

The XMM-Newton serendipitous survey. XI. The fifth XMM-Newton serendipitous source catalogue

The 5XMM catalogue

N. A. Webb^{1,*}, I. Traulsen^{2,**}, M. Coriat¹, F. J. Carrera³, L. Michel⁴, A. Nebot Gómez-Morán⁴, J. Ballet⁵, M. J. Page⁶, A. Ruiz⁷, M. Watson⁸, M. Freyberg⁹, R. M. Batalha^{5,11}, D. Bogensberger⁵, M. T. Ceballos³, S. Chakraborty^{5,12}, I. de la Calle Perez¹³, I. Georgantopoulos⁷, A. Georgakakis⁷, N. P. M. Kuin⁶, G. Lamer², F. Mernier¹, C. Motch⁴, G. Mountrichas³, A. M. Pires^{10,2}, E. Pouliaxis⁷, E. Quintin¹, D. Rawat⁴, A. Schwoppe², H. Tranin¹, J. Vicente¹⁴, and R. Webbe¹

¹ Université de Toulouse, CNES, CNRS, IRAP, 9 avenue du Colonel Roche, 31028 Toulouse, France

² Leibniz-Institut für Astrophysik Potsdam (AIP), An der Sternwarte 16, 14482 Potsdam, Germany

³ Instituto de Física de Cantabria (CSIC-Universidad de Cantabria), Avenida de los Castros, 39005 Santander, Spain

⁴ Université de Strasbourg, CNRS, Observatoire astronomique de Strasbourg, UMR 7550, 67000 Strasbourg, France

⁵ Université Paris-Saclay, Université Paris Cité, CEA, CNRS, AIM, F-91191 Gif-sur-Yvette Cedex, France

⁶ Mullard Space Science Laboratory, University College London, Holbury St Mary, Dorking, Surrey RH5 6NT, UK

⁷ IAASARS, National Observatory of Athens, I. Metaxa & V. Pavlou, 15236, Greece

⁸ Department of Physics & Astronomy, University of Leicester, Leicester, LE1 7RH, UK

⁹ Max-Planck-Institut für extraterrestrische Physik, Giessenbachstraße 1, 85748 Garching, Germany

¹⁰ Center for Lunar and Planetary Sciences, Institute of Geochemistry, Chinese Academy of Sciences, 99 West Lincheng Rd., 550051 Guiyang, China

¹¹ Observatório Nacional, Rua General José Cristino, 77 – Bairro Imperial de São Cristóvão, Rio de Janeiro 20921-400, Brazil

¹² Science and Technology Institute, Universities Space and Research Association, Huntsville, AL 35805, USA

¹³ TELESPAZIO for the European Space Agency. European Space Astronomy Center (ESAC-ESA). Madrid. 28691. Spain

¹⁴ Starion Group for the European Space Agency. European Space Astronomy Center (ESAC-ESA). Madrid. 28691. Spain

Received September 30, 20XX

ABSTRACT

Aims. In order to make it easy to access good quality, homogeneous information on X-ray, ultra-violet and optical sources and detections made with the *XMM-Newton* observatory, we provide a new version of the X-ray catalogue, 5XMM-DR15, using improved calibration and software, compared to previous versions. We also include new data that allow the user to choose homogeneous populations of sources, carry out spectral studies, find sources of similar luminosities, identify sources that vary on the short- and long-term and carry out multi-wavelength studies without the need for cross-correlating sources with multi-wavelength catalogues.

Methods. We have improved the effective areas used for the EPIC cameras, as well as the astrometry of the MOS cameras. For the first time we provide a single catalogue of stacked sources. To do this we have developed a new stacking procedure which has allowed us to increase the signal to noise of the detected sources, reach deeper fluxes and include all of the EPIC data. We have developed new methods and software which have allowed us to include new quantities including the *XMM-Newton Optical Monitor* counterparts to the X-ray sources, a measure of the long-term X-ray and optical variability when sources have been observed multiple times, source classification for both the OM and X-ray sources, WISE and Gaia counterparts where they exist, along with the Gaia distance estimate and further tables with other multi-wavelength data, spectral fitting parameters for absorbed power laws and photometric redshifts for extragalactic sources. Automated and manual screening indicate the reliability of the sources in the catalogue.

Results. The 5XMM-DR15 catalogue contains 818 656 unique X-ray sources detected in the 0.2–12.0 keV band. They have 2 578 752 individual detections or upper limits. Some sources are observed as many as 98 times with *XMM-Newton*. The catalogue covers ~3.5% of the sky. 87039 X-ray sources have an OM counterpart, 51624 X-ray sources have a Gaia DR3 counterpart, and 34279 have a Gaia parallax. 12330 sources are variable within an observation (5σ significance) and 41187 have a long-term variability of greater than a factor five. Spectra were extracted and fitted successfully with a power law for 234 133 sources, and 154 734 sources have a photometric redshift. 556 337 X-ray sources are classified as active galactic nuclei (AGN), 119 661 as stars, 26100 as Galactic X-ray binaries, 1276 as cataclysmic variables, 49969 as background AGN, 22732 as extra-galactic X-ray binaries and 42581 as extended sources. The catalogue is available to be downloaded as a FITS file, with and without the (non-)detections from each observation, and can also be accessed through dedicated servers, which also provide images, lightcurves, spectra and other value-added products.

Key words. Catalogs – Astronomical data bases – Surveys – X-rays: general

1. Introduction

Catalogues of homogeneous products provide an efficient way to access information on detections and sources made with a tele-

* Corresponding author: Natalie.Webb@irap.omp.eu

** Corresponding author: itraulsen@aip.de

scope, without having to reduce and analyse the observations, which can sometimes run to very large volumes of data. Catalogues can also be used to provide quick access to data products (fluxes, variability, spectral fits, images, etc), find new objects, carry out population studies or carry out cross correlation for multi-wavelength studies, for example.

10 *XMM-Newton* (Jansen et al. 2001) was launched over twenty six years ago on 10 December 1999. It has the largest effective area of any X-ray satellite (Ebrero 2019) thanks to the three X-ray telescopes aboard, each with ~ 1500 cm² of geometric effective area at 1.5 keV. This fact, coupled with the large field of view (FOV) of 30' diameter, means that a single pointing with the mean duration in the catalogue of 33 ks detects 70-75 serendipitous X-ray sources. To date, four major catalogue versions have been produced, following the re-reduction of all of the data available at the time of the catalogue production, using the best software and calibration available at that time. These were called 1XMM (33026 detections), 2XMM (2XMMi, 289083 detections Watson et al. 2009), 3XMM (3XMM-DR4, 531261 detections Rosen et al. 2016) and 4XMM (4XMM-DR9 810795 detections Webb et al. 2020), with incremental versions of these catalogues indicated by successive data releases, denoted -DR in association with the catalogue number. Since 3XMM-DR7 we have also released a catalogue made from stacking overlapping detections, where the first one was called 3XMM-DR7s, to indicate the stacked version of the catalogue (Traulsen et al. 2019). Since then we have released a further six versions of the stacked catalogue, with the most recent version of the catalogue being 4XMM-DR14s (Traulsen et al. 2020). Stacking data achieves better sensitivity and improves source parameters, and provides direct access to long-term flux variability, where some sources have almost 3 Ms of observations. 4XMM-DR14s was built from 1751 groups drawn from 10332 observations and contains 427524 sources. There are 1.8 million individual flux measurements, as these include upper limits where no detection can be made.

40 This paper details the most recent version of the *XMM-Newton* catalogue, 5XMM-DR15. With this version, only a single catalogue is provided as all data are treated in the same way as a stack and therefore all data from 14 616 observations are included. The format is the same as in the previous stacked catalogues, information on the source is provided and then also for each (non-)detection of the source in subsequent lines. The stacking procedure has been significantly improved, allowing to reach deeper fluxes than in previous versions, and the new procedure is described in Sec. 4. 5XMM-DR15 also provides new value-added products. For the first time *XMM-Newton* Optical Monitor (XMM-OM, Mason et al. 2001) measurements are included when an XMM-OM counterpart to the X-ray source is identified. Also available for the XMM-OM counterparts is a measure of the long-term variability for every ultraviolet / optical XMM-OM passband in which it has been observed multiple times. Further, source classification is provided for each XMM-OM counterpart, see Sec. 9. Also provided are the WISE and Gaia counterparts where they exist, along with the Gaia distance estimate and link to further tables with other multi-wavelength data, see Sec. 12. Other new information in 5XMM includes spectral fitting parameters for absorbed power laws, see Sec. 6, long-term variability of the X-ray source, see Sec. 5 and the photometric redshift for extragalactic sources, along with the redshift determined from spectra, when it is available, see Sec. 8. The methodology used to derive the data, as well as the content of the catalogue are described in detail in the following sections.

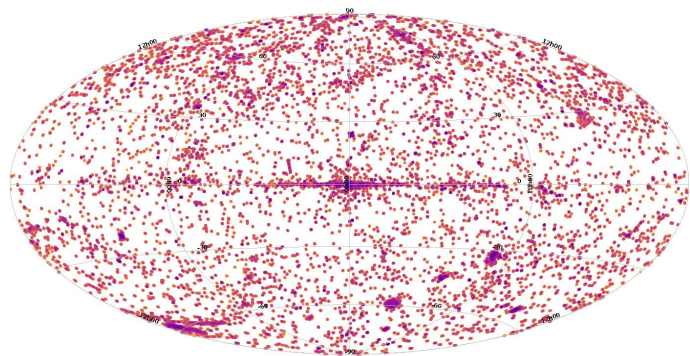


Fig. 1: Hammer-Aitoff equal area projection in Galactic coordinates of the 5XMM-DR15 fields. The colours represent the density of the observations, with yellow the lowest density to deep blue, the highest density.

2. Observations

A total of 17 554 *XMM-Newton* EPIC observations were publicly available as of 1 November 2024, of which 14 616 had EPIC event lists taken in thin, medium, or thick filter, and 14 480 had detected sources. All 14 616 observations are included in the catalogue, as the addition of these observations increased the signal to noise of the stack and therefore contributed photons to the detections, see Sec. 4. The repartition of data modes for each camera and observation can be found in Table 1. The Hammer-Aitoff equal area projection in Galactic coordinates of the 5XMM-DR15 fields can be seen in Fig. 1. Open filter data were processed but not used in the source detection stage of pipeline processing and are not included in the numbers provided above. The same *XMM-Newton* data modes were used as in 2XMM (Watson et al. 2009) and versions since and are included in Table A.1 of this paper, for convenience.

Observations that overlap by at least 1 arcmin in radius are grouped into so-called stacks for simultaneous source detection. The overlap is estimated initially from the nominal aim points of the observations. Due to the small offset between the EPIC instruments and the aim point, observations can end up in the same stack while their actual instrument footprints do not overlap. These observations with zero pixels of overlap are identified from their detection masks in a second step and treated separately. Of the 14 616 5XMM observations, 11 158 overlap with at least one other observation. They were grouped into 5303 stacks for simultaneous source detection. On the other 3 458 non-overlapping observations, source detection is performed individually.

3. Data Processing

Data processing for the 5XMM-DR15 catalogue was based on the SAS version 21 and carried out with the pipeline version 21.51¹ and the latest set of current calibration files at the time of processing (from November 2024). This new version includes a number of improvements compared to previous versions. Improvements to the EPIC (and RGS) effective areas were made using an empirical correction from MOS to pn (and RGS to pn), as well as a further correction above 3 keV to align the pn to the NUSTAR spectral fits. A correction was also included to update

¹ <https://www.cosmos.esa.int/web/xmm-newton/pipeline-configurations>

Table 1: Characteristics of the 14 616 *XMM-Newton* observations included in the 5XMM-DR15 catalogue.

Camera	Modes			Filters			Total
	Full ^a	Window ^b	Other ^c	Thin	Medium	Thick	
pn	11214	893	1777	7376	5453	1056	13894
MOS1	11373	2616	422	6657	6467	1288	14422
MOS2	11408	2495	563	6700	6489	1278	14477

^a Prime Full Window Extended (PFWE) and Prime Full Window (PFW) modes; ^b pn Prime Large Window (PLW) mode and any of the various MOS Prime Partial Window (PPW) modes; ^c other pn modes such as the Small Window, MOS modes e.g. Fast Uncompressed (FU). pn timing and burst modes and MOS Refresh Frame Store modes are not used in the catalogue.

the MOS CCD positions to improve the astrometry, see Sec. 4.1 for further details. Improvements were also applied to the XMM-OM data and will be included in the upcoming version of the XMM-OM catalogue, SUSS version 7.

