

# Seven Decades of Quasar Photometry: Paper I - construction of 7-DQ database and initial science results

Harry Rendell-Bhatti,<sup>1</sup> Andy Lawrence,<sup>1</sup>

<sup>1</sup>*Institute for Astronomy, University of Edinburgh, Royal Observatory, Blackford Hill, Edinburgh EH9 3HJ, UK*

Accepted XXX. Received YYY; in original form ZZZ

## ABSTRACT

We describe the construction of the 7-DQ quasar database, consisting of multi-colour photometry for half a million quasars spread over seven decades, and present some initial scientific analyses using this database. The sample is based on the SDSS DR14 quasar sample, and combines SDSS, PanSTARRS, ZTF, and SuperCOSMOS photometry. We describe the transformation to a common photometry system, some cleaning to produce a final light curve database, and the construction of pair-wise dataset to speed up later analyses. We also describe the construction of a control sample of stars observed over the same period. We use the database for three initial example analyses. First, we examine the claim by [Caplar et al. \(2020\)](#) that the mean magnitude of quasars is drifting with epoch, suggesting that variability of quasars is not a stationary process. We do not confirm this claim with our database. Second, we calculate the skewness (third moment) and the excess kurtosis (i.e. fourth moment compared to Gaussian expectation), and examine how these quantities depend on lag between observations. Skewness is consistent with zero at all lags. For the star control sample excess kurtosis is greater than zero, but does not change with lag, almost certainly indicating somewhat non-Gaussian photometric errors and/or remaining contamination by outliers. For the quasars, excess kurtosis systematically declines with lag. On most timescales, variability is distinctly non-Gaussian, in the sense of having heavier tails. On the longest timescales, variability seems to be consistent with Gaussian. Finally, we take an initial look at the Structure Function (SF; i.e. behaviour of the second moment with lag). The ensemble SF is consistent with a single power law with index  $\beta = 0.35 \pm 0.004$  over many decades, and strongly inconsistent with the form expected from a Damped Random Walk. There is no evidence for a long timescale turnover, i.e. no sign of a characteristic timescale.

**Key words:** quasar – structure function – ensemble

## 1 INTRODUCTION

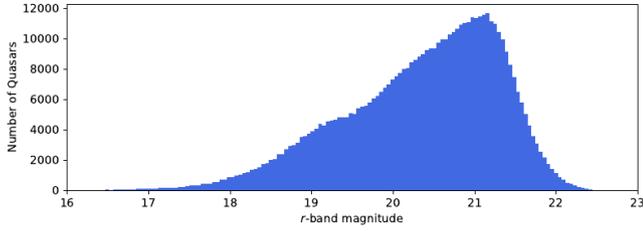
Variability is a fundamental aspect of the emission from Active Galactic Nuclei (AGN), including quasars, probing size scales and physical processes that are otherwise impossible to access. (See for example review by [Lawrence \(2016\)](#)). As well as the study of individual AGN, considerable progress has been made through statistical studies of large samples observed over long time periods ([Hook et al. \(1994\)](#); [de Vries et al. \(2003, 2005\)](#); [Morganson et al. \(2014\)](#); [Li et al. \(2018\)](#)).

In the work described here we present the most ambitious massive variability study so far, using half a million quasars, and combining data from multiple surveys, spread over seven decades, reduced to a common photometric system. We refer to our dataset as the 7-DQ database.

In this paper (Paper I) we describe the construction of the database, some initial science analyses, and the public availability of the data. In Paper II we analyse the behaviour of the quasar Structure Function in some detail, including a sub-ensemble analysis versus luminosity and black hole mass, and an examination of the time asymmetry of the Structure Func-

tion. In later papers we plan to estimate variability parameters for individual quasars and test dependence on a variety of astrophysical parameters; and to define a sample of Extreme Variables, and examine their properties.

In Section 2 of this paper we define the 7-DQ quasar sample, the accompanying control sample of stars, and the bright sub-sample used in some circumstances. In Section 3 we set out the various data sources used to construct our database, and the way these data sets were treated. In Section 3 we describe the various photometry data sources we have used, and in section 4 we describe how colour transformations were used to transform the photometry to a common system. Section 6 sets out the construction of a database of all  $(\Delta m \Delta t)$  pairs, which is crucial to further analysis, and shows the resulting  $\Delta m$  distributions versus rest-frame lag. Finally, in Sections 7, 8 and 9 we describe some initial scientific applications of the 7-DQ database: testing the claim of a mean drift of quasar brightness with time, looking at how skewness and kurtosis depend on rest-frame lag, and construction of the ensemble structure function, which we show clearly does not follow the form expected for the popular “damped



**Figure 1.** Distribution of  $r$ -band magnitudes for all DR14Q quasars, from the DR14Q catalogue.

random walk” model. More detailed analysis of the structure function is presented in Paper II.

## 2 SAMPLE AND SUB-SAMPLE DEFINITIONS

In this section, we present the two main samples that form the starting point for the 7-DQ database. The main quasar sample is identical to the SDSS Data Release 14 Quasar catalogue (DR14Q, Pâris et al. 2018). However we also define a control sample of non-variable stars which is a subsample of those defined by Ivezić et al. (2007). Throughout the paper, various sub-samples are defined for scientific purposes, but in technical testing it is also useful to define a bright quasar subsample.

### 2.1 7-DQ quasar sample

We define the 7-DQ quasar sample to be the same as the SDSS Data Release 14 Quasar catalogue (DR14Q, Pâris et al. 2018), a catalogue of 526,356 spectroscopically confirmed quasars. (DR14 was the most recent release at the time this project was started.)

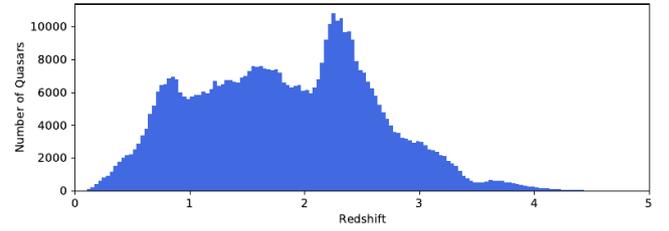
This sample constitutes the super-set for our analysis, and we use it in two ways; first, in our ensemble analysis, we use all available photometry from the full set of quasars in this sample, obtained as described in section 3. Second, we derive various subsamples based on quasar properties such as bolometric luminosity, for sub-ensemble analyses reported in later papers.

In order to track quasars from this sample, we assign each of them a unique identifier, `uid`, which is simply the row number of the DR14Q sample once sorted by RA. This identifier is arbitrary and only occasionally explicitly mentioned, but it is very convenient from a computational and analysis perspective.

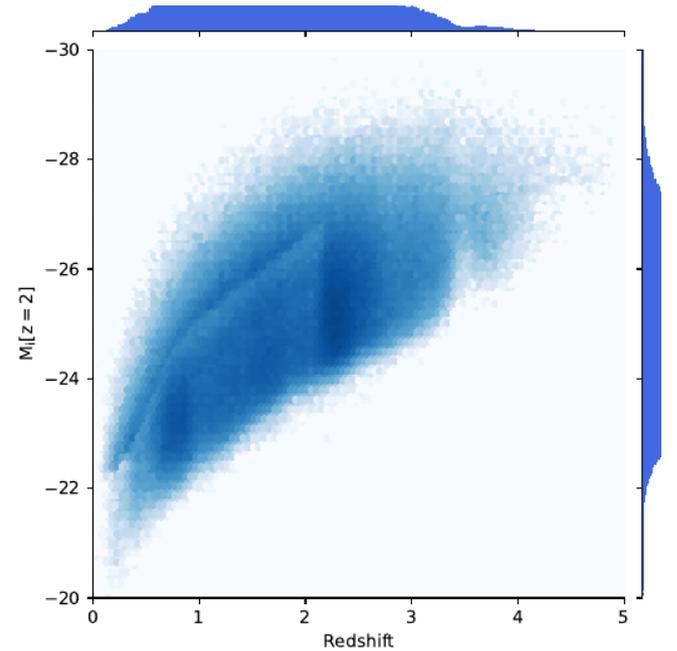
The apparent magnitude and redshift distributions of the DR14Q sample, as given by the DR14Q catalogue, are shown in Figure 1 and Figure 2. In Figure 3 we show the density map of the quasars in the luminosity-redshift ( $L - z$ ) plane.

### 2.2 7-DQ star sample

Here we define a control sample of non-variable stars. This enables us to cross-calibrate photometry between surveys and also to provide a comparison for variability studies of the quasar sample, as well as estimating photometric errors on the SuperCOSMOS data.



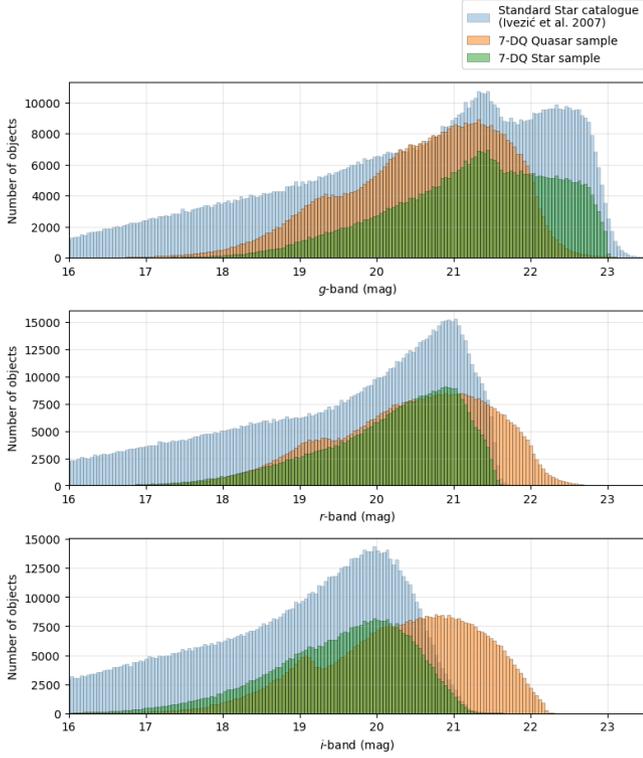
**Figure 2.** Distribution of redshifts for all DR14Q quasars, from the DR14Q catalogue. The peak at  $z = 2.2$  is due to the BOSS quasar selection criteria.



