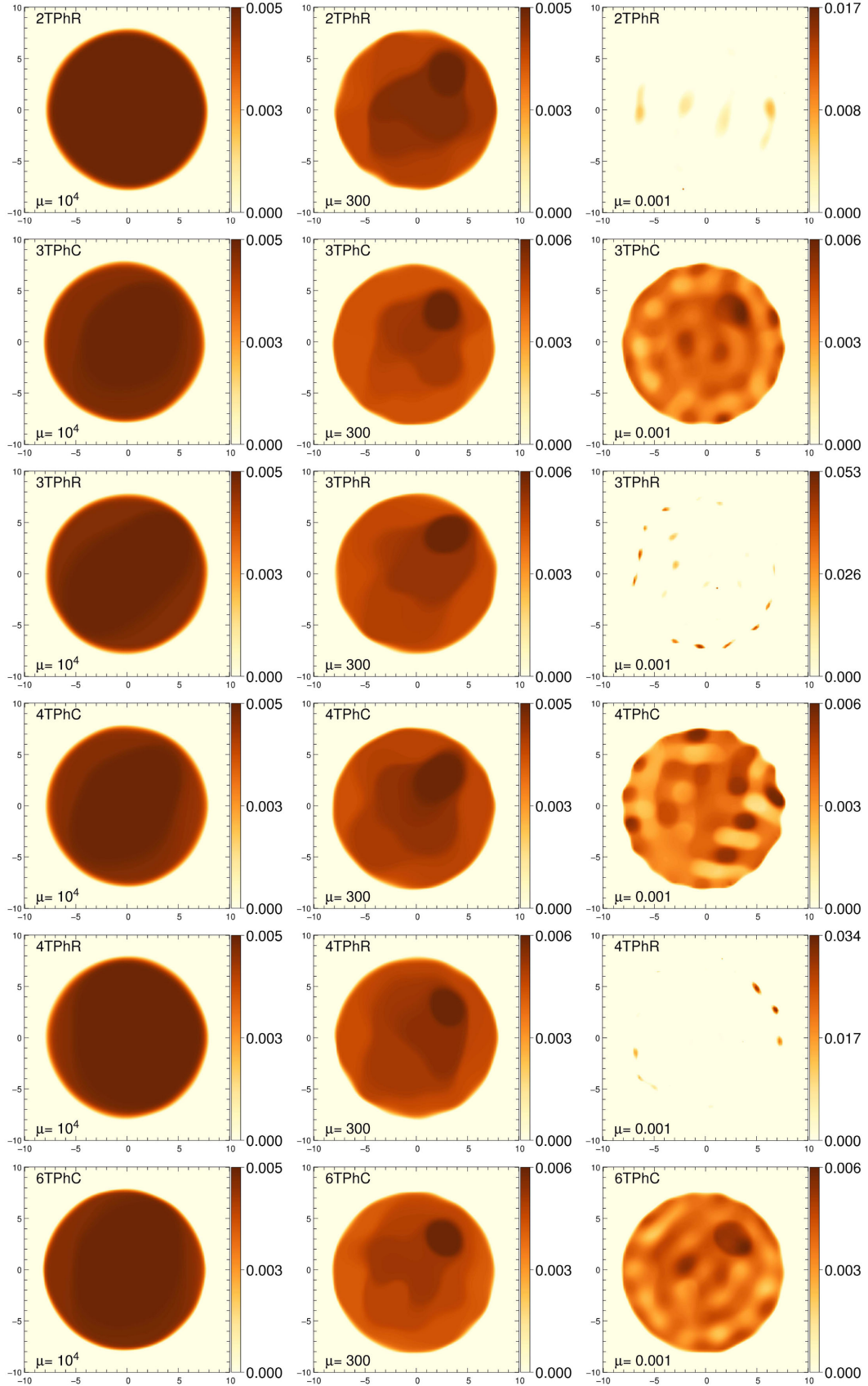


**Figure 3.** Examples of image reconstructions for the stellar cluster. Each column corresponds to a different level of regularization, and every row matches a different configuration of the synthetic observations. The lateral image size is 20 mas.



**Figure 4.** Same as in Fig. 3, but for the YSO.





**Figure 5.** Same as in Figs 3 and 4, but for the stellar photosphere.



**Table 2.** Value of the hyper-parameter  $\mu$  obtained by the human panel. The given values are the median of the values chosen by the experts, while the first and third quartiles are indicated between brackets.