110 The main data processing steps used to produce the 5XMM data products were similar to those outlined in Webb et al. (2020); Rosen et al. (2016); Watson et al. (2009) and described on the SOC webpages². For all the 5XMM data, the observation data files were processed to produce calibrated event lists. The optimised background time intervals were identified and using them, the filtered exposures (taking into account exposure time, instrument mode, etc.), multi-energy-band X-ray images, and exposure maps were generated. The initial detections were made on single observations, using simultaneously all images and bands (1: 0.2-0.5 keV, 2: 0.5-1.0 keV, 3: 1.0-2.0 keV, 4: 2.0-4.5 keV, 5: 4.5-12.0 keV) from the three cameras when available, as in Watson et al. (2009); Rosen et al. (2016). The probability, and corresponding likelihood, were computed from the null hypothesis that the measured counts in the search box result from a Poissonian fluctuation in the estimated background level. A detection mask was made for each camera that defines the area of the detector which is suitable for source detection. An initial source list was made using a ‘box detection’ algorithm. This slides a search box (20'' × 20'') across the image defined by the detection mask. Sources were cut-out using a radius that was dependent on source brightness in each band, and these areas of the image where sources had been detected were blanked out. The source-excised images, normalised by the exposure maps, and the corresponding masks are convolved with a Gaussian kernel to create the background map (see Traulsen et al. 2019). A second box-source-detection pass was then carried out, creating a new source list, this time using the background maps (‘map mode’) which increased the source detection sensitivity compared to the first pass. The box size was again set to 20'' × 20''.
130 A maximum likelihood fitting procedure was then applied to the sources to calculate source parameters in each input image, by fitting a model to the distribution of counts over a circular area of radius 60'', see Watson et al. (2009). 1.19 million detections were made before the stacking procedure. These detections were then used to produce detection level spectra and lightcurves, if the detection was comprised of more than 50 EPIC counts, where previously 100 EPIC counts were required. Detections were then matched to their counterpart stacked sources (see Sec. 4) after stacking. This resulted in almost 409000 spectra and lightcurves, an increase of 10% with respect to 4XMM-DR14. Automatic and visual screening procedures were carried out to check for any problems in the data products, as described in Sec. 10. The
150

² https://xmm-tools.cosmos.esa.int/external/xmm_products/pipeline/doc/19.16_20210326_1200/current_release

data from this processing have been made available through the *XMM-Newton* Science Archive³ (XSA), but see also Sec. 13.

4. Stacking and source detection

Source detection on *XMM-Newton* EPIC observations uses maximum-likelihood fits under Cash statistics as described for example by Watson et al. (2009); Traulsen et al. (2019). With 5XMM, we introduce a revised approach to stacked source detection in order to handle all 5XMM data from single observations to 99 directly overlapping observations. During the source-detection step, we assume that the flux of each source remains constant over all exposures and that its spectrum in the five standard energy bands can be described by a simple model. We choose an absorbed power-law as the spectral model, which is a reasonable approximation to most *XMM-Newton* sources (e.g. Watson et al. 2009; Mateos et al. 2009). Under these assumptions, the equations for the maximum-likelihood detection take the same form and the same degrees of freedom irrespective of the number of exposures in which a source is fitted. The degrees of freedom are the source coordinates, the mean source flux, and the spectral parameters – column density and power-law index – if the source is fitted as point-like, and additionally the radius of the extent model, if the source is fitted as extended. The results of the five-band spectral fit in the detection step are given in the catalogue in the columns with the prefix “STACK_”.
170

The photon flux is related to the measured count rates in each input image to source detection by energy conversion factors (ECFs). In the new *XMM-Newton* source detection, the ECFs for each fitted pair of spectral parameters, for each fitted detector position, and for each instrumental setup (EPIC/pn, MOS1, MOS2 with their respective filters) are extrapolated on the fly over a grid of pre-compiled values. They cover column densities between 10¹⁹ cm⁻² and 10²³ cm⁻² and power-law indices between 0 and 5. Time-dependence of the EPIC instrumental cross-calibration is taken into account over six different epochs.
180

Once a source is reliably detected with a log-likelihood STACK_DET_ML ≥ 6 (equivalent to a 3 σ detection), the assumptions of constant flux and power-law spectrum are dropped, and image-level count rates and related parameters are determined by forced PSF photometry at the detected source position and extent radius. During PSF photometry, the count rate in each contributing image is treated as a free fit parameter: the method used in source detection in the previous Serendipitous Source Catalogues from EPIC data. The photometry results are given in the catalogue in the RATES, FLUX, DET_ML⁴ and related columns without the prefix “STACK”.
190

³ <https://nxsas.esac.esa.int/nxsa-web>

⁴ see: http://xmmssc.irap.omp.eu/Catalogue/5XMM-DR15/5XMM-DR15_Catalogue_User_Guide.html#Catalogue

For each fit parameter, the lower and upper confidence limits are calculated, searching for the parameter values for which the minimum Cash statistics value plus one, $C_{\min} + 1$ (1σ for Gaussian distributed data), is reached. For an efficient and robust search, the source-detection task `emldetect` employs the so called false-position method, which is a numerical bracketing approach. This method converges robustly in most cases, but still fails occasionally (e.g. Chapra & Canale 2015), with *XMM-Newton* test data and the most recent `emldetect` implementation for example in about two out of a thousand error calculations on a coordinate and in less than one out of ten-thousand error calculations on another fit parameter. If the calculation of an error component does not converge, this component is now set to undefined in all cases. Previously, a count-rate dependent fall-back value was used for coordinates, extent, and count rates. The total 1σ error on a parameter is the arithmetic mean of the lower and the upper errors if both are defined. If one component does not converge, the other component is taken as the total error. In addition to the total errors, 5XMM also includes the asymmetric upper and lower errors on the image coordinates, the extent radius, and the spectral fit parameters `STACK_FLUX`, `STACK_NH`, and `STACK_GAMMA`.

The mean symmetric position error `RADEC_ERR` represents the radius of a circular confidence region, which can actually have an elliptical shape. For a simple approximation of elliptical position errors, three more pairs of confidence limits are determined by `emldetect`: The parameters of an ideal ellipse containing 68% of the positions could be calculated from any three points on this ellipse. We thus determine $C_{\min} + 2.3$ (1σ for two parameters) along the x-axis, the y-axis, and their diagonal, and derive the semi-axes and orientation of an ideal error ellipse from them. These ellipse parameters are provided in the catalogue. Exploring a larger parameter range, the false-position method fails more often than for $C_{\min} + 1$. We thus included Brent's method (e.g. Chapra & Canale 2015) as an alternative in the ellipse fits, which reduces the cases in which the $C_{\min} + 2.3$ fit does not converge.

Whilst noisy pixels are identified before source detection (often referred to as hot pixels), some pixels can become noisy for a period of time only (referred to as warm pixels). These can be identified by stacking all observations in the same orientation, but removing duplicate observations, with very similar pointing direction and orientation, to avoid sources becoming significant in the stack. The warm pixels can then be identified as they appear as bright in these images. The time dependence can then be identified. For 5XMM-DR15 this work is carried out following the source detection on stacked data and the warm pixels coinciding with detections are indicated in the 10th detector flag (`PN_FLAG`, `M1_FLAG` or `M2_FLAG`) as T (true). This is then propagated to the `SUM_FLAG` to indicate a possibly spurious detection/source.

4.1. Astrometry

4.1.1. MOS to pn comparison

In preparation for 5XMM, we carried out astrometry checks on 4XMM-DR10, similar to those described in Webb et al. (2020). We projected all EPIC sources and their counterparts to detector coordinates via `ecoordconv`, and averaged the offsets over a large number of neighbouring sources. We saw obvious patterns emerge, reminiscent of the structure of the EPIC MOS CCDs. We then rebuilt the source lists corresponding to 4XMM-DR12 independently for each camera, confirming obvious CCD pat-

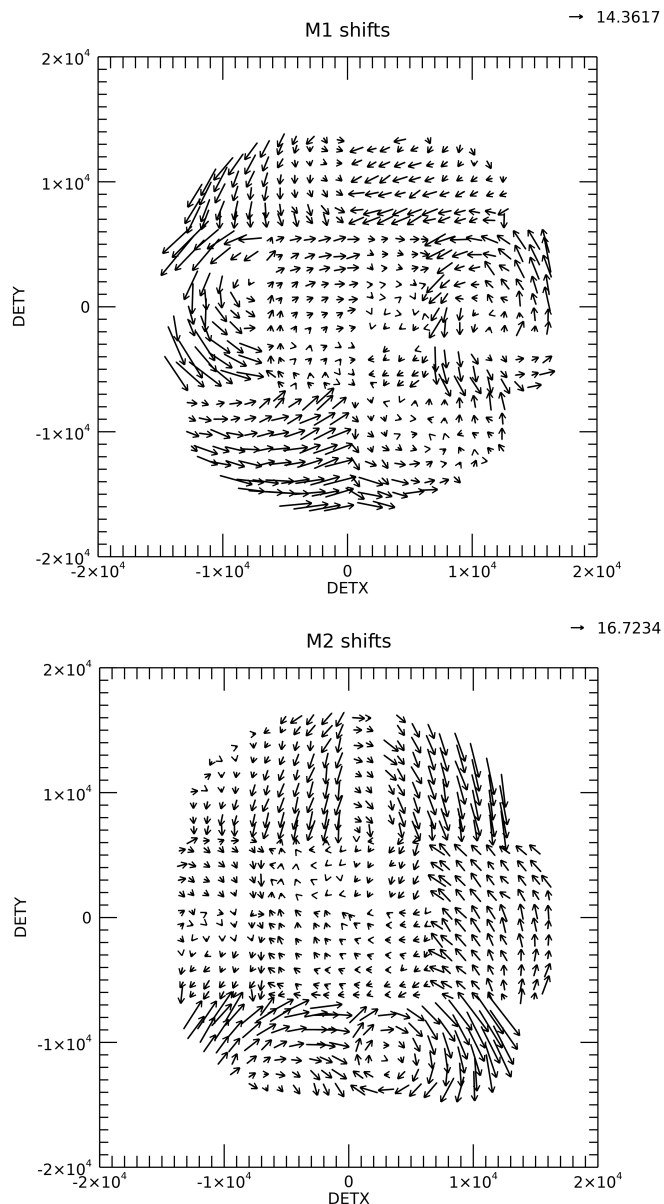


Fig. 2: Offsets between MOS1 (top) and MOS2 (bottom) source positions with respect to corresponding pn sources in detector coordinates of each instrument, for 4XMM DR12. Each point is an average over $1 \times 1'$ inside a given CCD. CCD 1 is at the center, CCD 2 in the lower right corner and CCD number increases counterclockwise up to CCD 7 on the lower left. The vector scale is given on the top on the right, in detector units of $0.05''$.

terns for both MOS1 and MOS2, see Fig. 2. This indicated that, to first order, repositioning the lateral MOS CCDs (by a suitable translation and rotation) should reduce the offsets considerably. Indeed the positioning of lateral CCDs was not reassessed since the initial release in October 2000⁵. Only the global alignment of the telescopes was considered since then, ensuring good overlap between MOS and pn source positions in the central MOS CCD but not in the lateral CCDs.

⁵ <https://xmmweb.esac.esa.int/docs/documents/CAL-SRN-0003-1-0.pdf>

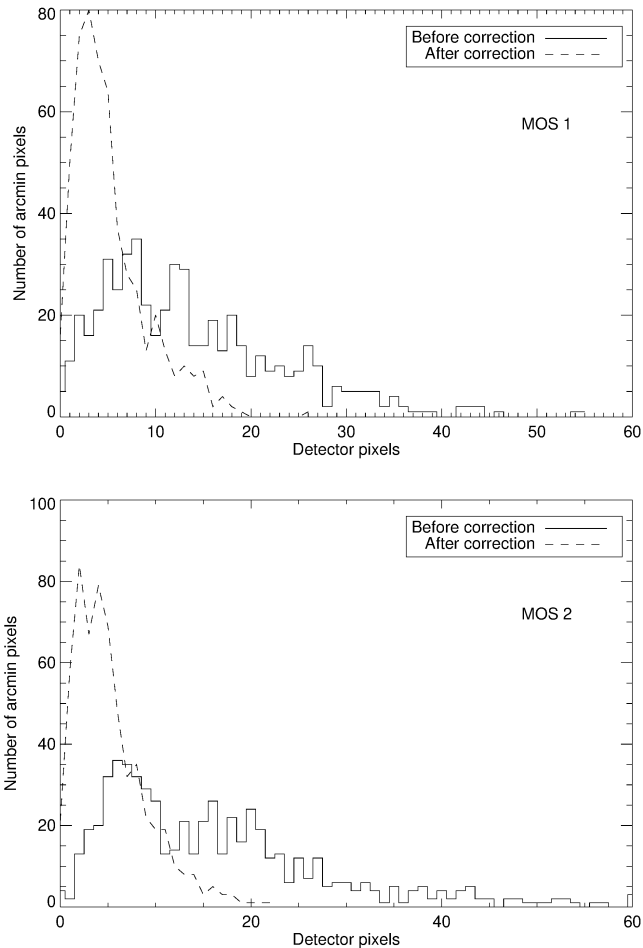


Fig. 3: Histograms of the offsets between MOS and pn sources averaged over arcmin pixels for MOS1 (top) and MOS2 (bottom). The solid histogram reflects directly Figure 2. The dashed line is what remains after optimally translating and rotating each MOS CCD. The scale is in detector units of $0.05''$.