**Figure 3.** Density map of the DR14Q quasars in the  $L - z$  plane, where luminosity is expressed as absolute  $I$ -band magnitude at  $z = 2$ . This absolute magnitude assumes  $H_0 = 70 \text{ km s}^{-1} \text{ Mpc}^{-1}$

Lost some text here???? of the quasar sample simultaneously in all of  $g$ ,  $r$  and  $i$  bands, or at least come closer to doing so. By performing this matching, we effectively balance both random and systematic uncertainties in the photometry of both datasets. This ensures that any observed noise in star photometry can serve as an appropriate measure for the baseline level of uncertainty in quasar photometry.

We achieved this by taking a weighted sample of stars determined by the quasar magnitude distribution across all three bands. First, the 7-DQ quasars are grouped into predefined magnitude bins, *separately* for the  $g$ ,  $r$ , and  $i$  bands (i.e., not simultaneously in  $g, r, i$ ), based on their corresponding mean magnitude in SDSS, resulting in bin counts for each band denoted as  $(n_g, n_r, n_i)$ . These bins specify the number of quasars per magnitude bin, per band. Then, for every star, we identify the corresponding three magnitude bins based on their SDSS mean magnitude in each band. By using the combined counts from the quasars  $(n_g + n_r + n_i)$ , we establish the weight for each star. Finally, we sample 400,000 stars from the full catalogue, using these weights. We then removed 78



**Figure 4.** Magnitude distributions of the quasar and star samples. The full Stripe 82 Standard Star Catalogue is shown in light blue, while our matched subset is shown in green.

objects that had not been observed in SDSS at least four times to ensure the mean magnitudes of the stars were reliable, resulting in a final sample of 399,922 stars.

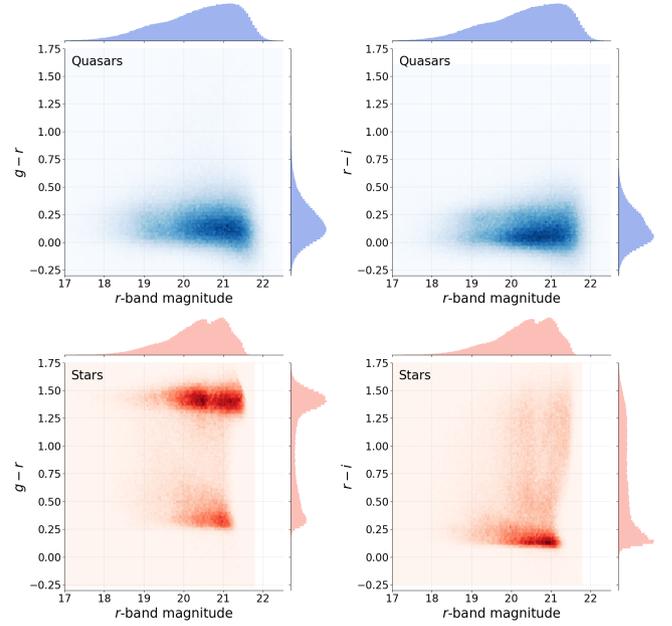
The magnitude distribution showing the result of the matching process is presented in Figure 4. The magnitude distribution is well-matched in the  $r$ -band, but not as well-matched in the  $g$  and  $i$  bands. This is because the colour distributions of the quasar and stars are significantly different. To illustrate this, colour-magnitude plots are shown for the two samples in Figure 5.

### 2.3 Bright Quasar sample

Our quasar sample is flux limited. Therefore, we expect there to be some selection effects. We define a “Bright Quasar Sample” whose median SDSS magnitude is  $g_{\text{SDSS}} < 20$ , and a corresponding bright star sample, which we used for testing and assessing the effects of Malmquist bias. This leaves 139,233 quasars, and 64,626 stars.

## 3 DATA SOURCES

In this section we describe how we compiled photometry spanning 70 years from multiple surveys, for both the quasars and the star control sample, as well as associated quasar properties.



**Figure 5.** Colour-magnitude distributions of the quasars (top row, in blue) and stars (bottom row, in red), for  $g - r$  and  $r - i$  colors against the  $r$ -band.

Filter	[H]			
	5 $\sigma$ depth			
	SDSS	Pan-STARRS1	ZTF	SuperCOSMOS
$u$	22.15	-	-	-
$g$	23.13	22.0	20.8	21.2
$r$	22.70	21.8	20.6	20.3
$i$	22.20	21.5	19.9	18.9
$z$	20.71	20.9	-	-
$y$	-	19.7	-	-

**Table 1.** Single-epoch 5 $\sigma$  imaging depths for SDSS, Pan-STARRS1, ZTF, and SuperCOSMOS in  $ugrizy$  bands. Note that, since each plate in SuperCOSMOS has a different limiting magnitude, the SuperCOSMOS depths represent an average over the plate-emulsion combinations corresponding to each  $grizy$  band.

### 3.1 SDSS photometry

SDSS photometry was obtained from the Catalogue Archive Server Jobs System (CasJobs)

### 3.2 PanSTARRS photometry

Pan-STARRS photometry was also obtained via CasJobs using a similar query to SDSS. Since the photometry is provided in fluxes (Jy), I convert to  $AB$  magnitudes using the following relations,

$$m = -2.5 \log_{10}(f) + 8.9 \quad (1)$$

$$\sigma_m = \frac{2.5}{\ln 10} \left( \frac{\sigma_f}{f} \right). \quad (2)$$

Each source in the Pan-STARRS database is assigned a unique Pan-STARRS object ID, denoted `ps_objID`. There were 164 instances where multiple Pan-STARRS unique sources were mapped to a single uid. This is either caused

by two objects in close proximity, or a mishap in the Pan-STARRS stacking, which results in the same source being assigned multiple `ps_objID`. In either case, 164 makes up a negligible fraction of our total sources, and thus we take the `ps_objID` of the object closest to the query coordinates and omit the rest.

### 3.3 ZTF photometry

To collect photometry of our samples from ZTF, the IRSA service was used to find sources which match our objects. ZTF often has multiple identifiers per object, therefore we used the source table to query ZTF sources which correspond to our `uid` unique identifier. Obtaining photometry from ZTF is less trivial than the others, since there is a limit on the number of objects when using the query form. Therefore, I wrote a script to automate the data collection. This query was rerun in April 2023 once Data Release 17 was available to obtain the most up-to-date observations for our sources.

### 3.4 SuperCOSMOS photometry

Photometry from the SuperCOSMOS Sky Survey was acquired using the SuperCOSMOS Science Archive (SSA) (<http://ssa.roe.ac.uk>). There are several databases which I queried to get the photometry in the desired form. The Detection Table includes all observations across the SuperCOSMOS Sky Survey across multiple bands. I queried this table for all observations within  $1.5''$  of the quasar and star coordinates. While I used a pairing radius of  $1''$  for the other surveys, I decided to use  $1.5''$  for SSS due to its lower precision astrometry. I then joined this result to the Plate Table to find which observations belong to which plates, in order to separate observations by survey. Due to the sky coverage of each individual survey, they each return a different number of observations. Since the 7-DQ stars cover a narrow range of sky, there are only observations from four surveys that make up the SuperCOSMOS Sky Survey.

### 3.5 Quasar properties

In order to obtain properties such as black hole mass, luminosity and Eddington ratio for the 7-DQ quasars, we cross-matched these sources with the DR16 catalogue of quasar properties provided by [Wu & Shen \(2022\)](#).

## 4 TRANSFORMATION TO A COMMON PHOTOMETRY SYSTEM

In order to create 70 year lightcurves, we need to combine the photometry from the various surveys. However, they have subtly different photometric systems and so we need to convert the magnitudes to a common photometric system. We decided to shift everything to the Pan-STARRS system, since it has the most reliable photometry and required the fewest total transformations.

Transformations between one system and another are of course colour-dependent. Ideally we require the colour of the object to be known at the same epoch as the measurement being transformed. Unfortunately, we do not have colours for every epoch where we wish to make a transformation, and

Filterband	$a_0$	$a_1$	$a_2$
<i>g</i>	-0.011	-0.125	-0.015
<i>r</i>	+0.001	-0.006	-0.002
<i>i</i>	+0.004	-0.014	+0.001

**Table 2.** Polynomial coefficients for Equation 3 used to transform SDSS magnitudes, from [Tonry et al. \(2012\)](#).

therefore an approximation must be made. There are two possible approaches. The first involves using non-simultaneous observations to estimate the colour of an object at a particular epoch. The second approach uses the colour at one epoch for all future transformations. The first approach is not suitable for quasars since they are inherently variable, and it is likely the quasar will have varied significantly between the two measurements, resulting in a colour that may be significantly different from the true colour at either of the two measurements. We decided to use the second approach, as it is suitable under the assumption that the colours of our objects stay roughly constant. We use the colours supplied by SDSS because only this survey provides simultaneous imaging. Therefore, we used the mean colours from SDSS for all colour transformations. This assumes that the colour of the objects stay roughly constant throughout the light curve. Although this is a safe assumption for the stars, it will introduce some uncertainty for the quasars. To estimate this uncertainty, we investigated the spread of quasar colours in SDSS for quasars that had been observed at least twice. The mean standard deviation of quasar colours are  $0.081079$  mag and  $0.086875$  mag for  $(g-r)$  and  $(r-i)$ , respectively. This was calculated using 1,167,613 measurements of each  $(g-r)$  and  $(r-i)$  from 258,319 quasars which had been multiply imaged by SDSS. These deviations are small enough to justify our method.