Object	2TPhR	3TPhC	3TPhR	4TPhC	4TPhR	6TPhC
Stellar cluster	10 ( $\frac{30}{3}$ )	30 ( $\frac{150}{3}$ )	3 ( $\frac{30}{3}$ )	3 ( $\frac{30}{3}$ )	10 ( $\frac{30}{10}$ )	10 ( $\frac{150}{10}$ )
YSO	1 ( $\frac{1}{1}$ )	0.1 ( $\frac{1}{0.001}$ )	3 ( $\frac{10}{3}$ )	3 ( $\frac{3}{1}$ )	3 ( $\frac{3}{1}$ )	3 ( $\frac{10}{3}$ )
Stellar photosphere	300 ( $\frac{1000}{100}$ )	300 ( $\frac{1000}{100}$ )	300 ( $\frac{1000}{100}$ )	1000 ( $\frac{1000}{100}$ )	300 ( $\frac{300}{300}$ )	100 ( $\frac{300}{10}$ )

**Table 3.** Mean values of the  $h_{\text{ref}}$   $\sigma$  for the synthesized objects, observational scenarios and merit functions. The numbers between parenthesis correspond to the standard error of the mean on the last digit.

	$\sigma/\text{mas}$					
Metric	2TPhR	3TPhC	3TPhR	4TPhC	4TPhR	6TPhC
<i>Stellar cluster</i>						
ACC	0.14612(3)	0.1484(3)	0.1481(2)	0.1472(2)	0.1587(9)	0.1481(1)
L1N	0.14373(7)	0.1483(4)	0.1464(4)	0.1458(3)	0.1625(9)	0.1498(3)
L2N	0.14985(3)	0.1522(3)	0.1518(2)	0.1508(2)	0.1629(9)	0.1520(1)
IBC	0.15437(3)	0.1560(3)	0.1560(2)	0.1550(1)	0.1648(8)	0.15547(9)
<i>YSO</i>						
ACC	0.281(2)	0.273(4)	0.294(2)	0.281(2)	0.320(2)	0.263(2)
L1N	0.204(2)	0.191(5)	0.198(2)	0.216(4)	0.259(4)	0.207(2)
L2N	0.306(3)	0.298(4)	0.320(1)	0.301(2)	0.347(2)	0.282(2)
IBC	0.343(4)	0.333(4)	0.367(2)	0.334(2)	0.384(2)	0.305(3)
<i>Stellar photosphere</i>						
ACC	0.293(2)	0.216(3)	0.270(2)	0.255(2)	0.242(2)	0.189(4)
L1N	0.274(2)	0.198(3)	0.239(2)	0.232(2)	0.219(2)	0.166(4)
L2N	0.269(2)	0.198(2)	0.245(2)	0.233(2)	0.221(2)	0.170(3)
IBC	0.277(2)	0.201(3)	0.251(2)	0.239(3)	0.226(2)	0.173(4)

In order to check the robustness of Figs 6 to 8 to different realizations of the data, we carried out 12 simulations of the 18 synthetic observations. The statistics of the minima for the human determined  $\mu$  are presented in Table 4 (the errors in Table 3 were computed from this same data set). The standard error of the mean is very small, supporting the robustness of the results to the noise in the data set.

#### 4.5 Benchmarking the metrics

As explained in Section 3.4, a reconstructed image is a function of the final chosen  $\mu$ . Furthermore, the application of a given metric requires the convolution by  $h_{\text{ref}}$ , whose width is characterized by  $\sigma$ . In this subsection we present and discuss the results for the behaviour of the merit functions.