Rather than attempting to realign the MOS CCDs based on an optical reference (which is limited to sources that have AGN counterparts, and is a little blurred by possible pointing imprecision of the satellite itself), we decided to align MOS CCDs to pn, whose detector is a single Si wafer. This is purely internal to XMM-Newton (does not require any outside reference) and should ensure that MOS and pn sources are aligned. It does not account, however, for (smaller) distortions between pn-only sources and optical quasars due to the EPIC pn telescope. We show the resulting offsets, averaged over $1'$ pixels, in Figure 2.

We then proceeded to translate and rotate each MOS CCD to correct for the offsets⁶. This reduced the median offsets from $0.6''$ (M1) and $0.7''$ (M2) to $0.24''$ for both cameras, with much smaller extreme offsets, nearly all below $1''$, as illustrated in Figure 3. This corrected MOS CCD metrology was applied to the entire XMM archive during the full reprocessing in November 2024, ensuring that the 5XMM sources do not suffer from this effect.

⁶ <https://xmmweb.esac.esa.int/docs/documents/CAL-SRN-0406-1-1.pdf>

4.1.2. Systematic position error

The systematic uncertainty of the 5XMM-DR15 astrometry is estimated using a statistical approach based on the cross-matching of the X-ray sources with an external catalogue with accurate positions. The adopted methodology is similar to that described in Section 6.2 of Merloni et al. (2024) and outlined in Appendix B of this paper. It is designed for X-ray positional errors that are symmetric in the direction of the right ascension and declination and which are described by a normal distribution. Under these assumptions, the probability of a radial offset r of an X-ray source from its true position is given by the Rayleigh distribution with parameter σ that represents the astrometric standard deviation in right ascension or declination. This σ is considered to have both a statistical (σ_{stat}) and a systematic (σ_{sys}) component

$$\sigma = \sqrt{\sigma_{\text{stat}}^2 + \sigma_{\text{sys}}^2}, \quad (1)$$

with $\sigma_{\text{stat}} = \text{RADEC_ERR} / \sqrt{2}$ from source detection. The systematic uncertainty is inferred at the population level from the distribution of the angular separation between the X-ray sources and an external catalogue with vanishing astrometric errors. For a clean sub-sample of point-like XMM-Newton sources and QSOs selected from Gaia and unWISE (see Sect. B), we infer a systematic position error $\sigma_{\text{sys}} = 0''.88 \pm 0''.01$. This σ_{sys} is larger than the one derived for 4XMM-DR9s ($0.227''$) by Traulsen et al. (2020), because the 5XMM data were not rectified astrometrically when producing DR15. The astrometric correction will be included in DR16 to further improve the source positions.

5. Long-term variability

5XMM-DR15 is a stacked catalogue, containing all of the sources detected following the stacking of overlapping observations, but also includes the individual detections in each of the contributing observations and non-detections when no detection was made. This provides the user with long-term XMM-Newton variability over the 25 years of data used to construct the catalogue. These data can be visualised in lightcurves produced for each source.

To increase the timeframe over which variability can be examined and to increase the number of data points for each source, the STONKS algorithm was implemented (Quintin et al. 2024). This algorithm uses a master catalogue constructed from data from a variety of different observatories. To create the 5XMM-DR15 catalogue, this master catalogue was generated in February 2026 using the most recent versions of the XMM-Newton catalogue (4XMM-DR14), Chandra Source Catalogue version 2.1 (Evans et al. 2010), the Living Swift/XRT point-source catalogue (Evans et al. 2023), the eROSITA eRASS1 catalogue (Merloni et al. 2024), the XMM-Newton Slew survey catalogue version 3, XMMSL3 (Saxton et al. 2008), and the ROSAT catalogues 2RXS (Boller et al. 2016) and WGA-CAT (White et al. 1994). We also generated upper limits for the XMM-Newton non-detections using the RapidXMM (Ruiz et al. 2022) version of HILIGT (Saxton et al. 2022). Matching was done on a two by two basis using an algorithm based on Budavári & Szalay (2008) and implemented in NWAY (Salvato et al. 2018), see Quintin et al. (2024) for the details of the algorithm. We ensured that the fluxes estimated were comparable by converting each flux detection to a single, common

energy band. The common band we chose was the 0.1–12 keV band, as it contains the energy bands of every one of the missions we used and then assumed an absorbed powerlaw spectrum, with parameters $\Gamma = 1.7$ and $n_H = 3 \times 10^{20} \text{ cm}^{-2}$, values commonly employed in the majority of the catalogues used.

The pessimistic variability ratio was calculated, taking the ratio of the highest flux point minus the 1σ error and the lowest flux point plus the 1σ error. Alternatively, in the case of an upper limit, the ratio is calculated using the 3σ upper limit and the highest flux point minus the 1σ error. We provide significant long-term variability for sources that have a ratio of five or greater. The variability is calculated for the detections made with the standard pipeline before stacking. This is provided in the column 'APPROX_SOURCE_VAR'. There are 41187 sources with a variability ratio of five or greater, with the highest reaching a ratio of 78000. The mean variability is a factor 80.

For close to real time alerts of long-term variable *XMM-Newton* sources, the reader is invited to consult the recently implemented transient alert server at <http://flix.irap.omp.eu/stonks>. STONKS is available ⁷ as a web service that takes a tar ball containing the EPIC source list and the matching EPIC image as input, and returns plots for each detected alert as output.

6. Spectral fitting

The absorbed power-law fit during source detection is restricted to five energy points – the five standard energy bands – and therefore limited. For bright enough sources, full energy-resolved spectra were thus extracted and fitted separately. The procedure followed to select and analyse the stacked source spectra is based on Viitanen et al. (2025), with some changes in the procedure used to combine the spectra, and in the output quality flags. We refer the reader to that paper for the details. Here we outline the procedure used for 5XMM.

The spectra and background for the individual detections, see Sec.3 were checked to ensure that they had a positive number of total counts (in the extracted detection spectrum, including the source and the background), background counts (in the extracted background spectrum), and net counts (calculated by subtracting the background counts from the total counts, after scaling the former for the extraction area), all in the full 0.2–12keV band. If any of these conditions is not fulfilled, the spectrum is discarded from further processing.

The selected spectra are then separated by instrument (pn or MOS). For each instrument and each physical stacked source, we retained the single detection with the highest signal-to-noise ratio (defined as the net counts N divided by $\sqrt{2T - N}$, where T is the total counts; equivalently $N / \sqrt{T + B_{sc}}$, with B_{sc} the area-rescaled background counts), simplifying the spectral-combination procedure of Viitanen et al. (2025). This ensured that, for each physical stacked source, there was at most one pn spectrum and one MOS spectrum, which were then re-binned to have one or more counts per bin.

For the spectral fitting and modelling procedures, we employed Xspec (Arnaud 1996) through its Python interface together with the Bayesian X-ray Analysis (BXA) tool (Buchner et al. 2014), which connects Xspec to the nested-sampling package UltraNest (Buchner et al. 2021). The spectral models were implemented in Xspec and explored with BXA. For the power-law model used in 5XMM, we adopted fixed, source-independent priors: a log-uniform (Jeffreys) prior on N_H over

[0.001, 1000] (in units of 10^{22} cm^{-2} , i.e. $[10^{19}, 10^{25}] \text{ cm}^{-2}$), a uniform prior on Γ over [1.0, 3.0], and a uniform prior on $\log_{10}(\text{Flux})$ over [−15, −9]. Following Viitanen et al. (2025), the log-uniform prior on N_H avoids the starting-value artefact produced by tightening priors around a preliminary fit. The posterior probability distributions were then obtained from the BXA/UltraNest sampling and stored as chains for further summary and export.

For 5XMM we fitted an absorbed powerlaw (in Xspec notation `cflux * phabs * powerlw`) to the selected spectra, as a flexible general-purpose model. This model has three free spectral parameters: the flux (SPEC_FLUX_PL in the catalogue, observed flux not corrected for absorption), the column density (SPEC_NH_PL, not constrained by the column density of our Galaxy in the direction of the source) and the spectral slope (SPEC_GAMMA_PL). In addition, when pn and MOS spectra are fitted jointly, we included an inter-instrument normalisation parameter (SPEC_IIN_PL), implemented as a multiplicative constant, with the pn normalisation fixed to unity and the MOS normalisation left free. Future versions of the catalogue will include additional model fits.

All spectral fitting was performed using the Cash statistic (C-stat; Cash 1979), which is appropriate for Poisson-distributed data and particularly effective in the low-count regime. As a pre-screening step, a quick Levenberg–Marquardt fit of the source model to the selected spectrum was performed over 0.3–10 keV with standard background subtraction. The resulting χ^2 test statistic was converted to a p -value via the χ^2 cumulative distribution; sources with $p < 0.01$ were discarded and excluded from the Bayesian fitting.

For sources passing the previous filter, the pn and MOS spectra (when both available) were jointly fitted with BXA, using standard background subtraction. The goodness-of-fit was estimated with a permutation test, following Buchner et al. (2014). In short, the standard method of comparing the C-stat between the fit to the real data and fits to simulated data using the best fit model parameters is incorrect, because the model and the real data are not statistically independent. We instead estimated the p -values using a permutation test. For each source, we generated 1000 resampled datasets by randomly redistributing the combined data+model counts into two equal-size subsamples, allowing each energy-bin count to originate from either the observed or modelled spectrum. For each resampling, we computed the corresponding Kolmogorov–Smirnov (KS) statistic. The p -value was then defined as the fraction of permutations yielding a KS statistic larger than that of the original data–model comparison. Models with KS p -values ≥ 0.01 were considered acceptable fits.

The posterior probability distribution as sampled by BXA is provided as chains of parameter combinations. In the catalogue, we consider each parameter independently and provide the median and the 16 and 84% percentiles (i.e. $\pm 1\sigma$). In addition, we also provide the degrees of freedom (SPEC_DOF_PL) and the KS GoF p -value (SPEC_PVALUE_PL) of the fit. The SPEC_INFO field gives the URL of the corresponding source page in the online catalogue. Finally, a flag (SPEC_FLAG_PL) is also provided, with the following possible values:

- 0: no issues detected
- 1: zero or negative background counts
- 2: zero or negative source (total) or net counts
- 3: could not create the spectrum for fitting, or the pre-screening fit failed (including $p < 0.01$)
- 4: BXA spectral fit failed (no posterior produced)

⁷ <https://xcatdb.unistra.fr/stonks>

- 5: poor goodness-of-fit, with KS p -value < 0.01
- 6: photon index pegged, with PhoIndex median within 0.05 of the hard prior limits (≤ 1.05 or ≥ 2.95 , for priors in the range [1.0, 3.0])
- 7: poor goodness-of-fit and photon index pegged
- 470 – 8: NH pegged, with median NH ≥ 100 in units of 10^{22} , cm^{-2} (i.e. Compton-thick); flagged as poorly constrained at high column density
- 9: poor goodness-of-fit and NH pegged
- 10: photon index pegged and NH pegged
- 11: poor goodness-of-fit, photon index pegged, and NH pegged

7. OM data

The XMM-OM observes the sky simultaneously with the X-ray instruments onboard *XMM-Newton* in a reduced field of view up to $17 \times 17'$. The XMM-OM utilises a filter wheel between the telescope and the detector to select either an imaging filter or a grism for dispersed spectroscopy (Mason et al. 2001). Principal investigators of *XMM-Newton* observations choose the filters and/or grisms which are most suitable for their scientific objectives, and therefore there is considerable diversity in the combination of filters and grisms that are employed in an *XMM-Newton* observation. For 5XMM, XMM-OM counterparts to X-ray sources are drawn from version 6.2 of the *XMM-Newton* Serendipitous Ultraviolet Source Survey (XMM-SUSS) catalogue. 87039 X-ray sources have an OM counterpart. The XMM-SUSS is compiled from images obtained through the six primary photometric filters of XMM-OM, which have effective wavelengths from 2120 Å (UVW2) to 5430 Å (V). The construction of XMM-SUSS 6.2 is broadly as described in Page et al. (2012), but whereas the first version of XMM-SUSS was restricted to sources that are detected in the ultraviolet, XMM-SUSS 6.2 includes sources detected in any of the six optical and ultraviolet photometric filters.