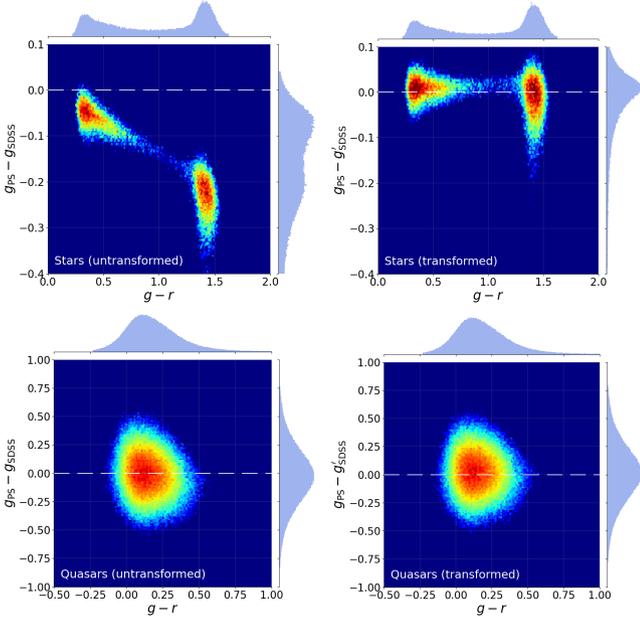
In the remainder of this section, we describe the colour transformations of each survey into the Pan-STARRS system and assess the validity of each transformation using colour-magnitude plots.

### 4.1 Transformation of SDSS magnitudes

The colour transformations from SDSS magnitudes to Pan-STARRS magnitudes are provided by [Tonry et al. \(2012\)](#). The transformation is a second order polynomial in  $g-r$ ,

$$m_{\text{PS1}} = m_{\text{SDSS}} + a_0 + a_1(g-r)_{\text{SDSS}} + a_2(g-r)_{\text{SDSS}}^2, \quad (3)$$

Coefficients are provided in Table 2. These coefficients were derived for stellar SEDs - however, we find that they work well for the quasar population. To illustrate the effectiveness of the transformations in each band, we produced magnitude-colour diagrams before and after the transformations were applied. An example, using the *g*-band, is shown in Figure 6 for the star and quasar population. Figure 6 shows that the colour dependence on magnitude is almost entirely removed in the star sample. This is also true for the quasar sample, however, it is less obvious since quasars are variable and therefore the *y*-axis limits are much larger.



**Figure 6.** 2-D histograms illustrating the effectiveness of the colour transformations between SDSS and PanSTARRS systems, for the star sample (top row) and quasar sample (bottom row). Left and right columns are untransformed and transformed, respectively. In each case,  $g - r$  is the SDSS colour, and the y-axis is the difference between SDSS and PanSTARRS magnitudes.

#### 4.2 Transformation of ZTF magnitudes

The transformations from ZTF to Pan-STARRS are provided by Masci et al. (2019),

$$g_{PS1} = g_{ZTF} + c_g(g - r)_{SDSS}, \quad (4)$$

$$r_{PS1} = r_{ZTF} + c_r(g - r)_{SDSS}, \quad (5)$$

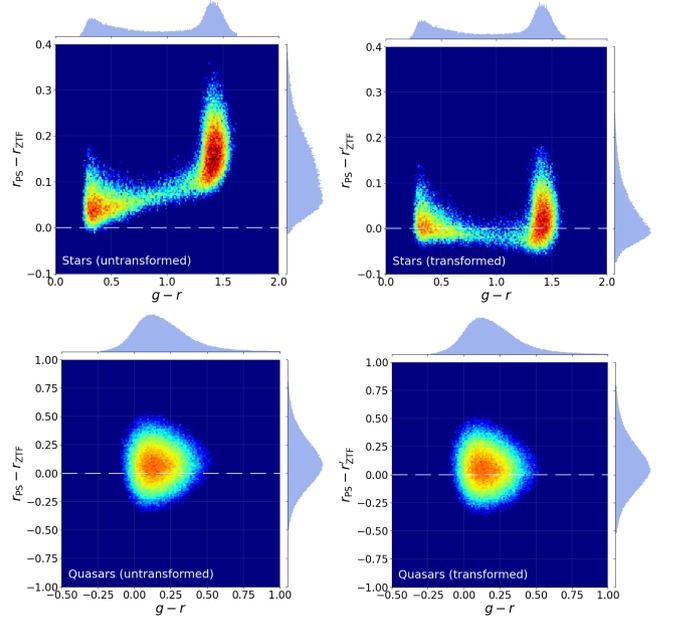
$$i_{PS1} = i_{ZTF} + c_i(r - i)_{SDSS}, \quad (6)$$

$$(7)$$

where  $c_f$  is the colour coefficient for a given filter ( $f = g, r, i$ ), which corresponds to the column `clrcoeff`. Each image is assigned its own unique solution for  $ZP_f$ ,  $c_f$ , which respectively correspond to the intercept and slope from a linear fit using matches of ZTF sources to a subset of calibrators in PS1. As with SDSS, we produced magnitude-colour diagrams before and after the transformations were applied. An example, using the  $r$ -band, is shown in Figure 7 for the star and quasar population. The transformation removes most of the magnitude-colour dependence for both populations.

#### 4.3 Transformation of SuperCOSMOS magnitudes

The case of transforming SuperCOSMOS magnitudes requires more care. One possibility is to use transformations produced by Peacock et al. (2016). They provide transformations of  $B_J$ ,  $R_F$ , and  $I_N$  bands to standard SDSS  $gri$  bands. However, since their colour corrections are designed for galaxies, we found that they performed poorly on 7-DQ stars and quasars; the transformations shifted 7-DQ star and quasar magnitudes away from their corresponding magnitudes measured by Pan-STARRS. In the absence of any other suitable



**Figure 7.** As per Fig 6, but for the transformation between PanSTARRS and ZTF magnitudes.

published transformations, we decided to create our own. We used the 7-DQ star photometry to produce linear transformations of  $B_J$ ,  $R_F$ ,  $E$ , and  $I_N$  to Pan-STARRS  $gri$  magnitudes. An additional benefit to creating our own transformations is that we can fit over a colour range that is more suitable for the quasars. For example, when creating transformations which use the  $g - r$  colour term, we pick stars in the range  $-0.1 < g - r < 0.7$  which is the extent of the quasar  $g - r$  colours. This results in a more accurate transformation for the 7-DQ quasars.

To create our own transformations, we calculated magnitude offsets between stars in the SuperCOSMOS and Pan-STARRS data, and performed linear regression of the offsets as a function of colour. The coefficients of the linear regression were found to be:

$$g_{PS1} = B_J - 0.222(g - r)_{SDSS} - 0.200, \quad (8)$$

$$r_{PS1} = R_F + 0.091(g - r)_{SDSS} + 0.151, \quad (9)$$

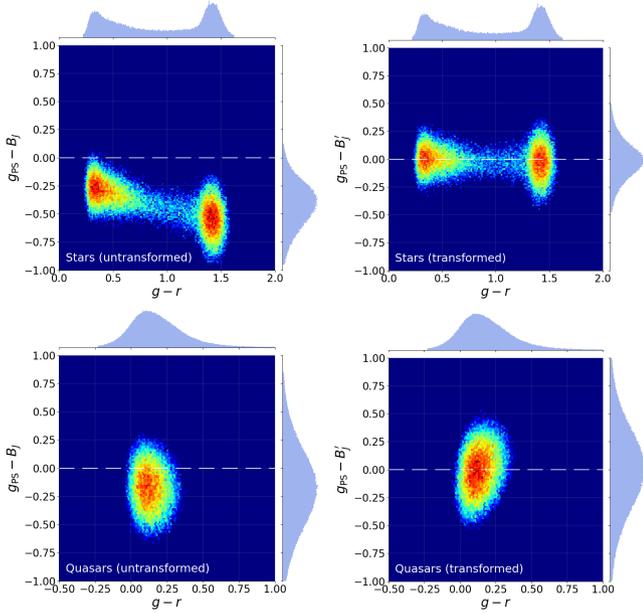
$$r_{PS1} = E - 0.210(g - r)_{SDSS} + 0.338, \quad (10)$$

$$i_{PS1} = I_N + 0.126(r - i)_{SDSS} + 0.436. \quad (11)$$

Note that, since the sky footprint of the 7-DQ stars is quite narrow, we were only able to use photometry from four sub-surveys of SuperCOSMOS to generate these transformations. Fortunately, these 4 sub-surveys covered all four bands of the photometric plates (i.e.,  $B_J$ ,  $R_F$ ,  $E$ , and  $I_N$ ). The results of the transformation for the 7-DQ stars and quasars are shown in Figure 8, using the  $g$  and  $B_J$  band as an example.

#### 4.4 Effectiveness of transformations

To evaluate the performance of colour transformations, we examined the distribution of the residuals between median magnitudes for corresponding objects between surveys. This procedure is conducted separately for all three  $gri$  bands, and both untransformed and transformed magnitudes. The



**Figure 8.** As per Fig 6, but for the transformation between PanSTARRS and SuperCOSMOS magnitudes.

Survey pair	Band	Residual (untransf.)	Residual (transf.)
PS-SDSS	<i>g</i>	$-0.237 \pm 0.207$	$-0.076 \pm 0.170$
	<i>r</i>	$-0.015 \pm 0.042$	$-0.007 \pm 0.042$
	<i>i</i>	$+0.001 \pm 0.030$	$+0.011 \pm 0.031$
PS-ZTF	<i>g</i>	$-0.012 \pm 0.119$	$+0.034 \pm 0.121$
	<i>r</i>	$+0.106 \pm 0.065$	$+0.007 \pm 0.040$
PS-SSS	<i>g</i>	$-0.338 \pm 0.186$	$+0.039 \pm 0.161$
	<i>r</i>	$+0.328 \pm 0.220$	$+0.134 \pm 0.239$
	<i>i</i>	$+0.608 \pm 0.200$	$+0.062 \pm 0.183$