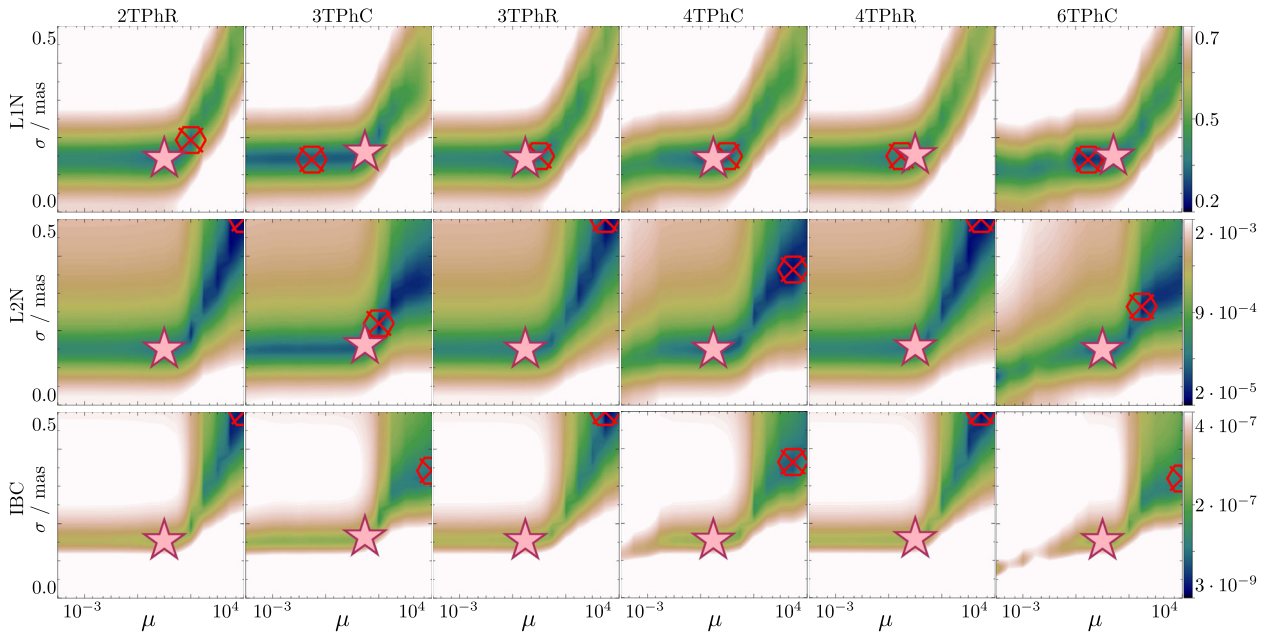
Table 4, where  $\mu$  is determined by human selection, provides an initial benchmark. The values of the quality functions show that IBC mimics the behaviour of L2N in most objects and configurations. On the one hand, this is explained by the quadratic nature of both metrics and, on the other hand, by the fact that the weighting function of IBC is the reference image itself, which makes the metric disregard pixels where the latter is zero. The failure of ACC in properly characterizing the quality of restored images in some scenarios is related to the fact that it applies a mask to the reference image before comparison, thus eliminating parts containing reconstruction artefacts that are important to determine the quality of the image. This however could be an interesting merit function when we are focused on certain parts of the image and want to eliminate others that we safely identify as artefacts of the reconstruction. For all

objects and configurations, the L1N metric appears to properly characterize the quality of the restored images.

We also conducted a systematic study of the metric behaviour as a function of  $\mu$  and  $\sigma$ . We varied  $\mu$  logarithmically between  $10^4$  and  $10^{-3}$ , and  $\sigma$  linearly between 0 and 0.5 mas. The average values of the merit functions for three realizations of the simulated observations versus  $\mu$  and  $\sigma$  are plotted in Fig. 6 (for the stellar cluster), Fig. 7 (for the YSO) and Fig. 8 (for the stellar photosphere). The top, middle and bottom rows present the results for the quality functions L1N, L2N and IBC, respectively. The columns are organized as the rows of Figs 3–5. The colour palette is inverted, such that the minima (darker colours) indicate a better agreement between the restored images and the references. All merit functions exhibit regions of minima, which is also verified in the ACC metric (not depicted). The red crossed circles point to the global minima of the panels. The pink stars are located at the position of the aforementioned values of  $\mu$  determined by human selection. The position of the corresponding  $\sigma$  was obtained by minimizing the merit function for the fixed  $\mu$ , using the `NEWUOA` algorithm (Powell 2006).

The first result is that, generally, the merit functions are reasonably convex (i.e. they depict regions with a clear minima). Overall, the effective resolution worsens with the hyper-parameter  $\mu$ , as expected (i.e. the dark regions bend towards larger values of  $\sigma$  and  $\mu$ ). This is expected because increasing  $\mu$  amounts to smooth the image.

The shape of the minima regions of Figs 6 to 8 depends on the object. In the case of the stellar cluster (Fig. 6), the minima regions exhibit a horizontal branch up to a certain level of regularization. This is compatible with the aforementioned limiting value of regularization, below which restored images present no noticeable differences in quality and the (super-)resolution becomes limited by



**Figure 6.** Average scores of the metrics L1N (top row), L2N (central row) and IBC (bottom row) for three sets of simulated observations as function of the standard deviation  $\sigma$  of  $h_{\text{ref}}$  and the level of regularization  $\mu$ . The object is the stellar cluster of Fig. 1. From left to right, the panels are organized as follows: 2TPhR, 3TPhC, 3TPhR, 4TPhC, 4TPhR, and 6TPhC. The red crossed circles correspond to global minima, while the pink stars are positioned at the human determined value of  $\mu$  and the value of  $\sigma$  that minimizes the merit function. A logarithmic and a linear scale were respectively used for  $\mu$  and  $\sigma$ .