For XMM-OM counterparts, 5XMM contains the corresponding source ID in XMM-SUSS 6.2, the match-probability to the X-ray source and the following information for each and every XMM-OM passband in which the counterpart is detected: AB magnitude and magnitude uncertainty, a quality flag, an extended flag, a χ^2 value and the degrees of freedom for which it is calculated. The AB magnitude and uncertainty provided for each band is a weighted mean of the measurements over all *XMM-Newton* observations in which the source is detected, and the corresponding magnitude uncertainty. The quality flag is an integer equivalent to a binary number in which each bit corresponds to a different quality issue; a bit is set to 1 when a data quality concern is identified or otherwise set to 0. Sources with the highest quality in the corresponding photometric band will thus have a value of 0. For more details of the quality flagging applied in XMM-SUSS see Page et al. (2012). The extended flag is set to 0 if the counterpart is consistent with the shape of a point source in the corresponding band or 1 if the source appears extended. Note that it is possible for a source to appear point-like in one passband but have measurable extent in another. Where a source has been detected in multiple *XMM-Newton* observations in the same XMM-OM passband a χ^2 value is computed for the sequence of magnitude measurements compared to a single, constant magnitude. The corresponding degrees of freedom is one fewer than the number of measurements in that band. The χ^2 divided by the degrees of freedom can be used as an indicator as to whether there is evidence for variability between observations in that XMM-OM passband. Where variability is suggested by the χ^2 ,

the individual measurements can be consulted in XMM-SUSS 6.2⁸. Caution is advised in inferring variability from χ^2 when the corresponding quality flag is other than 0, or for sources which appear point-like in some *XMM-Newton* observations and extended in others, because the photometry is measured differently for extended and point-like sources (see Page et al. 2012).

XMM-OM counterparts have been classified probabilistically into Galactic and extragalactic source types and this classification information is included in 5XMM; see Sec. 9 for more details.

8. Redshifts

8.1. Photometric redshifts

Photometric redshift estimation for X-ray selected samples requires a robust identification of the nature of each source, and in particular a reliable separation of AGN from stars and normal galaxies. Misclassified sources can lead to the use of incorrect SED libraries, priors or training sets, resulting in biased or catastrophic photometric redshift estimates. A dedicated classification step, and in particular the identification of AGN-dominated systems, is therefore essential to select the correct photometric models, improve the accuracy and reliability of the photometric redshifts, and control contamination in AGN samples derived from X-ray surveys. We selected all 5XMM sources classified as AGN (see Sect. 9) and outside the Galactic plane (CLASSX_CLASS="AGN" or CLASSOPT_CLASS="QSO", and $|b| > 20$ deg), a total of 464 280 sources. We compiled the optical and NIR-MIR photometry provided in their SEDs (Sect. 12) and calculated photometric redshifts for these sources. We used two different algorithms for estimating redshifts: TPZ (Carrasco Kind & Brunner 2013), a machine learning algorithm, and LePhare (Arnouts et al. 1999; Ilbert et al. 2006), a template fitting algorithm. The 5XMM catalogue provide photo-z for 154 734 sources, ~ 33% of the 464 280 selected sources, where a good fit was found.

8.2. Machine learning algorithm

We calculated photometric redshifts for these catalogues using the MLZ-TPZ package in combination with the training samples described in Sect. 8.4. MLZ-TPZ is a machine learning algorithm based on a supervised technique with prediction trees and random forest. The photometric redshifts and the corresponding cross-validation of the results was done through the photo-z pipeline we developed for the project, which includes a k-fold cross-validation method to evaluate the accuracy and reliability of our method, the selection of the optimal feature set for photo-z calculations using a Recursive Elimination Feature algorithm, and the quality evaluation of the individual photo-z, by using the shape of the redshift probability distribution given by TPZ and LePhare (see Sect. D.2).

8.3. Template fitting algorithm

We also calculated photometric redshifts using LePhare, an SED template fitting code. By using a template-based algorithm we can find potential high-redshift sources that would be otherwise missed by TPZ, since the results of machine learning methods are limited to the redshift range of the corresponding training

⁸ <https://www.cosmos.esa.int/web/xmm-newton/xsa>

sample (redshifts below 3.5 in our case). Moreover, LePhare allows us to estimate redshifts for sources with only partial photometry in the optical or infrared bands.

We used two different sets of templates for LePhare, depending on the optical morphological classification of the sources. For extended objects we used the templates proposed by Salvato et al. (2009, 2011) for the COSMOS survey. For point-like objects we used the eFEDS templates (Salvato et al. 2022).

8.4. Data: Training and validation samples

590 We compiled a large spectroscopic sample of X-ray selected extragalactic sources that can be used for the training and cross-validation of the machine learning and template fitting algorithms we used for calculating photometric redshifts. The training sample was selected from the second version of the Millions of Optical-Radio/X-ray Associations (MORX) Catalogue (Flesch 2024). We selected sources with secure spectroscopic redshifts, with an X-ray counterpart and classified as AGN or galaxies. We defined four different subsamples based on the photometry available in these large area optical surveys: SDSS
600 sample (~ 55,000 sources), PanSTARRS sample (~ 47,000 sources), SkyMapper sample (~ 6000 sources), and DES sample (~ 14,000 sources). Ancillary photometry in the near- (from the 2MASS, UKIDSS and VHS surveys) and mid-infrared (All-WISE catalogue) was included if available.

This training sample is an order of magnitude larger than those previously used in similar efforts to estimate photometric redshifts for X-ray sources using machine learning techniques (e.g. Mountrichas et al. 2017; Ruiz et al. 2018), thus improving on previous work. The spectral redshifts are provided in the
610 5XMM catalogue under the column REDSHIFT_ZSP.

9. Classification

Both the X-ray sources and the optical / ultraviolet sources in XMM-SUSS 6.2 have undergone a classification using an adapted version of the Naive Bayes classifier CLAXBOI, presented in Tranin et al. (2022). For the X-ray sources this algorithm uses the *XMM-Newton* X-ray properties of each source such as the hardness ratios, spectral fits (with a power law, but also an APEC model, used only in the classification), along with the X-ray to r-band flux ratio, the X-ray to W1 infra-red band ratio when these complementary data are available, the maximum X-ray variability, the X-ray luminosity when the distance is known from Gaia or the Glade+ catalogue (Dályá et al. 2022) and the distance to the centre of the galaxy in case of extra-galactic sources associated with a galaxy. For the X-ray sources, the most-likely classifications are given in the column CLASSX_CLASS. These are AGN, star, Galactic X-ray binary, cataclysmic variable, background AGN, extra-galactic X-ray binary and extended sources. We also provide an outlier measure, for the case when none of the above source-types matches the
620 source. Seven further columns provide the probability attributed to each classification. This allows the user to make an informed decision about the reliability of the classification. For the X-ray sources, there are 556337 AGN, 119661 stars, 26100 Galactic X-ray binaries, 1276 cataclysmic variables, 49969 background AGN, 22732 extra-galactic X-ray binary and 42581 extended sources. The higher the outlier value (maximum 10), the more likely the source does not fit any of the designated categories, but see also Tranin et al. (2022). 49404 sources have an outlier value greater than five. Off these sources that may be extreme
630 types of each classification, or indeed, different objects, 3943

have the best classification as AGN and maybe for example tidal disruption events, merging massive black holes or some extreme type of AGN, for example, 17834 have the best classification as a star, 15554 as a Galactic X-ray binary, 3657 as an extra-galactic X-ray binary and 2394 as an extended source.

For the XMM-OM sources only three source classes were retained, quasars (QSO), galaxies and stars. The most probable classification is given in the 'CLASSOPT_CLASS' column. Again the probability attributed to the three classes for each source is provided in the three subsequent columns. A total of
650 201536 sources have an optical classification. There are 66306 QSO, 29971 galaxies and 105259 stars. Of the 20564 sources with both an X-ray classification of star and an XMM-OM classification, all of the sources are classified as stars, implying that the classification is reliable. The other classifications are more complicated to compare as an AGN does not necessarily have the same definition as a quasar, however, of the X-ray sources classified as an AGN and with an XMM-OM classification, two-thirds are classified as a quasar.

10. Screening

The detection quality of each catalogue source is determined automatically using the SSC-internal task `dpssflag` and visually by inspection of all images. It is given as boolean strings in the catalogue columns PN,M1,M2_FLAG on instrument level (10 possible entries true/false) and EP_FLAG for the whole observation / stack (11 possible entries), and as an integer summary in the catalogue column SUM_FLAG. EP_FLAG on an observation level is defined as the worst of all instrument flags, and the flags on a stacked level take the detection with the most flags. For a clean detection, all strings are "F", and the summary flag is "0". The first ten components of the string flag encode low PSF coverage, a detection close to a bright point-like or close to an extended source, a detection close to a bad CCD area, each of which is summarised by a SUM_FLAG=1; a probably spurious detection which is close to a bright detection or which is significant in only one band or which is on a warm pixel each of which is summarised by a SUM_FLAG=2.

The eleventh character of the EP_FLAG indicates the result of visual screening, during which the screeners mark problematic detector areas like single reflection patterns and largely extended background emission. It translates into a SUM_FLAG=3 and to SUM_FLAG=4 if it is associated with one of the spurious or warm pixel flag (see Table 3). The flagging results for 5XMM-DR15 are summarised in Table 2.
680

11. Catalogue properties

The 5XMM-DR15 catalogue contains 818656 sources detected a total of 2578752 times, extracted from 14616 public *XMM-Newton* observations. Figure 4 shows the distribution of the source fluxes in the total EPIC band and in the soft (0.2-2.0 keV) and the hard band (2.0-12.0 keV). Also shown in the figure is the distribution of the EPIC counts.
690

As for previous versions of the stacked catalogue, data for detections and upper-limits are extracted and provided in 5XMM for each observation contributing to an individual source. Some sources are observed up to 98 times amounting to 2.7 Ms of observations. The distribution of the number of detections and upper limits per source is shown in Fig. 5.

New for 5XMM are columns providing information on *XMM-Newton* Optical Monitor (XMM-OM, Mason et al. 2001)

Table 2: Meaning of the characters in the quality flags PN_FLAG, M1_FLAG, M2_FLAG, EP_FLAG and their distribution in 5XMM-DR15. A source can have multiple flags.

Flag	Description	EP		PN		M1		M2	
all		818 656	100%	753 847	100%	692 614	100%	777 900	100%
0	No warning issued	537 679	66%	550 693	73%	594 018	86%	679 272	87%
1	PSF coverage < 50%	160 445	20%	97 178	13%	42 025	6%	41 749	5%
2	Near a bright point-like source	4 090	0%	3 875	1%	3 588	1%	3 876	0%
3	Near a bright extended source	60 955	7%	56 180	7%	54 894	8%	58 519	8%
4	Extended near a bright point source	728	0%	683	0%	629	0%	675	0%
5	Extended near a bright extended source	12 201	1%	8 203	1%	9 917	1%	11 088	1%
6	Extended, significant in one band	6 146	1%	5 467	1%	5 644	1%	5 862	1%
7	Extended, flag 4, 5, or 6	14 605	2%	11 699	2%	12 654	2%	13 694	2%
8	On a bad pixel or CCD area	24 433	3%	24 342	3%	170	0%	34	0%
9	Near a bad CCD area	65 913	8%	65 711	9%	476	0%	268	0%
10	On a warm CCD pixel	13 662	2%	8 589	1%	4 482	1%	492	1%
11	Flagged during visual screening	54 516	7%						

Table 3: Meaning of the SUM_FLAG values

Value	Description
0	Good
1	if the warning flags EP_FLAG 1, 2, 3 or 9 set to true but not 7, 8, 10 or 11
2	if the possibly-spurious or warm pixel flags EP_FLAG 7, 8 or 10 set to true but not the manual flag 11
3	if the manual flag EP_FLAG 11 is set to true but not the spurious or warm pixel flags 7, 8 or 10
4	if the manual flag 11 as well as one of the spurious or warm pixel flags 7, 8 or 10 are set to true

700 counterparts to the X-ray sources when identified, along with a measure of the long-term variability when observed multiple times, see Sec. 7 for column details and source classification is provided for both the OM and X-ray source, see Sec. 9, the WISE and Gaia counterparts where they exist, along with the Gaia distance estimate and link to further tables with other multi-wavelength data, see Sec. 12, the spectral fitting parameters for absorbed power laws, see Sec. 6, the long-term X-ray variability, see Sec. 5 and the photometric redshift for extragalactic sources, see Sec. 8.