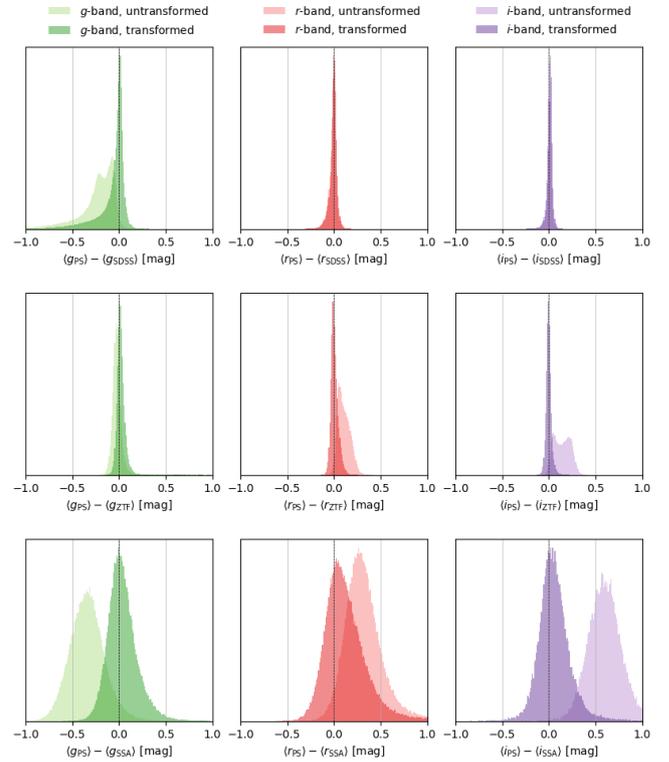
**Table 3.** Mean residuals of the star sample in the *g*, *r* and *i* bands before and after applying colour transformations

Survey pair	Band	Residual (untransf.)	Residual (transf.)
PS-SDSS	<i>g</i>	$-0.035 \pm 0.301$	$+0.003 \pm 0.296$
	<i>r</i>	$-0.026 \pm 0.275$	$-0.025 \pm 0.275$
	<i>i</i>	$-0.021 \pm 0.256$	$-0.022 \pm 0.257$
PS-ZTF	<i>g</i>	$+0.034 \pm 0.264$	$+0.047 \pm 0.265$
	<i>r</i>	$+0.075 \pm 0.235$	$+0.051 \pm 0.235$
	<i>i</i>	$+0.074 \pm 0.210$	$+0.047 \pm 0.210$
PS-SSS	<i>g</i>	$-0.099 \pm 0.380$	$+0.097 \pm 0.388$
	<i>r</i>	$+0.281 \pm 0.375$	$+0.062 \pm 0.373$
	<i>i</i>	$+0.547 \pm 0.383$	$+0.140 \pm 0.382$

**Table 4.** Mean residuals of the quasar sample in the *g*, *r* and *i* bands before and after applying colour transformations

results are shown in Figures 9 and 10. Since quasars are variable, these distributions are much wider than for the star sample. Therefore, the effect of the transformation is not as clear as it is for the stars (with the exception of SSS). The mean residuals and their errors (calculated as the standard deviation of the residuals) are presented in Tables 3 and 4 respectively.

Since the mean residuals in each band are reduced in the majority of cases, it is clear that the transformations are per-



**Figure 9.** Residuals between Pan-STARRS and the other surveys for the star sample. Residuals are calculated as median magnitudes per object before and after transformation. Note that the angled brackets denote the median average.

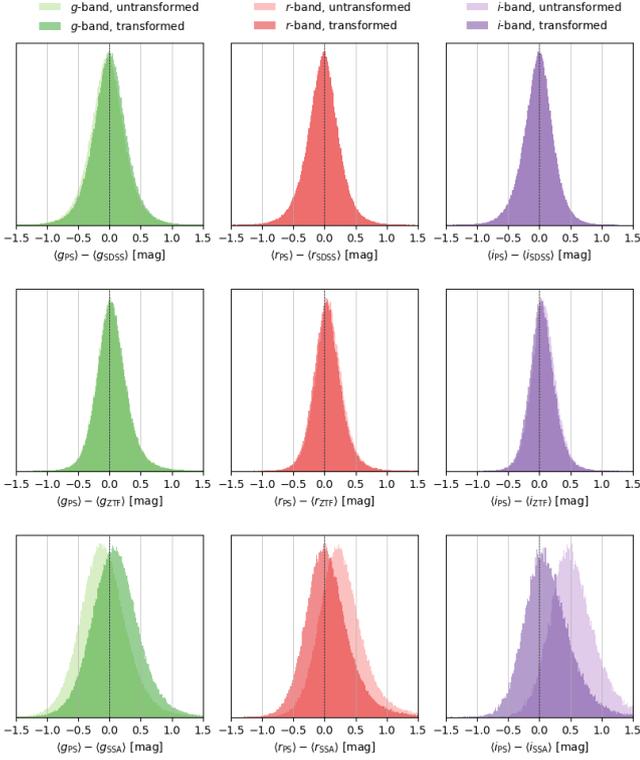
forming as expected and are now within the same photometric system with sufficient precision.

## 5 FINAL CLEANED LIGHT CURVE DATABASE

A few further steps are required to produce a final light-curve database. First, the inclusion of estimated magnitude errors. Second, removal of bad data which may for example be due to instrumental or pipeline errors (e.g. issues with photometric or astrometric calibration, or mismatching of nearby objects), environmental conditions (e.g. poor seeing conditions, or cosmic rays), or data and computational errors (e.g. issues with data transfer, storage or retrieval). Any such data cleaning is a compromise between keeping too much bad data, and throwing away too much good data, because of course we cannot tell the difference with complete reliability. Here we describe our approach.

### 5.1 Estimating SuperCOSMOS magnitude errors

The SDSS, PanSTARRS and ZTF databases come with their own estimated photometric errors, and we have simply adopted these. The SuperCOSMOS data however does not provide magnitude errors on individual measurements. We have therefore estimated these empirically, using the 7-DQ star database. We calculated the magnitude dependence on error empirically, using the star sample. Since the stars are non-variable, the observed variance of magnitude differences

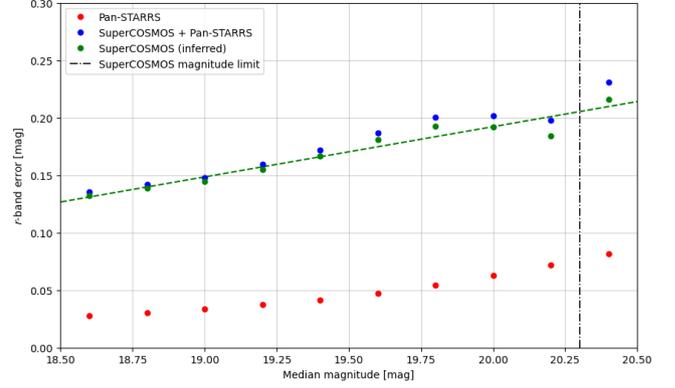


**Figure 10.** Residuals between Pan-STARRS and the other surveys for the quasar sample. Residuals are calculated as median magnitudes per object before and after transformation. Note that the angled brackets denote the median average.

is solely due to photometric noise. Thus, the variance,  $\sigma_{\text{tot}}^2$ , of the magnitude differences,  $\Delta m_{\text{SSS-PS}} = m_{\text{SSS}} - m_{\text{PS}}$  between the SuperCOSMOS and Pan-STARRS observations is due to the combined errors of these two surveys. By approximating photometric noise as a function of magnitude, we may write the total combined photometric error as

$$\sigma_{\text{PS-SSS}}^2(m) = \sigma_{\text{PS}}^2(m) + \sigma_{\text{SSS}}^2(m). \quad (12)$$

We were able to estimate  $\sigma_{\text{PS-SSS}}^2(m)$  as the width of the  $\Delta m_{\text{PS-SSS}}$  distribution. There are two methods to estimate  $\sigma_{\text{PS}}^2(m)$ . The first involves calculating the average photometric error with increasing magnitude bins. The second involves calculating the width of the  $\Delta m_{\text{SSS-SSS}}$  distribution with increasing magnitude bins. After trying both of these methods, we found that the average photometric noise was on average 50% smaller than the observed scatter of the  $\Delta m_{\text{SSS-SSS}}$  distribution. We opted to take the second method in case the photometric errors were slightly underestimated. Having approximated both  $\sigma_{\text{PS-SSS}}^2(m)$  and  $\sigma_{\text{PS}}^2(m)$ , the difference represents the average photometric noise as a function of magnitude for SuperCOSMOS. We added these errors back into each SuperCOSMOS observation depending on the magnitude of the object, using a linear fit to interpolate between magnitude bins. This process gives a good estimate of the errors and allows us to weight observations in favour of objects which have better photometry. Figure 11 illustrates the process of calculating  $\sigma_{\text{SSS}}^2(m)$  in magnitude bins, with a linear fit.



**Figure 11.** Photometric error-magnitude plot for Pan-STARRS and SuperCOSMOS for the star population

## 5.2 Magnitude cuts

First, we remove unphysical observations by omitting observations whose magnitudes are outside the range  $15 < m < 25$  and magnitude errors are outside the range  $0 < \sigma < 2$ . Although these are generous bounds, they are designed to remove observations that have had measurement issues, such as those which contain NaNs or values such as  $-999$  (a typical entry for a null measurement). This step removes a very small fraction of the photometry ( $< 1\%$ ).

Second, we imposed magnitude constraints on each survey based on the limiting magnitude (shown previously in Table 1). Since the surveys have different imaging depths, some surveys are able to provide reliable photometry of fainter objects while others are not. For example, the  $5\text{-}\sigma$  depths for ZTF and SDSS are 20.6 and 22.7 in the  $r$ -band respectively, therefore a quasar of magnitude 21 would have reliable photometry from SDSS, but not ZTF. Moreover, because quasars are variable, ZTF would more likely report observations of this object if it exhibited a period of brightening during the observations. The resulting effect is a systematic overestimation of brightness for dimmer quasars in shallower surveys, known as Malmquist bias. To remove this bias, we decided to remove photometry of an object from a particular survey if the median magnitude of that object across the light curve is fainter than the limiting magnitude of that survey. We opted for the median as it is more robust to outliers. Compared to other data cleaning steps outlined in this section, these magnitude constraints remove the most amount of photometry by far. Nevertheless, it is a necessary trade-off to drastically reduce the effect of Malmquist bias in my analyses. On average, these constraints remove up to 30% of quasars for SuperCOSMOS and ZTF, and up to 3% of quasars in SDSS and Pan-STARRS. The equivalent fractions for the stars are roughly half of those for the quasars.