**Table 4.** Mean values of the merit functions at the positions of  $\mu$  determined by human selection (pink stars in Figs 6–8). The scores were obtained by computing the statistics for at least 12 realizations in each scenario. The smaller the values, the better the agreement. The numbers between parenthesis correspond to the standard error of the mean of the last digit.

	2TPhR	3TPhC	3TPhR	4TPhC	4TPhR	6TPhC
<i>Stellar cluster</i>						
ACC	0.03760(9)	0.0364(2)	0.065(5)	0.060(4)	0.066(4)	0.063(3)
L1N	0.239(1)	0.231(2)	0.19(1)	0.199(9)	0.191(8)	0.195(7)
L2N	$6.08(1) \times 10^{-8}$	$5.73(6) \times 10^{-8}$	$2.3(3) \times 10^{-8}$	$2.7(2) \times 10^{-8}$	$2.0(2) \times 10^{-8}$	$2.3(2) \times 10^{-8}$
IBC	$4.76(1) \times 10^{-5}$	$4.49(4) \times 10^{-5}$	$2.0(2) \times 10^{-5}$	$2.3(2) \times 10^{-5}$	$1.8(1) \times 10^{-5}$	$2.0(1) \times 10^{-5}$
<i>YSO</i>						
ACC	0.064(7)	0.092(7)	0.067(4)	0.077(4)	0.066(3)	0.072(3)
L1N	0.254(6)	0.274(5)	0.207(9)	0.220(8)	0.200(7)	0.208(7)
L2N	$4.3(4) \times 10^{-8}$	$2.9(3) \times 10^{-8}$	$2.5(2) \times 10^{-8}$	$2.2(2) \times 10^{-8}$	$2.2(2) \times 10^{-8}$	$2.0(2) \times 10^{-8}$
IBC	$3.5(2) \times 10^{-5}$	$2.5(2) \times 10^{-5}$	$2.2(2) \times 10^{-5}$	$2.0(1) \times 10^{-5}$	$2.0(1) \times 10^{-5}$	$1.8(1) \times 10^{-5}$
<i>Stellar photosphere</i>						
ACC	0.083(6)	0.068(5)	0.074(4)	0.068(4)	0.070(3)	0.066(3)
L1N	0.24(1)	0.19(1)	0.208(8)	0.187(8)	0.201(7)	0.187(7)
L2N	$2.5(3) \times 10^{-8}$	$1.9(3) \times 10^{-8}$	$2.0(2) \times 10^{-8}$	$1.8(2) \times 10^{-8}$	$2.0(2) \times 10^{-8}$	$1.8(1) \times 10^{-8}$
IBC	$2.2(2) \times 10^{-5}$	$1.7(2) \times 10^{-5}$	$1.8(1) \times 10^{-5}$	$1.6(1) \times 10^{-5}$	$1.8(1) \times 10^{-5}$	$1.6(1) \times 10^{-5}$

the size of the pixel. A single pixel encompasses the totality of the flux emanating from a restored unresolved star lying inside of it. The value of  $\sigma \sim 0.15$  mas indicated by the branch is compatible with the pixel size of 0.25 mas. For sources with extended emission, the branch is not visible because the image degrades rapidly below a certain level of regularization (cf. Figs 4 and 5 for some examples). Nevertheless, regions of minima are also evident, the position of which largely depends on the merit function.

#### 4.5.1 L1N as the most robust metric

For L1N, the global minima typically lie well inside the limits defined by the plots. That is not the case for many L2N and IBC

observations (especially for the cluster and YSO), suggesting that if the study was extended to larger values of  $\sigma$  and  $\mu$ , the global minima would point to more blurred images. The minima valley oriented in the direction of increasing  $\mu$  and  $\sigma$  is less pronounced for L1N than for L2N and IBC. For L2N and IBC, this would indicate a better agreement between the restored and the reference images in those extreme regions of the plots, where the restored images are more blurred. This clearly shows that these metrics are biased and are not robust to over-smoothing by large values of the  $\mu$  hyper-parameter. They will consider that an image with lower ‘angular resolution’ is a better image than one with higher ‘angular resolution’. These results support L1N as the most robust of the merit functions used for the variety of cases considered.