710 *11.1. Extended sources*

5XMM-DR15 contains 42,669 X-ray sources that are identified as extended objects, that is, with a core radius parameter, r_c , as defined in section 4.4.4 of Watson et al. (2009), $> 6''$ and $EXTENT_ML \geq 4$. This is 88 more than the number of sources considered as extended using the classification algorithm. 41,850 of the extended sources have $r_c < 80''$, where $80''$ is the highest extension considered in the extended source fit, indicating that 819 sources may have an extension beyond $80''$.

720 In Section 9.5 of the 3XMM paper (Rosen et al. 2016) we compared the results of the detection of extended sources (in terms of extension and count rate) with the independent X-CLASS catalogue of galaxy clusters. This end-to-end comparison between two different data analyses on nearly the same data allowed us to quantify systematic uncertainties. Since then, our background modelling has changed (Section 3.4 of the 4XMM paper, Webb et al. 2020) so we re-evaluate the uncertainties here.

730 The latest release of X-CLASS (Koulouridis et al. 2021) contains 1,559 targeted and serendipitous galaxy clusters, about 3.7 times more than those available at the time of the 3XMM comparison. It is based on observations up to August 2015 and exposures limited to 20 ks, and is restricted to galaxy clusters, so it is not expected to contain all extended sources in 5XMM but

we would expect all X-CLASS sources to have counterparts in 5XMM. A specific difficulty is that the current X-CLASS catalogue does not include uncertainties on source positions, extensions, and count rates. In what follows we applied the same uncertainties as those in 5XMM (probably an underestimate, since 5XMM data are deeper).

Among the 1,559 X-CLASS clusters, only 44 are not matched to any 5XMM source (extended or not) within the search radius, defined as the quadratic sum of their extensions in both catalogues and their 95% positional errors. Those few missing X-CLASS sources are fainter and smaller than the others in X-CLASS. Among those that do have a 5XMM counterpart, 191 nearest counterparts are classified as point sources in 5XMM. There is a 5XMM extended source within the search radius in more than half of those cases, but we did not consider those in the comparisons because of the risk of confusion. For the correlations we also discarded X-CLASS sources with extensions larger than $80''$, since this is the maximum allowed in 5XMM, and 5XMM sources with errors on extension larger than $50''$ (unconstraining).

After applying all selections, there remained 1,199 common clusters. We show the results of the log-log correlations in Figures 6 (extension) and 7 (count rate). We removed outliers from the final correlations, with little impact. The 5XMM parameters tend to be smaller than the X-CLASS ones. The intrinsic (beyond statistical) scatter is strong (32% for extension and 39% for count rate), but the effect is clearly significant with so many points. The two parameters (extension and count rate) are strongly linked, with a close to linear relation between the 5XMM to X-CLASS ratios, as reported in the 3XMM paper.

From visual inspection of the extension plot, it seems that most points are above the optimal correlation at low flux and most points are below it at high flux, pointing to a log-log slope less than one (so that the 5XMM extension increases more slowly than the X-CLASS extension). This could be biased by the lack of uncertainties in X-CLASS, so we have not tried to

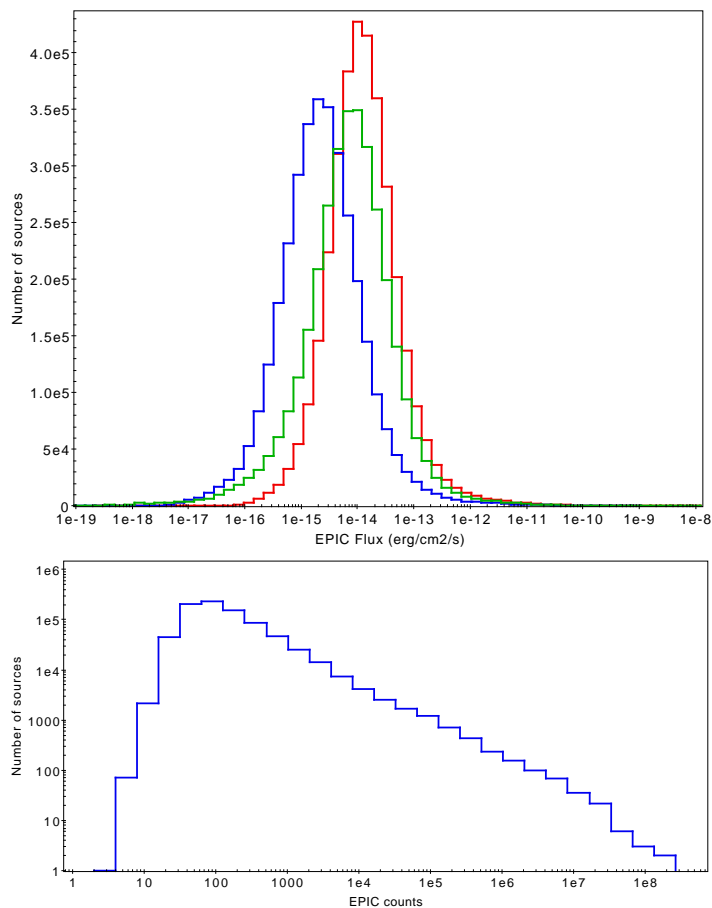


Fig. 4: Top: Distribution of source fluxes for the 5XMM-DR15 catalogue in the soft (0.2–2.0 keV, blue), hard (2.0–12.0 keV, green), and total band (red) energy bands. Bottom: distribution of total EPIC counts for the same sample of 5XMM-DR15 sources.

quantify this. The most likely reason for that difference is the background modelling. The 5XMM approach, which must be robust in all circumstances (in particular close to the Galactic plane), involves a smoothing radius of 40'' that could prevent detecting the outskirts of extended sources, resulting in smaller extension and count rate, particularly for large sources.

12. Catalogue products

12.1. ACDS: External catalogue cross-correlation

Cross-correlation with archival catalogues is performed by a distinct pipeline module running at the Observatoire Astronomique de Strasbourg and referred to as the Astronomical Catalogue Data Subsystem (ACDS). For each individual EPIC detection the ACDS lists all possible multi-wavelength identifications located within a 3σ combined XMM and catalogue error radius from the EPIC position. Finding charts and overlays with ROSAT all-sky survey images of the field are also produced. A detailed description of ACDS is given in Rosen et al. (2016). The list of the 227 archival catalogues used, which is updated regularly during operations, has been frozen for the bulk reprocessing, ensuring that all products associated with the 5XMM catalogue are built using the same configuration. Among the catalogues providing the largest sky coverage are GALEX GR6+7 (Bianchi et al. 2017), UCAC4 (Zacharias et al. 2013), SDSS DR12 (Alam et al.

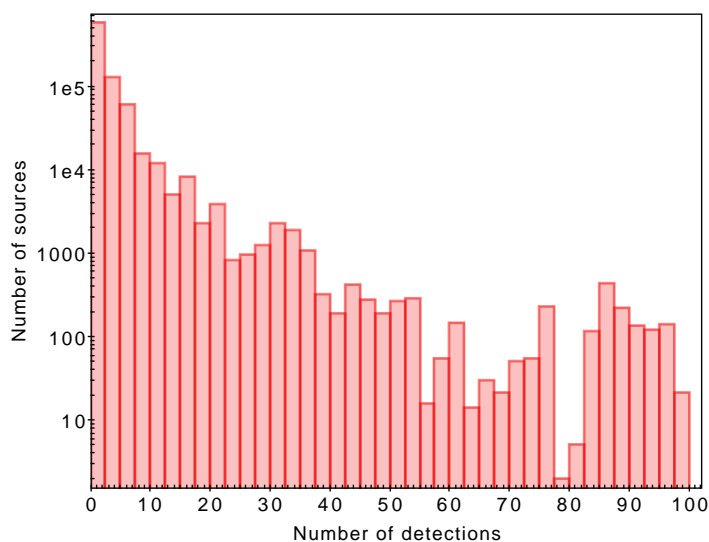


Fig. 5: 5XMM-DR15 unique sources plotted as a function of the number of detections/upper limits

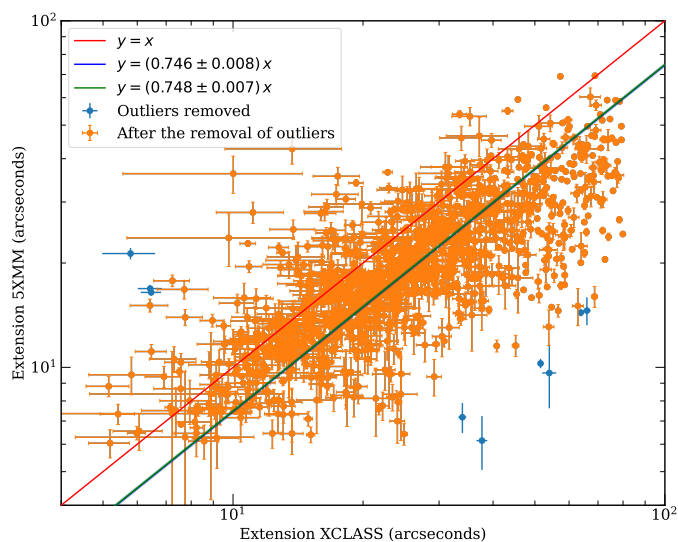


Fig. 6: Comparison of the first (blue) and final (green) regression fits before (blue) and after (orange) the outlier removal procedure against the ideal relationship ($y = x$ in red), indicating that on average the X-CLASS extensions are larger than the 5XMM extensions.

2015), panStarrs-DR1 (Chambers et al. 2016), IPHAS DR2 (Barentsen et al. 2014), Gaia DR3 (Gaia Collaboration et al. 2021), 2MASS (Cutri et al. 2003), AllWISE (Cutri et al. 2021), Akari (Ishihara et al. 2010), NVSS (Condon et al. 1998), FIRST (Helfand et al. 2015) and GLEAM (Hurley-Walker et al. 2017). The XMM-OM Serendipitous Source Survey Catalogue XMM-SUSS4.1 (Page et al. 2012), Chandra Source Catalogue v. 2.0 (Evans et al. 2024) and the second ROSAT all-sky survey (Voges et al. 1999) are also queried. Apart from the Chandra Catalogue Release 2.0 whose entries are extracted from the CXC server and Simbad which is served by a specific facility, all other ACDS catalogues are queried using the VizieR catalogue server. As for previous releases, the tentative identifications of 5XMM sources are not included in the catalogue itself, but are distributed to the community by the XSA and through the

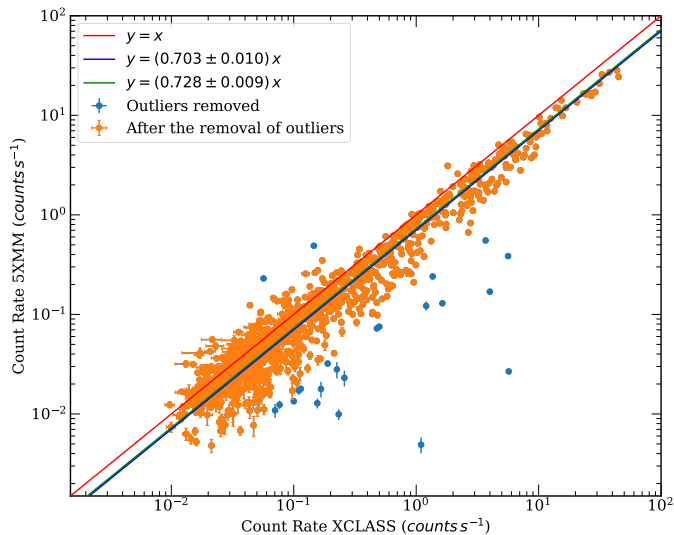


Fig. 7: Same as Figure 6 for count rate values.