## 5.3 Outlier detection

We developed an outlier detection algorithm to automatically remove unphysical observations. We take a few steps to clean the data; almost all quasars (96.8%) have at least two observations within the same night at some point in their light curve. However, since observed variability on timescales this short will be dominated by noise, we combine nightly

observations into a single data point by taking the median of the measured magnitudes. The median was chosen over the mean as it is less sensitive to outliers. We then remove observations within a light curve whose magnitude deviates  $7\sigma$  from the median magnitude of that light curve, where  $\sigma$  is the standard deviation of the magnitudes in that light curve. Although this is a generous constraint, quasars are known to exhibit dramatic changes in brightness, and we did not want to unintentionally remove observations of extreme variability. This only removes a small fraction of obvious bad data, so we do a second pass with a more sophisticated outlier detection algorithm to further remove unphysical observations.

Bad data points are often easy to spot by eye. By leveraging contextual information from adjacent data points, examining data in different wave bands, and applying an intuitive understanding of physical plausibility, a person can confidently determine the authenticity of an observation with a relatively high degree of certainty by eye. However, producing an algorithm to carry out the same process with similar accuracy is not a simple task. We created our own outlier detection algorithm to filter out bad points which were not removed from the first pass. Our algorithm involves a rolling window of a width of 5 data points to calculate the absolute deviation of a particular point away from the median of the window. If this deviation is greater than 2.6 times the interquartile range of the distribution of magnitudes in the light curve, it is deemed an outlier. The number 2.6 seems arbitrary, but we calculated it empirically after looking at the distribution of absolute deviations across a sample of 5000 light curves which are representative of the rest of the light curves. This number roughly marks the boundary between high variability and obvious outliers. A demonstration of our algorithm is shown in Figure 12.

#### 5.4 Final light curve database

Table 5 shows an extract of the resulting final light curve database. Each line represents a measurement at single epoch and a single filter, from a specific survey. The “uid” column is the unique quasar identifier as explained in section 2.1. The “sid” column is a survey identifier, which takes on the integers 3,5,7 and 11 for SuperCOSMOS, PanSTARRS, SDSS and ZTF respectively. These integers are intentionally prime, for reasons which become clear later. The magnitudes are the transformed magnitudes, in the PanSTARRS system, as explained in section 4. The “mjd” column is the Modified Julian Date of the measurement concerned. The column mjd\_rf is the observation date in the rest-frame of the quasar, i.e.  $\text{mjd}/(1+z)$  where  $z$  is the redshift of the quasar.

Each quasar typically has something like  $\sim 200$  epochs in a given filter, but there is wide variety. Fig 13 shows the distribution of the number of observations per object. Table shows the total number of unique objects and number of observations in each survey and band, after removing outliers and objects outside the magnitude thresholds. Finally, Fig 14 displays a four-way Venn diagram to illustrate the number of 7-DQ quasars with at least one observation in each combination of surveys

**DEFINITELY A GOOD IDEA.** Maybe consider showing one or more example light curves? Both the whole thing, and the PS and ZTF bits zoomed in?

uid	band	mag	magerr	mjd	mjd_rf	sid
1	g	21.917612	0.067707	54741.371094	16543.176517	5
1	i	21.933083	0.148910	54741.371094	16543.176517	5
1	g	22.171011	0.186231	55477.343750	16765.591946	7
1	i	21.863562	0.171900	55484.289062	16767.690862	7
1	r	21.584364	0.179572	55806.593750	16865.093306	7
1	g	21.753338	0.153714	55879.226562	16887.043385	7
1	i	21.352200	0.181100	56231.410156	16993.475417	7
1	r	21.763350	0.176923	56247.320312	16998.283564	7
1	g	20.977379	0.094571	56549.582031	17089.628900	7
...	...	...	...	...	...	...

**Table 5.** An example of quasar photometry from the 7-DQ database. Here, uid refers to my unique quasar identifier, mjd and mjd\_rf are the time of observation (in MJD), except the latter is in the rest-frame, sid is the survey ID.

Survey	Band	Quasars		Stars	
		$N_{\text{obs}}$	$N_{\text{uniq}}$	$N_{\text{obs}}$	$N_{\text{uniq}}$
SDSS	<i>g</i>	1,403,558	524,966	6,925,451	399,806
	<i>r</i>	1,403,448	525,311	6,931,668	399,819
	<i>i</i>	1,383,692	522,078	6,933,457	399,819
PS	<i>g</i>	2,587,265	500,380	2,036,345	382,275
	<i>r</i>	3,023,809	504,321	2,830,731	399,223
	<i>i</i>	3,620,154	505,146	3,503,282	399,324
ZTF	<i>g</i>	62,635,876	263,221	11,573,951	71,071
	<i>r</i>	75,540,780	269,808	28,076,200	132,170
	<i>i</i>	11,203,322	144,125	4,853,779	101,747
SSS	<i>g</i>	464,138	324,465	112,231	95,411
	<i>r</i>	476,091	213,150	208,257	106,878
	<i>i</i>	67,467	49,703	44,375	37,792
Total	<i>g</i>	67,090,837	525,507	20,647,978	399,914
	<i>r</i>	80,444,128	525,965	38,046,856	399,926
	<i>i</i>	16,274,635	524,512	15,334,893	399,927

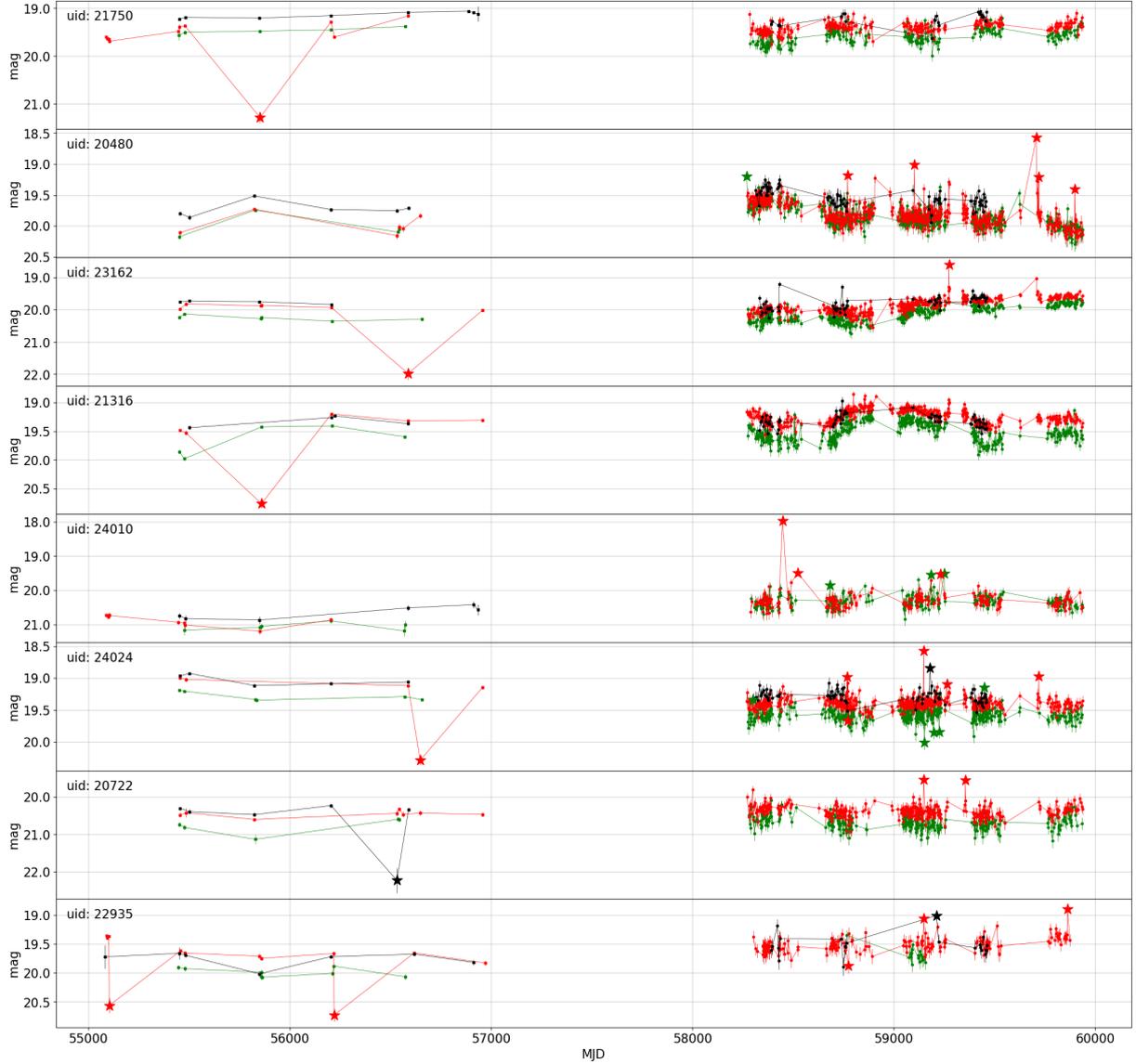
**Table 6.** Cumulative counts of the number of observations,  $N_{\text{obs}}$ , and the number of unique objects,  $N_{\text{uniq}}$ , for 7-DQ quasars and stars in the *g*, *r* and *i* bands after removing outliers and objects outside the magnitude thresholds

## 6 PAIRWISE DATASET AND $\Delta m - \Delta t$ DISTRIBUTIONS

Although the behaviour of individual quasars as a function of absolute time has its own interest, many scientific analyses involve the distribution of magnitude differences of the population as a whole, or or sub-populations grouped by luminosity and other properties, and especially how such distributions depend on the *lag time*. We will look at examples in the following sections.