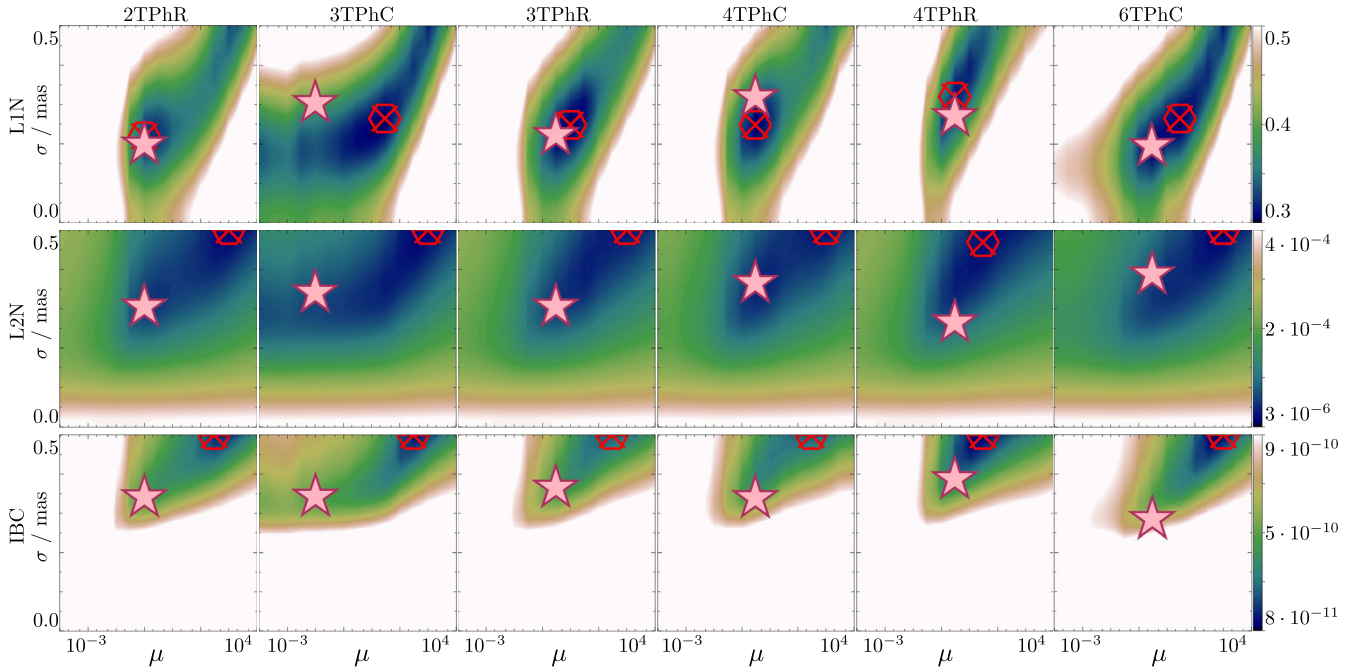


Figure 7. Same as in Fig. 6, but for the YSO.