XCatDB (Michel et al. 2015) as regular pipeline products. Finding charts are extracted from several imaging surveys with the following decreasing priority order. First the Sloan digital sky survey (Alam et al. 2015) with colour images made from the g, r and i is queried at the SDSS server. Second colour images are extracted from the z, g and z+g surveys of Pan-STARRS-DR1 (Chambers et al. 2016) third, images from the MAMA/SRC-J and MAMA/POSS-E plate collections and as a last choice from the DSS2 photographic plates are extracted. For the one-colour photographic surveys, we select the blue image at Galactic latitude $> 20^\circ$, while the red images are preferred in the Galactic plane. Apart from the SDSS, all images are extracted in HEALPix format from Hierarchical Progressive Surveys (HiPS) Aladin server (Fernique et al. 2014). Because the ACDS is operated by the regular pipeline, it does not process sources detected solely by stacking observations.

12.2. Methodology for producing multi-wavelength Spectral Energy Distributions

Spectral energy distributions (SEDs) are provided for each of the unresolved 5XMM-DR15 sources. For that purpose, we use ARCHES, a tool developed under The Astronomical Resource Cross-matching for High Energy Studies (ARCHES) framework under the EU Horizon ARCHES project (Motch et al. 2017). The ARCHES algorithm (Pineau et al. 2017) performs a simultaneous cross-match of all selected archival catalogues, computing a probability for every possible combination of catalogue entries. These probabilities are derived from the likelihood that sources from different catalogues share the same sky position, taking into account their astrometric uncertainties. The method is based exclusively on positional information and associated astrometric uncertainties, without the inclusion of photometric or other auxiliary source properties in the matching likelihood.

For each combination of catalogue entries, excluding high optical density regions such as the LMC, SMC, and M31, a cross-match probability is computed based on the likelihood that all sources share the same true sky position. This process led to the removal of 19 stacks comprising 708 observations. The final association probability further depends on a prior, representing the likelihood that an X-ray source has a genuine counterpart

in a longer-wavelength catalogue. This prior is estimated empirically from observed association rates while accounting for spurious matches. The method assumes ideal conditions, such as accurate astrometry, negligible systematics, and uniform source densities, while ignoring effects such as proper motion or blending. To reduce these limitations, additional preprocessing steps, such as clustering by source density, are performed. The stacks were grouped using a cumulative binning strategy. They were first sorted by the N_{optical}/N_x ratio, then grouped according to similar ratios while maintaining a threshold of 10^4 X-ray sources per group for the probabilistic cross-match. This resulted in 10 groups in total, with the last group corresponding to the highest N_{optical}/N_x ratio, requiring further rebinning due to the hardware limitation of the crossmatch node.

As for 4XMM (Webb et al. 2020) we select archival catalogues which cover the largest sky area and energy bands from UV to radio. The selected catalogues are listed in Table C.1. We grouped catalogues by wavelength coverage to produce single master UV, optical and infrared catalogues and perform a statistical cross-match of four catalogues (X-ray to IR). Some of the catalogues are all sky (e.g. Gaia DR3, 2MASS), while other cover only partially the areas covered by the 5XMM-DR15 stacks, notably GALEX in the UV. Probabilities depend on the source density of each of these catalogues, which in turn depends on the depth and the area covered by each catalogue. We therefore calculated the intersected area covered by each of the several groups of stacks with each of the catalogues using mocpy and set the area to the corresponding value.

The set of SEDs are available as individual FITS files and graphical output for the three highest probability SEDs via the SEDFinder service⁹.

13. Catalogue access

The catalogue of sources is provided in several formats. A Flexible Image Transport System (FITS) file is provided containing all of the sources, detections and upper limits in the catalogue. For 5XMM-DR15 there are 3 397 248 rows (relating to 818 656 individual sources) and 421 columns. This can be found on the XMM-Newton Survey Science Centre webpages¹⁰. A FITS table of the 818656 sources only is also provided. Ancillary tables to the catalogue, also available from the XMM-Newton Survey Science Centre webpages, include the table of observations incorporated in the catalogue.

The XMM-Newton Survey Science Centre webpages provide access to the 5XMM catalogue, as well as links to the different servers distributing the full range of catalogue products. These include XMM-SSC catalogue server which provides catalogue values and products for each source and detection in the catalogue¹¹, the ESA XMM-Newton archive (XSA), which provides access to all of the 5XMM data products, and the ODF data, the XCatDB¹² produced and maintained by the XMM-Newton SSC, which contains possible EPIC source identification produced by the pipeline by querying 227 archival catalogues, see Sec. 12. Finding charts are also provided for these possible identifications. Other source properties as well as images, time series and spectra are also provided. Multi-wavelength data taken as a part

⁹ <https://xcatdb.unistra.fr/sedfinder/> and through a link in the catalogue in the INFO_COUNTERPARTS column.

¹⁰ <http://xmmssc.irap.omp.eu/> and <https://xmmssc.aip.de>

¹¹ [://xmm-catalog.irap.omp.eu/](https://xmm-catalog.irap.omp.eu/)

¹² <https://xcatdb.unistra.fr/5xmm/>

of the XID (X-ray identification project) run by the SSC over the first fifteen years of the mission, e.g. Carrera et al. (2007).

The XCatDB operated by the Strasbourg Observatory provides both a classical HTML interface and an all-sky browser based on VO protocols. The latter is able to plot sources detected by XMM upon any of the 1100 image surveys available across the entire electromagnetic spectrum, and combine them with data from over 22,000 Vizier catalogues (Fernique et al. 2019). A key XCatDB feature, Amora (Asynchronous Multi-Observation Region-based Analysis), allows users to access and process all the events/photons collected by the *XMM-Newton* mission in the region of interest interactively. The catalogue can also be accessed through HEASARC¹³ and VIZIER¹⁴. The results of the external catalogue cross-correlation carried out for the 5XMM catalogue (section 12) are available as data products within the XSA or through the XCatDB. The *XMM-Newton* Survey Science Centre webpages also detail how to provide feedback on the catalogue.

Where the 5XMM catalogue is used for research and publications, please acknowledge the use by citing this paper and including the following:

'This research has made use of data obtained from the 5XMM serendipitous source catalogue compiled by the *XMM-Newton* Survey Science Centre, the XMM2ATHENA project and in collaboration with the *XMM-Newton* SOC.'

The 5XMM catalogue assumes that a source is significantly detected if it has a maximum log-likelihood value of six (STACK_DET_ML, 3σ detection) in the detection step that is based on the assumptions of constant source flux and a power-law spectrum. A good spectrum-based fit reaches a higher sensitivity and a higher detection likelihood than the fit to the image-level count rates during the forced PSF photometry. Thus, there are sources with EPIC log-likelihoods EP_DET_ML from the photometry step below STACK_DET_ML from the detection step. Further, some sources have been flagged as possibly spurious, see Sec. 10. In order to create the cleanest catalogue possible, where statistically almost all sources are real, it is necessary to filter the catalogue to include only EPIC sources with for example, a 4σ significance (Maximum likelihood of ~ 10) and to keep only those without flags or sources where some of the parameters may be affected by a neighbouring source or bad pixel, for example,

$$\text{STACK_DET_ML} > 10 \ \&\& \ \text{SUM_FLAG} < 2$$

Filtering with these criteria for 5XMM-DR15 leaves 607392 sources (74% of sources). 92.6% or 758497 of the sources have negligible pileup (XX_PILEUP < 1, where XX is either pn, M1 or M2 for the pn, MOS 1 or MOS 2 detectors).

14. Limitations of the catalogue

As indicated in Sec. 4.1.2, the astrometric correction was not propagated to the detections to stack and as a result, the systematic position error is greater (by a few tenths of an arcsecond) than expected if the correction had been properly propagated. This will be rectified from version DR16 of 5XMM. Equally, the warm pixels will be implemented before source detection from version DR16, as described in Sec. 4. This is expected to reduce the number of sources flagged by a few percent. Finally,

¹³ <http://heasarc.gsfc.nasa.gov/db-perl/W3Browse/w3table.pl?tablehead=name%3Dxmmssc&Action=More+Options>

¹⁴ <http://vizier.u-strasbg.fr/cgi-bin/VizieR>

additional catalogues will be added to those used in the photometric redshift determination, which could potentially increase the number of redshifts by a factor two thanks to the wider area coverage and increased infra-red data.

15. Summary

This paper describes the improvements made to the software and calibration used to produce the new major version of the *XMM-Newton* catalogue, 5XMM, a single catalogue in which all sources are processed through a unified stacking framework, including both overlapping and non-overlapping observations. The 5XMM-DR15 catalogue contains 818656 unique X-ray sources (0.2-12.0 keV) which were compiled from 2578752 individual detections or upper limits, with some sources observed as many as 98 times. The catalogue covers $\sim 3.5\%$ of the sky. In terms of the number of X-ray sources, 5XMM-DR15 is 88% of the eROSITA DR1 catalogue that covers half of the sky and more than twice the number of sources and detections that are in the Chandra source catalogue version 2.1.

In this new catalogue, many value-added products have been added including *XMM-Newton* Optical Monitor counterparts to the X-ray sources, a measure of the long-term X-ray and optical variability when sources have been observed multiple times, source classification for both the OM and X-ray source, WISE and Gaia counterparts where they exist, along with the Gaia distance estimate and link to further tables with other multi-wavelength data, spectral fitting parameters for absorbed power laws and photometric redshift for extragalactic sources. These quantities will facilitate the selection of homogeneous populations of sources, carry out spectral studies, find sources of similar luminosities, identify sources that vary on the short- and long-term and carry out multi-wavelength studies without the need for cross-correlating sources with multi-wavelength catalogues.

Acknowledgements. We are grateful for the strong support provided by the *XMM-Newton* SOC. We also thank the CDS team for their active contribution and support. The French teams are grateful to Centre National d'Études Spatiales (CNES) for their outstanding support for the SSC activities.¹⁵ This work was carried out in the framework of the project XMM2ATHENA, which has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement n°101004168. Support of the XMM-SSC work at AIP by Deutsches Zentrum für Luft- und Raumfahrt (DLR) through grants 50 OX 2301 and 50 OX 2601 is gratefully acknowledged.

References

- Alam, S., Albareti, F. D., Allende Prieto, C., et al. 2015, *ApJS*, 219, 12
- Arnaud, K. A. 1996, in *Astronomical Society of the Pacific Conference Series*, Vol. 101, *Astronomical Data Analysis Software and Systems V*, ed. G. H. Jacoby & J. Barnes, 17
- Arnouts, S., Cristiani, S., Moscardini, L., et al. 1999, *MNRAS*, 310, 540
- Barentsen, G., Farnhill, H. J., Drew, J. E., et al. 2014, *VizieR Online Data Catalog*, II/321
- Bianchi, L., Shiao, B., & Thilker, D. 2017, *ApJS*, 230, 24
- Boller, T., Freyberg, M. J., Trümper, J., et al. 2016, *A&A*, 588, A103
- Buchner, J. 2016, *Statistics and Computing*, 26, 383
- Buchner, J. 2019, *PASP*, 131, 108005
- Buchner, J. 2021, *The Journal of Open Source Software*, 6, 3001
- Buchner, J., Brightman, M., Baloković, M., et al. 2021, *A&A*, 651, A58
- Buchner, J., Georgakakis, A., Nandra, K., et al. 2014, *A&A*, 564, A125
- Budavári, T. & Szalay, A. S. 2008, *ApJ*, 679, 301
- Carrasco Kind, M. & Brunner, R. J. 2013, *MNRAS*, 432, 1483
- Carrera, F. J., Ebrero, J., Mateos, S., et al. 2007, *A&A*, 469, 27
- Cash, W. 1979, *ApJ*, 228, 939
- Chambers, K. C., Magnier, E. A., Metcalfe, N., et al. 2016, *arXiv e-prints*, arXiv:1612.05560