### 6.1 Construction of pairwise dataset

To speed up the computation of a variety of such analyses, we pre-computed  $\Delta m$  and  $\Delta t$  (rest-frame) for all possible pairs of observations for a given quasar in a given band. For a light curve of length  $N$ , there are  $N(N+1)/2$  unique unordered  $(\Delta m, \Delta t)$  pairs. **INCLUDE THESE TABLES FROM THESIS - thesis tables 3.3, 3.4, 3.5?.** An example of the data from the resulting pairwise dataset is shown in Table 7, and Tables 8 and 9 show statistics of the pairwise dataset.



**Figure 12.** Demonstration of our outlier detection algorithm. Outlier points are marked with a ‘star’, and will be removed from the light curve. Here, we chose to show only parts of the light curve covering photometry from Pan-STARRS and ZTF for clarity, but the same algorithm applies to the entire light curve which includes all four surveys.

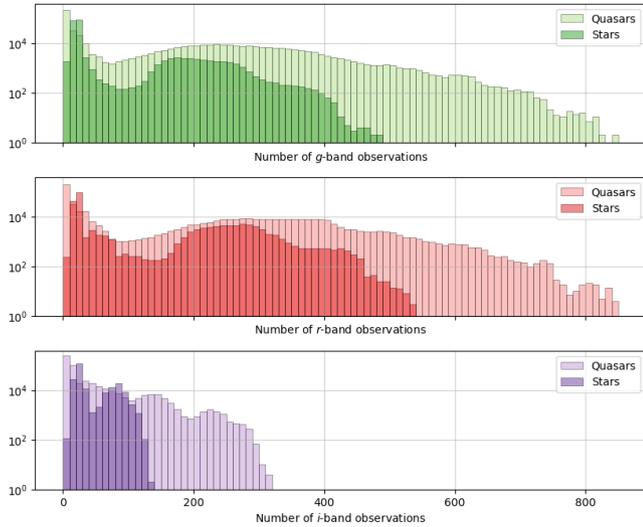
## 6.2 Full $\Delta m$ distributions versus lag

Fig 17 shows the distribution of  $\Delta r$  values in a set of fourteen of rest-frame  $\Delta t$  lag ranges. The  $\Delta t$  bin edges are manually chosen, but spaced logarithmically. We found that 14 bins was an optimal number resulting in gradual, but visible changes between the resulting distributions. The same set of  $\Delta t$  bins is used in later analyses. Comparison of the star and quasar distributions in Figure 17 illustrates clearly that the 7-DQ quasars are variable on timescales  $\Delta t > 20$  days. Additionally, the amplitude of variability (width of the distribution) increases with  $\Delta t$ , as expected. For small time-lags, ( $\Delta t < 20$  days), the quasar and star distributions are very similar, as observed variability is dominated by photometric noise.

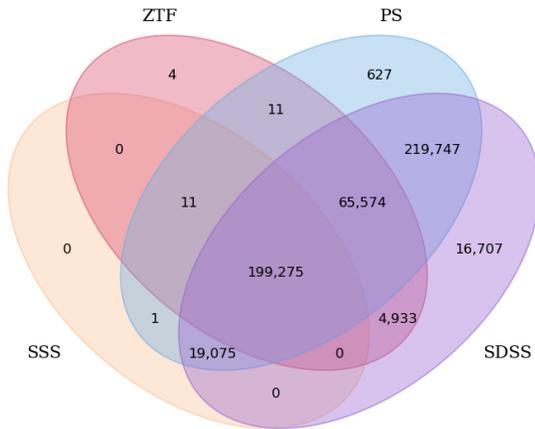
## 6.3 The star sample $\Delta m$ distribution

The  $\Delta m$  distributions for the stars do not look Gaussian, and are not quite the same for all  $\Delta t$ . As the stars should not be variable, the error on any given  $\Delta m$  would be given by the photometric errors of the individual measurements, added in quadrature. However, some points have larger photometric errors than others. The distributions we find are therefore a sum-of-gaussians, which tends to produce the long-tailed distributions that we see. The most important lesson is that when constructing statistics such as the structure function, it is important to take the errors of individual points carefully into consideration - some points have considerably safer variability information than others.

The second important issue is that as well as forming differences within specific surveys, we are also using differences for the same object in different surveys. We have transformed

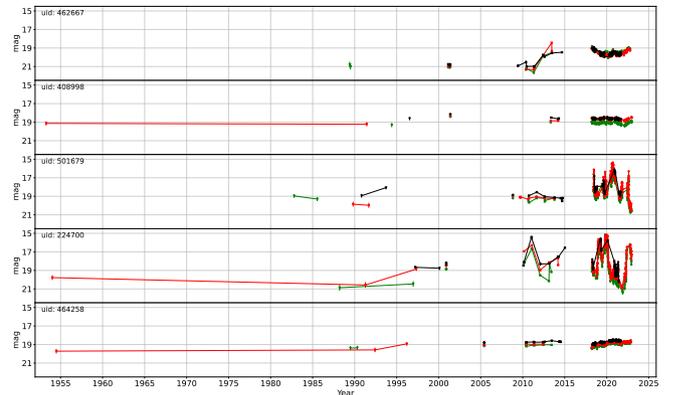


**Figure 13.** Number of observations per object for 7-DQ quasars and stars. There are fewer observations in the *i*-band because only 1/3 of the ZTF *i*-band is public.

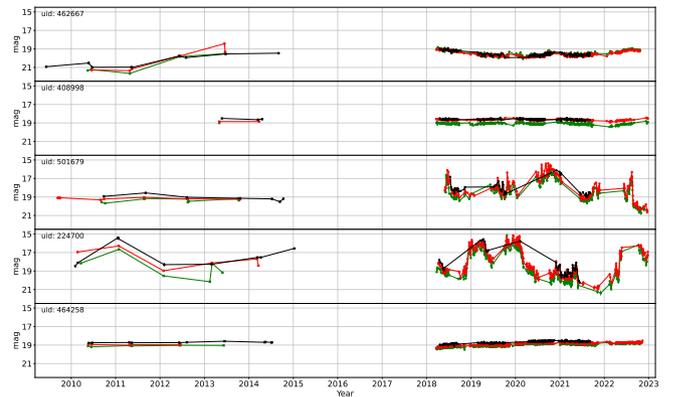


**Figure 14.** 4-way Venn diagram showing the number of 7-DQ quasars which contain at least one observation in each combination of surveys. The majority of quasars are present in either all four surveys, or just Pan-STARRS and SDSS

to a common system as accurately as we can, but systematic errors will remain. We can see that as well the distributions being non-Gaussian, there is some skewness and peak-offset. Therefore, we have annotated the fraction of pairs from each survey combination on each panel of Figure 17. These problems are clearest in the very last panel, which is dominated by pairs combining ZTF and SuperCOSMOS. As well as being the hardest pairing to minimise systematics, the photometric errors on the SuperCOSMOS magnitudes are considerable.



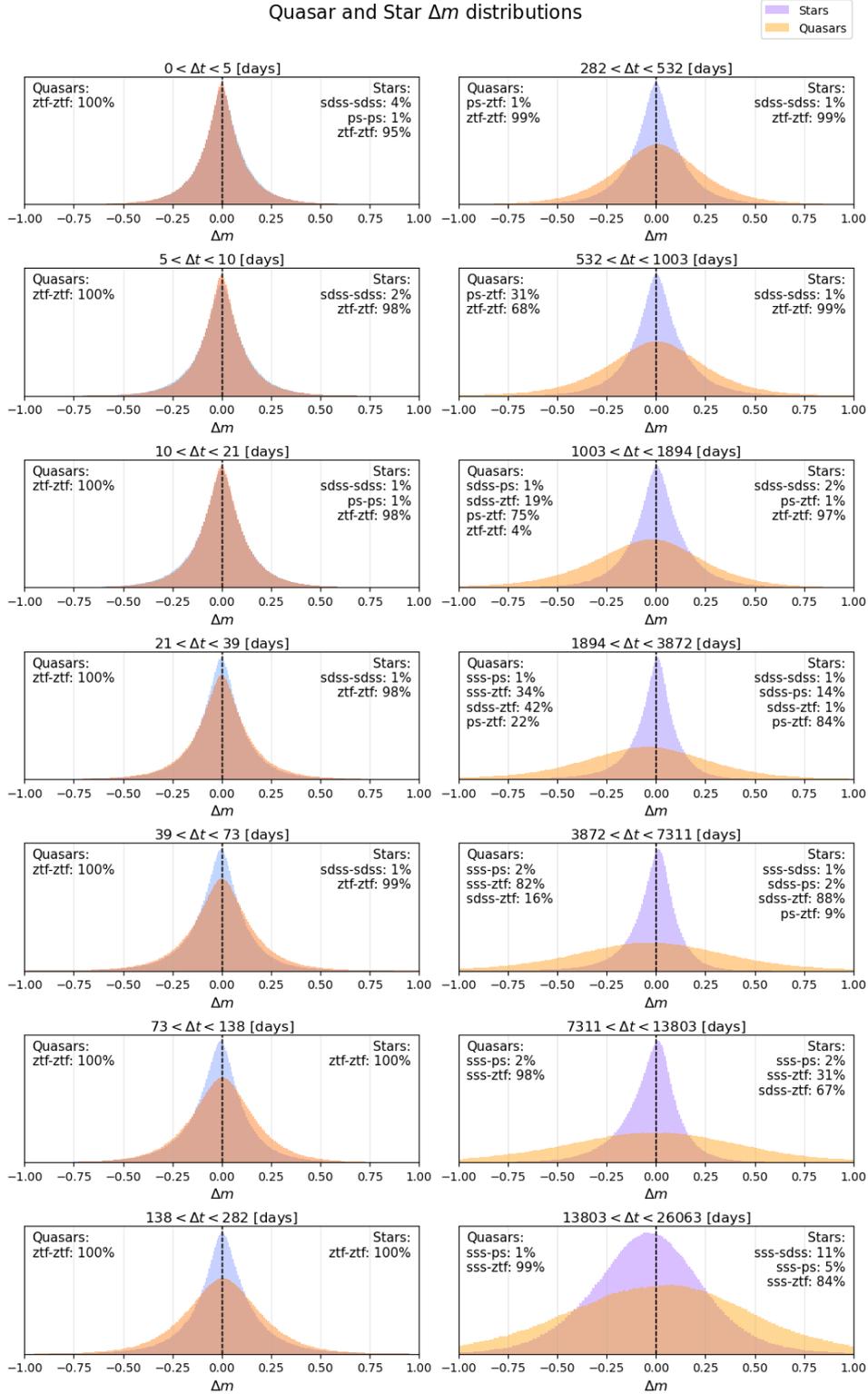
**Figure 15.** Example light curves of four quasars in our sample



**Figure 16.** Example light curves of four quasars in our sample (zoomed into PS/ZTF)

uid	dm	dt	de	dsid
1	0.178986	133.19026	0.252087	49
3	-0.025074	112.04023	0.135006	49
3	0.002586	209.57025	0.101637	49
3	-0.250851	412.65503	0.100035	49
3	0.027660	97.53003	0.133875	49
3	-0.225777	300.61480	0.132663	49
3	-0.253437	203.08480	0.098503	49
4	0.169956	2511.08180	0.239570	21
4	0.188280	2511.46680	0.238982	21
4	0.134626	2517.16720	0.255475	21
...	...	...	...	...