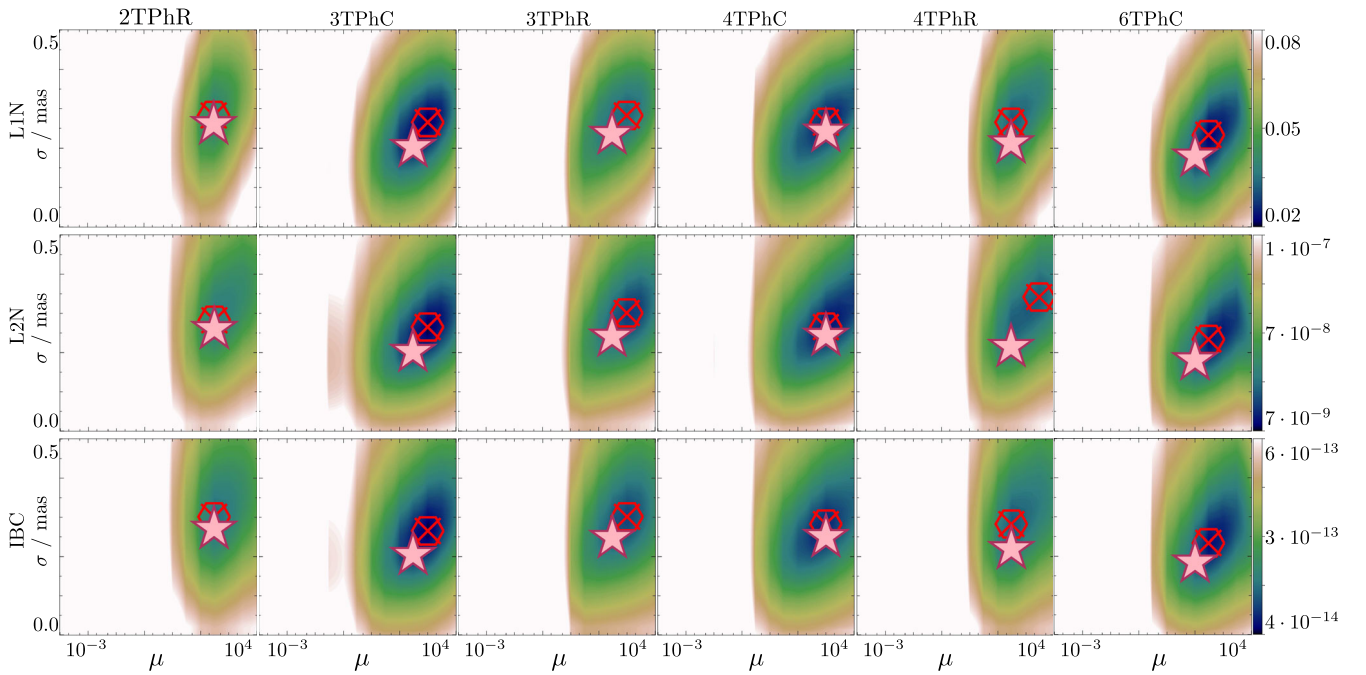


Figure 8. Same as in Figs 6 and 7, but for the stellar photosphere.

The morphology of the object has some impact on the behaviour of the metrics. The quality of extended resolved objects can be more easily assessed than that of unresolved sources. When the emitting source combines both types of objects (resolved and unresolved), the studied merit functions seem to have a harder job to evaluate the quality of the restored images. The great imbalance in intensity between the central star and the surrounding disc might explain the differences in quality.

#### 4.6 Automatic image quality assessment

The distance between the pink stars (minima obtained from human selection) and the circled red crosses (global minima) in Figs 6 to 8 indicates how well a given merit function translates the human perception of a ‘good’ restored image. In this regard, LIN is clearly the best of all studied metrics, as it is the only one where both beacons lie close together for the typology of objects and observing configurations.



This is not as well verified with the other metrics, being IBC the less robust of the tested merit functions. In the case of the stellar photosphere (Fig. 8), all metrics behave similarly.

Since we are truncating the intervals of  $\sigma$  and  $\mu$ , those distances most probably would increase in the cases where the global minima lie at extreme points of the plots.

These results open the possibility of automatic image quality assessment, thus removing human intervention in the process.

## 5 CONCLUSIONS AND FUTURE DEVELOPMENTS

This article addresses the question: what is the best metric to assess the quality of a reconstructed image?

Several merit functions are considered in the realistic context of the VLTI and using the MiRA image reconstruction software.

A semi-automatic pipeline is developed to reconstruct images, with the only human intervention being the determination of the final value of the hyper-parameter  $\mu$ . It is found that the image reconstruction process outputs images with an effective angular resolution, characterized by a Gaussian, whose standard deviation  $\sigma$  is significantly smaller than an equivalent Rayleigh-like criterion, based on the maximum baseline. Hence, a certain amount of super-resolution is achievable thanks to the constraints imposed by a regularized image reconstruction algorithm.

In order to cope with the mismatch between the effective resolution of the restored image and that of the simulated object, we advocate that convolution by an effective PSF is mandatory for proper image quality assessment. This effective PSF can be further used to compensate for image shift, which is unavoidable when image reconstruction is performed from power-spectrum and phase closure data.

Of all the merit functions considered, the  $\ell_1$ -norm is the most robust. The commonly used Interferometric Imaging Beauty Contest quadratic metric is biased, considering as best images those with higher smoothing (or hyper-parameter  $\mu$ ), and not fully exploiting the effective angular resolution of the data and image reconstruction process.

By minimizing the  $\ell_1$ -norm over the  $\mu$  and  $\sigma$  parameter space, it is possible to implement automated image quality assessment.