¹⁵ <https://ror.org/04h1h0y33>

- Chapra, S. C. & Canale, R. 2015, *Numerical Methods for Engineers*, 7th edn. (USA: McGraw-Hill, Inc.)
- 1020 Condon, J. J., Cotton, W. D., Greisen, E. W., et al. 1998, *AJ*, 115, 1693
- Cutri, R. M., Skrutskie, M. F., van Dyk, S., et al. 2003, *VizieR Online Data Catalog: 2MASS All-Sky Catalog of Point Sources (Cutri+ 2003)*, *VizieR On-line Data Catalog: II/246*. Originally published in: University of Massachusetts and Infrared Processing and Analysis Center, (IPAC/California Institute of Technology) (2003)
- Cutri, R. M., Wright, E. L., Conrow, T., et al. 2021, *VizieR Online Data Catalog: AllWISE Data Release (Cutri+ 2013)*, *VizieR On-line Data Catalog: II/328*. Originally published in: IPAC/Caltech (2013)
- 1030 Dállya, G., Díaz, R., Bouchet, F. R., et al. 2022, *MNRAS*, 514, 1403
- Ebrero, J. 2019, *XMM-Newton Users Handbook*, Tech. Rep. 2.17, ESA: XMM-Newton SOC
- Evans, I. N., Evans, J. D., Martínez-Galarza, J. R., et al. 2024, *VizieR Online Data Catalog: Chandra Source Catalog Release 2 (CSC 2.1) (Evans+, 2024)*, *VizieR On-line Data Catalog: IX/70*. Originally published in: 2024ApJS..274...22E
- Evans, I. N., Primini, F. A., Glotfelty, K. J., et al. 2010, *ApJS*, 189, 37
- Evans, P. A., Page, K. L., Beardmore, A. P., et al. 2023, *MNRAS*, 518, 174
- 1040 Fernique, P., Boch, T., Oberto, A., et al. 2019, in *Astronomical Society of the Pacific Conference Series*, Vol. 521, *Astronomical Data Analysis Software and Systems XXVI*, ed. M. Molinaro, K. Shorridge, & F. Pasian, 46
- Fernique, P., Boch, T., Pineau, F., & Oberto, A. 2014, in *Astronomical Society of the Pacific Conference Series*, Vol. 485, *Astronomical Data Analysis Software and Systems XXIII*, ed. N. Manset & P. Forshay, 281
- Flesch, E. W. 2024, *The Open Journal of Astrophysics*, 7, 6
- Gaia Collaboration, Brown, A. G. A., Vallenari, A., et al. 2021, *A&A*, 649, A1
- Helfand, D. J., White, R. L., & Becker, R. H. 2015, *ApJ*, 801, 26
- Hurley-Walker, N., Callingham, J. R., Hancock, P. J., et al. 2017, *MNRAS*, 464, 1146
- 1050 Ilbert, O., Arnouts, S., McCracken, H. J., et al. 2006, *A&A*, 457, 841
- Ishihara, D., Onaka, T., Kataza, H., et al. 2010, *A&A*, 514, A1
- Jansen, F., Lumb, D., Altieri, B., et al. 2001, *A&A*, 365, L1
- Koulouridis, E., Clerc, N., Sadibekova, T., et al. 2021, *Astronomy and Astrophysics*, 652, A12
- Mason, K. O., Breeveld, A., Much, R., et al. 2001, *A&A*, 365, L36
- Mateos, S., Saxton, R. D., Read, A. M., & Sembay, S. 2009, *A&A*, 496, 879
- Merloni, A., Lamer, G., Liu, T., et al. 2024, *A&A*, 682, A34
- 1060 Michel, L., Grisé, F., Motch, C., & Gomez-Moran, A. N. 2015, in *Astronomical Society of the Pacific Conference Series*, Vol. 495, *Astronomical Data Analysis Software and Systems XXIV (ADASS XXIV)*, ed. A. R. Taylor & E. Rosolowsky, 173
- Motch, C., Carrera, F., Genova, F., et al. 2017, in *Astronomical Society of the Pacific Conference Series*, Vol. 512, *Astronomical Data Analysis Software and Systems XXV*, ed. N. P. F. Lorente, K. Shorridge, & R. Wayth, 165
- Mountrichas, G., Corral, A., Masoura, V. A., et al. 2017, *A&A*, 608, A39
- Page, M. J., Brindle, C., Talavera, A., et al. 2012, *MNRAS*, 426, 903
- Pineau, F. X., Derriere, S., Motch, C., et al. 2017, *A&A*, 597, A89
- Quintin, E., Webb, N. A., Georgantopoulos, I., et al. 2024, *A&A*, 687, A250
- Rosen, S. R., Webb, N. A., Watson, M. G., et al. 2016, *A&A*, 590, A1
- 1070 Ruiz, A., Corral, A., Mountrichas, G., & Georgantopoulos, I. 2018, *A&A*, 618, A52
- Ruiz, A., Georgakakis, A., Gerakakis, S., et al. 2022, *MNRAS*, 511, 4265
- Salvato, M., Buchner, J., Budavári, T., et al. 2018, *MNRAS*, 473, 4937
- Salvato, M., Hasinger, G., Ilbert, O., et al. 2009, *ApJ*, 690, 1250
- Salvato, M., Ilbert, O., Hasinger, G., et al. 2011, *ApJ*, 742, 61
- Salvato, M., Wolf, J., Dwelly, T., et al. 2022, *A&A*, 661, A3
- Saxton, R. D., König, O., Descalzo, M., et al. 2022, *Astronomy and Computing*, 38, 100531
- Saxton, R. D., Read, A. M., Esquej, P., et al. 2008, *A&A*, 480, 611
- 1080 Shu, Y., Koposov, S. E., Evans, N. W., et al. 2019, *MNRAS*, 489, 4741
- Tranin, H., Godet, O., Webb, N., & Primorac, D. 2022, *A&A*, 657, A138
- Traulsen, I., Schwöpe, A. D., Lamer, G., et al. 2019, *A&A*, 624, A77
- Traulsen, I., Schwöpe, A. D., Lamer, G., et al. 2020, *A&A*, 641, A137
- Viitanen, A., Mountrichas, G., Stiele, H., et al. 2025, *A&A*, 704, A16
- Voges, W., Aschenbach, B., Boller, T., et al. 1999, *A&A*, 349, 389
- Watson, M. G., Schröder, A. C., Fyfe, D., et al. 2009, *A&A*, 493, 339
- Webb, N. A., Coriat, M., Traulsen, I., et al. 2020, *A&A*, 641, A136
- White, N. E., Giommi, P., & Angelini, L. 1994, in *American Astronomical Society Meeting Abstracts*, Vol. 185, *American Astronomical Society Meeting Abstracts*, 41.11
- 1090 Zacharias, N., Finch, C. T., Girard, T. M., et al. 2013, *AJ*, 145, 44

Appendix A: Data modes of *XMM-Newton* exposures

Table A.1: Data modes of *XMM-Newton* exposures included in the 5XMM catalogue.

Abbr.	Designation	Description
<i>MOS cameras:</i>		
PFW	Prime Full Window	covering full FOV
PPW2	Prime Partial W2	small central window
PPW3	Prime Partial W3	large central window
PPW4	Prime Partial W4	small central window
PPW5	Prime Partial W5	large central window
FU	Fast Uncompressed	central CCD in timing mode
RFS	Prime Partial RFS	central CCD with different frame time ('Refreshed Frame Store')
<i>pn camera:</i>		
PFW	Prime Full Window	covering full FOV
PLW	Prime Large Window	half the height of PFW/PFWE

Appendix B: Derivation of a systematic position error

To assess the systematic position uncertainty of an X-ray catalogue, its positional errors are circularised to have same errors on right ascension and declination and then compared with a reference catalogue. An example of such a distribution is shown in Fig. B.1. The linear part at large angular distances represents chance alignments and a pronounced peak at small separations corresponds to true associations. The linear term can be expressed as

$$N_{\text{rand}}(\theta) = N_X \cdot 2\pi \cdot \theta \cdot \rho \cdot d\theta, \quad (\text{B.1})$$

where $N_{\text{rand}}(\theta)$ is the number of random matches at angular separation θ , ρ is the sky density of the external catalogue and N_X the total number of X-ray sources. The number of true associations can be expressed as the superposition of N_X Rayleigh distributions, each of which describes the probability of an angular distance θ between the i X-ray source and its counterpart given the X-ray source's positional uncertainty σ_i

$$N_{\text{assoc}}(\theta) = F \cdot \sum_{i=1}^{N_X} \mathcal{R}_i(\theta | \sigma_i), \quad (\text{B.2})$$

where $\mathcal{R}_i(\theta | \sigma_i)$ is the Rayleigh distribution with scale parameter σ_i given by Equation 1 and the factor, F , is the fraction of X-ray sources with true associations in the external catalogue. Modeling the observed number of pairs at a given angular separation as the sum of the terms in Equations B.1, B.2 can constrain the parameters F , ρ , σ_i and hence σ_{sys}^2 via Equation 1.

The total number of X-ray vs external catalogue pairs at a given angular separation bin θ is a Poisson variate with expectation value $\lambda(\theta) = N_{\text{rand}}(\theta) + N_{\text{assoc}}(\theta)$. Therefore the likelihood of the model can then be expressed as the product of the Poisson probabilities at each angular separation bin

$$\mathcal{L} = \prod_{j=1}^{N_\theta} \mathcal{P}[N_j | \lambda(\theta_j)], \quad (\text{B.3})$$

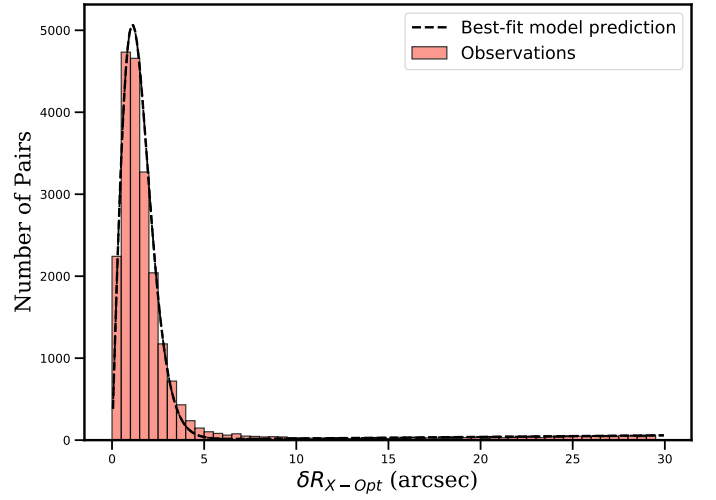


Fig. B.1: Distribution of the angular separation between 5XMM X-ray source positions and the Gaia/unWISE QSOs. The histogram is of the observed number of associations at a given angular separation bin, δR_{X-Opt} . The dashed line is the model described in the text for parameters fixed to the median of the corresponding posteriors.

where the index j is for all angular separation bins N_θ and $\mathcal{P}[N_j | \lambda(\theta_j)]$ is the Poisson probability of N_j pairs given the expectation value $\lambda(\theta_j)$.

The modeling assumes that the positional uncertainty of an X-ray source is given by Equation 1 and that the systematic uncertainty is the same for all sources. Although σ_{sys} can depend on source and instrumental details, these dependences are expected to be weak, and therefore assuming a single catalogue-wide value for this parameter is an acceptable approximation. The likelihood function of Equation B.3 is sampled using the ULTRANEST¹⁶ (Buchner 2021) nested sampling package based on the MLFriends (Buchner 2016, 2019) Monte Carlo algorithm. The inferred model parameters are ρ , F and σ_{sys} . We caution that this parametrisation assumes that the external catalogue has a homogeneous source density, ρ , across the sky. Although the sample selection (see below) minimises such variations, it is inevitable that ρ has an intrinsic scatter that is not accounted for in the current implementation.

The external astrometric catalogue are QSOs from Gaia and unWISE Data (Gaia/unWISE; Shu et al. 2019). We only consider Gaia/unWISE sources with probability being a QSO $\text{PROB_RF} > 0.8$ and G -band magnitude < 20.5 mag. The latter criterion is adopted to minimise variations in the sky density of QSO candidates because of the variable depth of the GAIA survey as a result of the scanning law of the mission. For this magnitude cut it is empirically found that the sky density of QSO candidates in the extra-galactic sky (Galactic latitudes $|b| > 20^\circ$) is nearly homogeneous. We limit the 5XMM catalogue to sources with emldetect detection likelihood $\text{EP_DET_ML} > 15$ (to increase the purity of the sample), that are not spatially extended (parameter $\text{EXTENT}=0$), are not close to CCD gaps or the edges of the field of view ($\text{PN_MASKFRAC} > 0.9$ or $\text{M1_MASKFRAC} > 0.9$ or $\text{M2_MASKFRAC} > 0.9$), have quality flags that do not indicate issues during the detection ($\text{SUM_FLAG}=0$) and lie outside the Galactic plane (Galactic latitude > 30 deg).