**Table 7.** An example of the data from the pairwise dataset for *r*-band observations of the 7-DQ quasars. Each row represents a unique pair. Here, *uid* is my unique quasar identifier to specify which quasar the pair came from, *dm* is the magnitude difference (equivalent to  $\Delta m$ ) *dt* is the time difference (equivalent to  $\Delta t$ ). For quasars, this is transformed to the rest-frame, whereas for stars, it is left in the observer-frame. *de* is the root-sum-square of the photometric errors, and *dsid* is the product of the survey IDs.



**Figure 17.**  $\Delta m$  distributions grouped by increasing time-lags (top to bottom, left to right) for the ensemble quasar (orange gradient) and star (blue gradient) population. The histograms are coloured using a gradient that changes progressively with  $\Delta t$ , to aid the eye. Each panel is annotated with the percentage of pairs from a specific survey combination compared to the total number of pairs. The two distributions have been normalised by their integrals so that their shapes may be directly compared. Note that, because the quasar  $\Delta t$  have been transformed to the rest frame, the survey combination fractions differ to the stars.

	$N^\circ$ pairs in $g$ -band	$N^\circ$ pairs in $r$ -band	$N^\circ$ pairs in $i$ -band
Quasars	4,959,745,834	7,001,145,864	362,894,257
Stars	835,486,406	2,416,938,570	231,827,462

**Table 8.** Total count for number of pairs calculated from photometry of quasars and stars in the 7-DQ database. **This table is too big!**

Survey combination	Number of pairs		
	$g$ -band	$r$ -band	$i$ -band
SSS-SSS	84,601	174,218	5,798
SSS-SDSS	405,253	397,488	32,983
SSS-PS	1,344,862	1,545,208	154,336
SSS-ZTF	42,099,734	63,621,981	1,365,843
SDSS-SDSS	2,705,746	2,848,155	2,644,237
SDSS-PS	2,709,202	3,471,226	3,787,033
SDSS-ZTF	48,529,353	61,150,318	6,613,372
PS-PS	3,424,657	4,854,153	6,913,230
PS-ZTF	168,019,874	241,167,030	39,727,966
ZTF-ZTF	4,690,422,446	6,621,915,981	301,649,353

**Table 9.** Number of pairs in each band for the quasar photometry, split into survey combinations.

## 7 INITIAL SCIENCE: MEAN DRIFT

In this section and the following two sections, we undertake some preliminary scientific analyses, which are of independent value, but which also demonstrate the power of the database. We begin by looking at the mean of the  $\Delta m$  distributions, and whether they change with time.

Some studies of AGN activity suggest large changes over Myr scales. Evidence for this can be seen in the large-scale extended structures surrounding AGN, such as the so-called ‘‘Voorwerp’’ objects (see e.g., [Sartori et al. 2016](#) and references therein), which clearly show that some AGN exhibit order-of-magnitude changes in their luminosity over  $10^4$ – $10^5$  year timescales. However, there has been comparatively little research to investigate average luminosity changes in large samples of quasars, over timescales for which we have continuous monitoring. Quasar variability is often assumed to be a weakly stationary process, i.e., that the mean luminosity of large ensembles of AGN do not change with time. A few initial studies looked at AGN that were brightening or dimming during the period of observation, but found little or no evidence for differences in their variability properties (see e.g., [de Vries et al. 2003](#), [de Vries et al. 2005](#), [Bauer et al. 2009](#), [Voevodkin 2011](#)). Studies that have observed bulk differences in mean brightness over a sample of quasars often attribute the result to bias; [MacLeod et al. \(2012\)](#) combined data from the Palomar Observatory Sky Surveys (POSS) and SDSS data and noted that objects from POSS are dimmer when observed in SDSS. They concluded that this may be explained by a Malmquist bias, i.e., the fact that a flux limited sample of variable objects will necessarily be dimmer in the later survey even if there is no change in the mean brightness of the underlying sample. [Rumbaugh et al. \(2018\)](#) came to a similar conclusion when studying a sample of extreme variability quasars (EVQs). [Morganson et al. \(2014\)](#) found decrease in mean brightness on decade timescales when comparing SDSS and Pan-STARRS data, but attributed this to filter differences.

One of the most precise studies of mean brightness changes to-date was conducted by [Caplar et al. \(2020\)](#). By comparing observations of  $\sim 6000$  quasars from SDSS and Hyper Suprime-Cam (HSC), they claimed a characteristic dimming of  $\sim 0.2$  mag over a timescale of 10 years in the rest-frame. However, [Shen & Burke \(2021\)](#) claim that the result obtained by [Caplar et al. \(2020\)](#) is biased due to using a flux limited sample of objects. Clearly, this topic is contentious and there is no clear consensus on whether the luminosity variations of AGN are stationary or not.

If quasar variability is not a weakly stationary process, as is often assumed, we would expect the ensemble mean magnitude to drift with time as [Caplar et al. \(2020\)](#) observed. Using 7-DQ, we have tested this assumption and measured the mean magnitude drift over the ensemble population by calculating the first moment of the  $\Delta m$  distributions. If quasar variability is non-stationary, it would manifest in a significant non-zero mean of our  $\Delta m$  distributions, provided that the result cannot be attributed to bias.

We calculated the mean using two methods: First, we used all  $\Delta m$  pairs from every combination of surveys. We refer to the mean using this method as  $\mu_{\text{all}}$ . Second, we used  $\Delta m$  pairs within surveys only, such that we can avoid biases introduced when comparing magnitudes between different photometric systems. We refer to the mean using this method as  $\mu_{\text{inner}}$ . Note that there are two exceptions for the inner pairs: First, we included observation pairs between SDSS and Pan-STARRS because their photometric systems are very similar. Second, we omitted pairs within SuperCOSMOS. This was done because the photometry of SuperCOSMOS does not use the same photometric system; each survey within SuperCOSMOS has its own filter profile and plate emulsion combination. Furthermore, photographic plate data is inherently more noisy than the other surveys and does not have the required precision to effectively constrain the mean. The survey combinations used for the inner pairs are therefore:  $\Delta m_{\text{SDSS-SDSS}}$ ,  $\Delta m_{\text{PS-PS}}$ ,  $\Delta m_{\text{SDSS-PS}}$  and  $\Delta m_{\text{ZTF-ZTF}}$ .

In order to obtain a precise measurement of the mean, we use an optimal weighted average of  $\Delta m$  values,

$$\hat{\mu}(\Delta t) = \frac{\sum_k \Delta m(\Delta t)_k \cdot w_k}{\sum_k w_k} \quad (13)$$

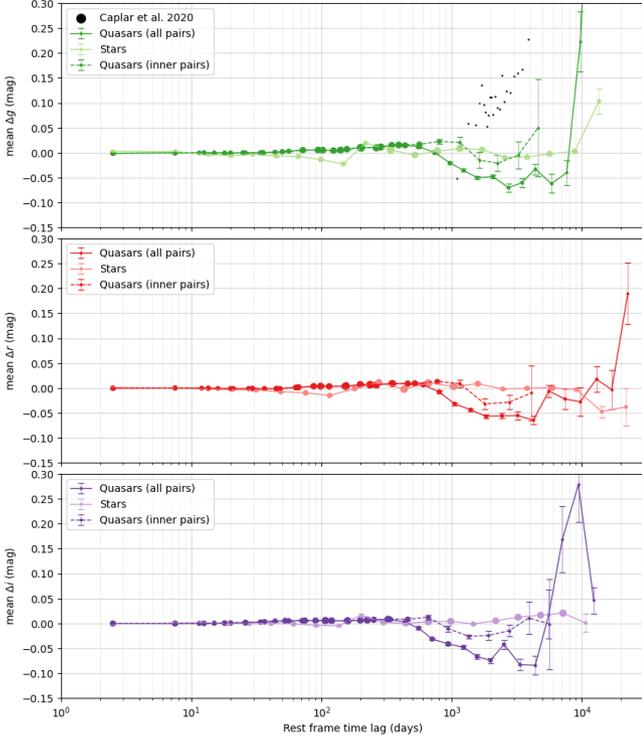
with weights defined as the variance of pair  $k$ , defined as sum of photometric variances from measurements  $i$  and  $j$ ,

$$w_k = \sigma_k^{-2} = (\sigma_i^2 + \sigma_j^2)^{-2}. \quad (14)$$

The optimal error on  $\hat{\mu}$  is then,

$$\sigma_{\hat{\mu}} = \left( \sum_k w_k \right)^{-1/2}. \quad (15)$$

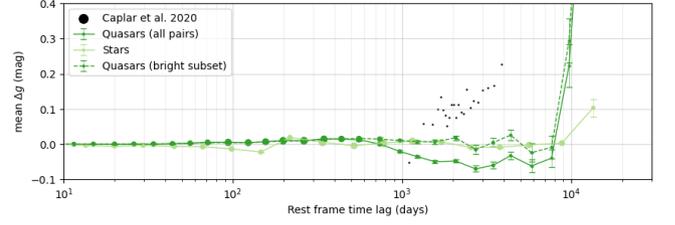
This process was repeated in the  $g$ ,  $r$  and  $i$  bands, and the results are shown in Figure 18. The values of  $\mu_{\text{all}}$  and  $\mu_{\text{inner}}$  for the stars are both consistent with zero on all timescales (with the exception of a few noisy points), which is reassuring, and confirms the effectiveness of our colour transformations. Note that  $\mu_{\text{inner}}$  for the stars is not shown in Figure 18 as it is visually identical to  $\mu_{\text{all}}$ , except that it does not extend beyond  $\Delta t > 5 \times 10^3$  days. For  $\Delta t < 5 \times 10^2$  days,  $\mu_{\text{all}}$  and  $\mu_{\text{inner}}$  for the quasars are both consistent with zero, however, on longer time-lags,  $\mu_{\text{all}}$  starts to deviate from zero, then fluctuate erratically at the longest timescales of  $10^4$  days. Conversely,



**Figure 18.** Mean magnitude drift for  $g$ ,  $r$  and  $i$  bands for 7-DQ quasars (dark line) and stars (light line). The mean calculated using inner pairs for the quasars is also shown (dotted line). Data from Caplar et al. (2020) is overplotted on the top panel (black points). Note that the size of the circular markers illustrates the relative number of points in that bin. Note that  $\mu_{\text{inner}}$  is not shown for the stars to prevent overcrowding. Negative magnitude change corresponds to brightening.