Based on this work, several developments are foreseen, the most obvious of which being algorithm comparison with the  $\ell_1$ -norm and proper convolution. The most ambitious is automated image reconstruction. To achieve this goal, two aspects must be addressed: (i) the determination of an initial image for the reconstruction algorithm (for phase closure only), and (ii) the determination of the final  $\mu$  in the reconstruction. The second aspect is clearly the most difficult. It opens the requirements for image reconstruction algorithms to output tables of images for different levels of regularization, allowing the end-user to determine the final values of  $\mu$ .

An important aspect is to identify the situations where phase referencing or phase closure are the best options for imaging. This choice is now possible with the *GRAVITY* and *PIONIER* instruments. Its study requires the inclusion of other ingredients not addressed in the present article, such as (i) compatible *uv*-coverages, (ii) noise models taking into account photon and detector statistics (e.g. Tatulli & Chelli 2005) and/or light splitting between telescopes (e.g. Gordon & Buscher 2012), and (iii) a span of SNRs.

## ACKNOWLEDGEMENTS

The research leading to these results has received funding from the European Community's Seventh Framework Programme, under Grant Agreements 226604 and 312430 (OPTI-CON), as well as from Fundação para a Ciência e Tecnologia grants SRFH/BD/44282/2008, PTDC/CTE-AST/116561/2010, UID/FIS/00099/2013 and COMPETE FCOMP-01-0124-FEDER-019965. This research has made use of the Jean-Marie Mariotti Center ASPRO service.<sup>12</sup> All the analysis was done with YORICK, a free interactive data processing language written by David Munro.<sup>13</sup>

The authors would like to thank P. Andrade, N. Anugu, T. Armstrong, J. Ascenso, P. Berlioz-Arthaud, W.-J. de Wit, G. Duvert, S. Ertel, D. Faes, H. Hummel, P. Kervella, J. Kluska, Y. Kok, M. Langlois, J. Léger, M. Loupias, D. Mourard, M. Ozon, E. Pinho, N. Scott, M. Silva, F. Soulez, M. Tallon, I. Tallon-Bosc, and N. Verrier, for their valuable contribution in the poll aiming at determining the best levels of regularization for image reconstruction.

The authors would also like to thank the referee for his/her insightful comments.

## REFERENCES

- Arsenault R. et al., 2004, Proc. SPIE Conf. Ser. Vol. 5490, Advancements in Adaptive Optics. SPIE, Bellingham, p. 47
- Baldwin J. E. et al., 1996, A&A, 306, L13
- Baron F. et al., 2012, Proc. SPIE Conf. Ser. Vol. 8445, The 2012 Interferometric Imaging Beauty Contest. SPIE, Bellingham, p. 84451E
- Baron F. et al., 2014, ApJ, 785, 46
- Benisty M. et al., 2011, A&A, 531, A84
- Benson J. A. et al., 1997, AJ, 114, 1221
- Berger J.-P. et al., 2012, A&AR, 20, 53
- Buscher D. F., 2015, Practical Optical Interferometry: Imaging at Visible and Infrared. Cambridge Univ. Press, Cambridge
- Charbonnier P., Blanc-Féraud L., Aubert G., Barlaud M., 1997, IEEE Trans. Im. Proc., 6, 298
- Che X. et al., 2011, ApJ, 732, 68
- Cotton W. et al., 2008, Proc. SPIE Conf. Ser. Vol. 7013, 2008 Imaging Beauty Contest. SPIE, Bellingham, p. 70131N
- Coudé du Foresto V., Ridgway S., Mariotti J.-M., 1997, A&A, 121, 379
- Delplanck F., 2008, New Astron. Rev., 52, 199
- den Dekker A. J., van den Bos A., 1997, J. Opt. Soc. Am. A: Optics, Image Sci. Vision, 14, 547
- Eisenhauer F. et al., 2008, Proc. SPIE Conf. Ser. Vol. 7013, Optical and Infrared Interferometry. SPIE, Bellingham, p. 70132A
- Eisenhauer F. et al., 2011, The Messenger, 143, 16
- Filho M. E. et al., 2008a, Proc. SPIE Conf. Ser. Vol. 7013, Optical and Infrared Interferometry. SPIE, Bellingham, p. 70131F
- Filho M. E. et al., 2008b, Proc. SPIE Conf. Ser. Vol. 7013, Optical and Infrared Interferometry. SPIE, Bellingham, p. 70133Z
- Glindemann A., 2011, Principles of Stellar Interferometry. Springer-Verlag, Berlin, Heidelberg
- Gomes N., 2016, Imaging with the VLTI, PhD thesis
- Goodman J. W., 1985, Statistical Optics. Methods. John Wiley and Sons, New York
- Gordon J. A., Buscher D. F., 2012, A&A, 541, A46
- Hillen M., Kluska J., Le Bouquin J.-B., Van Winckel H., Berger J.-P., Kamath D., Bujarrabal V., 2016, A&A, 588, L1
- Jennison R. C., 1958, MNRAS, 118, 256
- Kloppenborg B. K. et al., 2015, ApJS, 220, 14
- Kluska J. et al., 2014, A&A, 564, A80