¹⁶ <https://johannesbuchner.github.io/ULtraNest>

Appendix C: Catalogues used for Cross-matching

Table C.1: Catalogues used for cross-matching. The last column indicates the processing method: either the ARCHES multi-catalogue statistical cross-match (s) or a simple positional cross-match (x).

Catalogue	Band	Xmatch mode
Gaia DR2	opt	s
AllWISE	ir	s
2MASS	ir	s
GALEX GR67	uv	s
NVSS	radio	x
FIRST	radio	x
Akari	farir	x
PanSTARRS DR1	opt	s
XMM-OM-SUSS5	uv	s

Appendix D: Quality of photometric redshifts

D.1. Cross validation

To characterize the performance of the photo-z algorithms we used a k -fold cross-validation method for TPZ, where the training sample is split in k equal parts, $k-1$ parts are used for training the algorithm and the remaining part is used for validation. In our case we used $k = 6$ and we did 30 iterations, where we randomly shuffled the training sample for each iteration, which gave us a total of 180 runs of training/validation. In the case of LePhare, we simply used the full training samples with the selected templates to validate the method.

For validation we compared the estimated photometric redshifts with the corresponding spectroscopic redshifts. To this end we make use of the most widely used statistical indicators, which are the following:

$$x = \Delta(z_{norm}) = \frac{z_{spec} - z_{phot}}{1 + z_{spec}}, \quad (D.1)$$

$$MAD(x) = Median(|x|), \quad (D.2)$$

$$\sigma_{NMAD}(x) = 1.4826 \times MAD(x), \quad (D.3)$$

$$\eta = \frac{N_{outliers}}{N} \times 100, \quad (D.4)$$

where σ_{NMAD} is the normalised median absolute deviation (MAD), and η is the percentage of catastrophic outliers. A source is considered a catastrophic outlier if $|x| > 0.15$.

Figure D.1 shows the fraction of outliers estimated for TPZ and LePhare, for the different optical surveys we employed, the morphology of the source (i.e extended or point-like) and the available photometry.

D.2. Impact of quality flags

TPZ and LePhare both provide a full estimate for the probability density function of the photometric redshift (PDFZ). The overall shape of the PDFZ is a useful indicator for the reliability of each estimated photo-z. We calculated several PDFZ-derived parameters defined as follows:

- $zConf$: The integral of the PDFZ in the interval $\pm(1 + z_{phot}) \times rms$, centred in z_{phot} . z_{phot} is the mode of the PDF (the absolute maximum and the value chosen as the photometric

redshift of the source) and rms is the intrinsic dispersion of the method, which depends on the employed training sample. For our sample we have used $rms = 0.06$. A high value of $zConf$ means that the probability is highly concentrated around the estimated photo-z.

- NP (Number of peaks): Number of local maxima (peaks) in the PDF.
- PS (Peak strength): $1-P2/P1$, where P1 is the probability density of the highest local maximum in the PDF, and P2 is the second maximum peak. If the PDF is unimodal ($P2=0$) or $P2 \ll P1$, $PS \approx 1$.

Additional quality flags are included in the catalogue for each source based on the photometry available. For TPZ we estimate how representative the training sample is for a given source, based on the position occupied in the colour space (TPZ_FLAG_TSCS_ALL, TPZ_FLAG_TSCS_ANY). In the case of LePhare, we provide the number of photometric points used during the fitting (LPH_NBANDS).

Ruiz et al. (2018, see Sect. 5.1.2 and Figs. 9, 10) presented extensive tests showing the impact on the photometric redshifts statistics of selecting sources based on these quality flags. Fig. D.2 illustrates this effect for the 5XMM results. The top panel shows the photo-z performance for the full 5XMM sample with available spectroscopic redshifts, without any quality filtering, yielding a catastrophic outlier fraction of 16.5%.¹⁷ The middle panel shows the results after applying two simultaneous selection criteria: full photometric coverage across the optical, NIR, and MIR bands, and a Peak Strength value of $PS > 0.9$. This combined filtering reduces the outlier fraction substantially, to 2.8%, but at a significant cost in completeness, the sample size drops from 31 829 to 7 446 sources. Notably, 72.8% of the rejected sources are not catastrophic outliers, meaning that the photometric coverage cut discards a large number of well-recovered objects. The bottom panel shows the effect of applying the $PS > 0.9$ criterion alone, without imposing any requirement on photometric band coverage. This less restrictive selection still yields a meaningful reduction in the outlier fraction, to 4%, while retaining a larger portion of the sample: 23,003 out of 31,829 sources, with only 17% of the rejected sources being non-outliers.

¹⁷ Note that this value is not representative of η for the full catalogue. A more accurate method to estimate this value is explained in Sect. D.1.

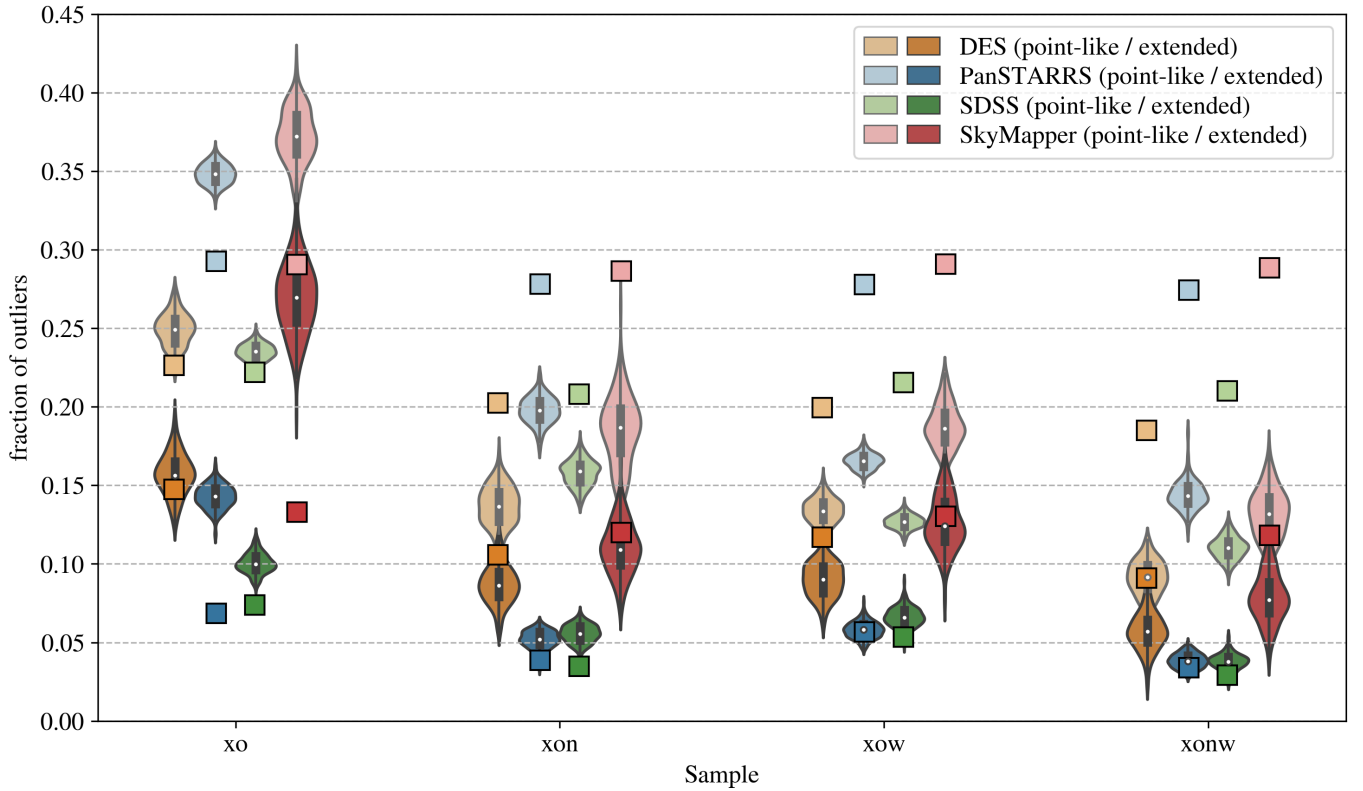


Fig. D.1: Fraction of outliers for our photometric redshifts estimated using our training samples. The violins show the resulting distribution for the TPZ algorithm and the square symbols show the results using LePhare. We show the results for each optical survey used in our catalogues: DES (orange), PanSTARRS (blue), SDSS (green) and SkyMapper (red). Results are split by the optical morphological classification: extended sources (bright colours) and point-like sources (dim colours). We group the results in four different classes, according to the total available photometric: Only optical magnitudes (xo), optical and near-infrared magnitudes (xon), optical and mid-infrared magnitudes (xow), and optical, near-, and mid-infrared magnitudes (xonw).

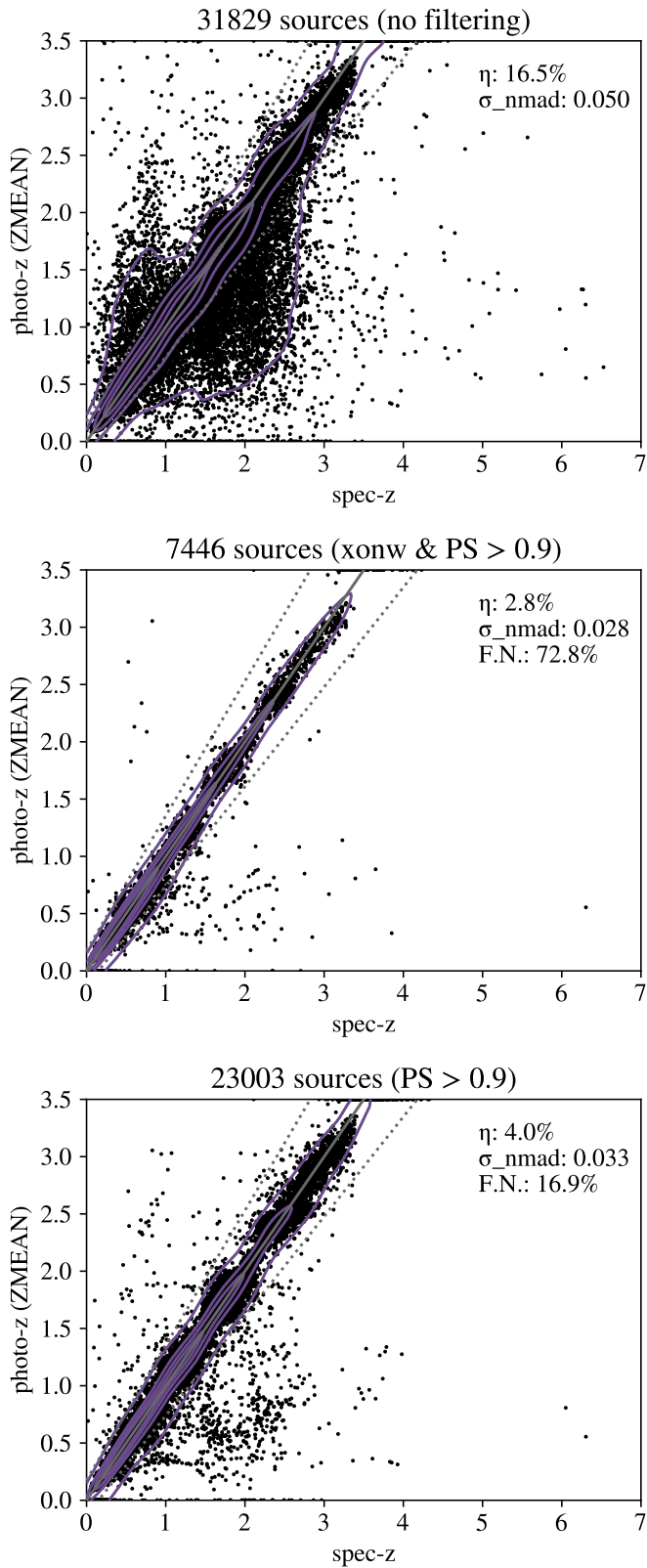


Fig. D.2: Comparison of spectroscopic redshifts with photometric redshifts estimated via TPZ for sources in the 5XMM catalogue. Top panel: full sample. Middle panel: sources with photometric coverage in the optical, NIR, and MIR bands, restricted to quality flag PS > 0.9. Bottom panel: sources with quality flag PS > 0.9, regardless of photometric band coverage. Each panel reports the outlier fraction, the normalised median absolute deviation, and the false negative rate (F.N.) for the corresponding sample.