<sup>12</sup> Available at <http://www.jmmc.fr/aspro>

<sup>13</sup> Available at <http://yorick.sourceforge.net/>

- Kraus S. et al., 2014, in Rajagopal J. K., Creech-Eakman M. J., Malbet F., eds, Proc. SPIE Conf. Ser. Vol. 9146, The science case for the Planet Formation Imager (PFI). SPIE, Bellingham, p. 914611
- Lawson P. R. et al., 2004, Proc. SPIE Conf. Ser. Vol. 5491, New Frontiers in Stellar Interferometry. SPIE, Bellingham, p. 886
- Lawson P. R. et al., 2006, Proc. SPIE Conf. Ser. Vol. 6268, Advances in Stellar Interferometry, SPIE, Bellingham, p. 62681U
- Le Bouquin J.-B., Lacour S., Renard S., Thiébaud E., Merand A., Verhoelst T., 2009, A&A, 496, L1
- Le Bouquin J.-B. et al., 2011, A&A, 535, A67
- Malbet F., Kern P. Y., Berger J.-P., 2006, Proc. SPIE Conf. Ser. Vol. 6268, Advances in Stellar Interferometry. SPIE, Bellingham, p. 62680Y
- Malbet F. et al., 2010, Proc. SPIE Conf. Ser. Vol. 7734, Optical and Infrared Interferometry II. SPIE, Bellingham, p. 77342N
- Millour F., Meilland A., Chesneau O., Stee Ph., Kanaan S., Petrov R., Mourard D., Kraus S., 2011, A&A, 526, A107
- Monnier J. D., 2007, New Astron. Rev., 51, 604
- Monnier J. D. et al., 2014a, Proc. SPIE Conf. Ser. Vol. 9146, Optical and Infrared Interferometry IV. SPIE, Bellingham, p. 91461Q
- Monnier J. D. et al., 2014b, Proc. SPIE Conf. Ser. Vol. 9146, Optical and Infrared Interferometry IV. SPIE, Bellingham, p. 914610
- Mourard D. et al., 2015, A&A, 577, A51
- Pauls T. A., Young J. S., Cotton W. D., Monnier J. D., 2005, PASP, 117, 1255
- Petrov R. G. et al., 2007, A&A, 464, 1
- Pety J., Gueth F., Guilloteau S., 2001b, ALMA Memo, 386, 10
- Powell M. J. D., 2006, Large-Scale Nonlinear Optimization. Springer-Verlag, Berlin, p. 255
- Renard S., Thiébaud É., Malbet F., 2011, A&A, 533
- Rudin L. I., Osher S., Fatemi E., 1992, Phys D Nonlinear Phenomena, 60, 259
- Schöller M., 2007, New Astron. Rev., 51, 628
- Schutz A., Vannier M., Mary D., Ferrari A., Millour F., Petrov R., 2014, A&A, 565, A88
- Strong D., Chan T., 2003, Inverse Problems, 19, S165
- Tatulli É., Chelli A., 2005, Opt. Soc. Am J, 22, 1589
- Tatulli É., Blind N., Berger J. P., Chelli A., Malbet F., 2010, A&A, 524, A65
- ten Brummelaar T. A. et al., 2005, ApJ, 628, 453
- Thiébaud É., 2008, Proc. SPIE Conf. Ser. 7013, 70131I
- Thiébaud É., 2013, EAS Publ., 59, 157
- Thompson A. R., Moran J. M., Swenson G. W., 2001, Radio Interference, in Interferometry and Synthesis in Radio Astronomy, 2nd edn. Wiley-VCH Verlag, GmbH
- Vincent F. H., Paumard T., Perrin G., Mugnier L., Eisenhauer F., Gillessen S., 2011, MNRAS, 412, 2653
- Wang Z., Bovik A. C., 2002, IEEE Signal Proc. LET, 9, 81
- Wang Z. et al., 2004, IEEE T Image Process, 13, 600

This paper has been typeset from a  $\text{\LaTeX}$  file prepared by the author.