$\mu_{\text{inner}}$  does not show these large variations on long timescales. Since  $\mu_{\text{all}}$  and  $\mu_{\text{inner}}$  only differ in the pairs used during calculation, differences between them can only be caused by bias introduced when comparing magnitudes from different surveys. Since the result using the star sample demonstrate the validity of the colour transformations, the cause of this bias is likely dominated by differing survey depths.

The behaviour of our result  $\mu_{\text{all}}$  (for the quasars), along with the result obtained by Caplar et al. (2020) (overplotted on Figure 18), can be attributed to biases described by Shen & Burke (2021). Caplar et al. (2020) compared quasar photometry from SDSS to Hyper Suprime-Cam (HSC), with HSC having a fainter limiting magnitude, than SDSS. Their quasar sample was defined and observed using SDSS, with HSC data of the same objects being taken at a later date. Shen & Burke (2021) claim that, by using a deeper survey at a later date, the mean of a flux-limited sample is biased such that dimming is more likely to be seen. My result, using  $\mu_{\text{all}}$ , behaves in the opposite sense. Every survey combination (ignoring inner pairs) involves a shallower survey following a deeper survey, with the exception of SuperCOSMOS. For example, SDSS-PS, SDSS-ZTF, and PS-ZTF comparisons are all deep-to-shallow comparisons, which explains the apparent brightening implied by negative  $\mu_{\text{all}}$  on timescales of  $3 \times 10^3$  days. However, pair combinations using SuperCOSMOS photometry involve shallow-to-deep comparisons and



**Figure 19.** Comparing the mean magnitude drift calculated for all quasars, and the bright subset of quasars.

explains the sudden apparent dimming on the longest timescale (final data points in each  $g$ ,  $r$  and  $i$  bands).

As a second test, to demonstrate that these effects are indeed caused by bias, we re-evaluated  $\mu_{\text{all}}$  using the bright subsample of quasars. The result is shown in Figure 19, with  $\mu_{\text{all}}$  overplotted for comparison. We have only plotted the  $g$ -band for brevity, but the same effect applies in all  $g$ ,  $r$  and  $i$  bands.

$\mu_{\text{inner}}$  shows no significant change in mean magnitude on timescales  $\Delta t < 15$  years. Since this is our most precise measurement of the mean, we conclude that there is no bulk change in quasar magnitude up to these timescales, and that the results of Caplar et al. (2020) are due to bias, consistent with Shen & Burke (2021). We are unable to constrain the mean effectively beyond these timescales, as  $\mu_{\text{all}}$  also suffers from the same bias.

## 8 INITIAL SCIENCE: HIGHER ORDER MOMENTS VERSUS LAG

After the mean, the second moment is of course the variance. Behaviour of this versus lag is the essence of the structure function, which we discuss in the following section. First however, we look at the third and fourth moments - skewness and kurtosis - which we believe have not previously been examined.

### 8.1 Skewness

The skewness,  $S$  is the third moment of the  $\Delta m$  distribution and represents the overall asymmetry of the distribution. It is calculated via

$$S = \frac{\mu_3}{\mu_2^{3/2}}, \quad (16)$$

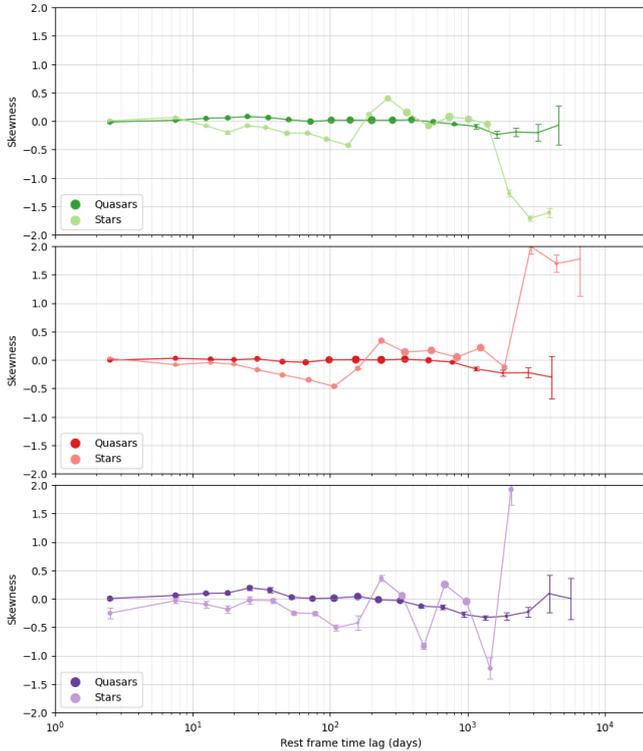
where  $\mu_n$  is the  $n^{\text{th}}$  central moment of the distribution,

$$\mu_i = \frac{1}{N} \sum_{n=1}^N (x[n] - \bar{x})^i. \quad (17)$$

The standard error on the skewness depends only on the sample size:

$$\sigma_{\text{skew}} = \sqrt{\frac{6n(n-1)}{(n-2)(n+1)(n+3)}}, \quad (18)$$

and is used for the error bars in Figure 22. For a Gaussian,  $S = 0$ . Distributions with a positive skew,  $S > 0$ , implies that the  $\Delta m > 0$  tail is longer and therefore a higher probability of large dimming changes, while  $S < 0$  implies that the



**Figure 20.** Skewness of  $\Delta m$  of 7-DQ quasars and stars in the  $g$ ,  $r$  and  $i$  bands. Note that the size of the circular markers illustrates the relative number of points in that bin.

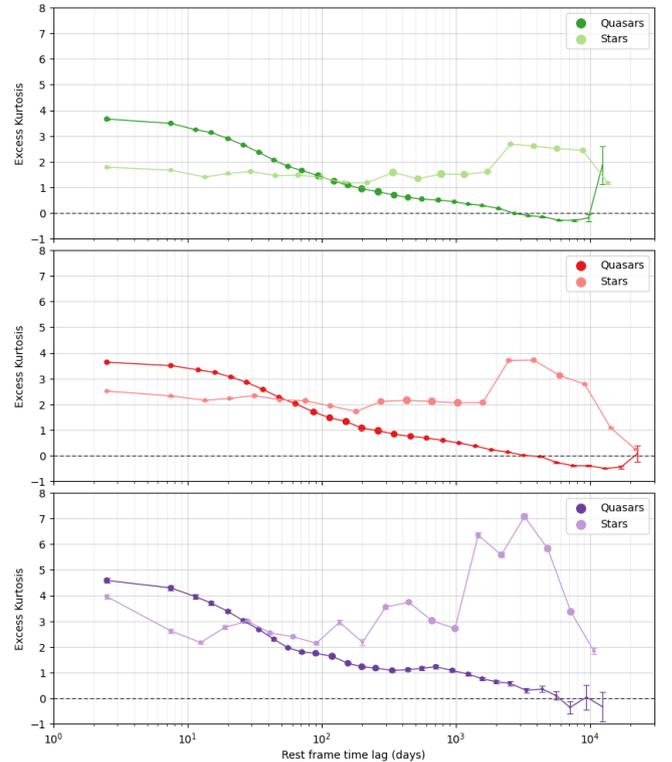
$\Delta m < 0$  tail is longer and therefore a higher probability of large brightening changes. However, because the skewness is a dimensionless parameter, it is hard to physically interpret the exact value of  $S$  over a given timescale.

We found that skewness was particularly sensitive to differences in photometric systems of the surveys, evidenced by erratic fluctuations in the skewness of the star sample, when using the full set of pairs. Therefore, as with the mean, we decided to compute the skewness using restricted pairs ( $\Delta m_{\text{SDSS-SDSS}}$ ,  $\Delta m_{\text{PS-PS}}$ ,  $\Delta m_{\text{SDSS-PS}}$  and  $\Delta m_{\text{ZTF-ZTF}}$ ), which offers the most precise measurement of  $S$ . In Figure 20, we present the skewness of the quasar and star  $\Delta m$  distributions, in the  $g$ ,  $r$  and  $i$  bands. This is a novel result, however, since this statistic has not yet been studied for quasar variability we present it tentatively as a reference that can be used when comparing results from different photometric studies or simulations. We see that the skewness of the 7-DQ quasars and stars are comparable for time-scales up to 1 year. This suggests fluctuations in brightness for quasar variability are equally probably for dimming and brightening over this timescale. For time-scales greater than 1 year, the data is not conclusive, and it is likely that inaccuracies of the colour transformations cause systematic errors.

## 8.2 Kurtosis

The kurtosis,  $K$ , is the fourth moment of the  $\Delta m$  distribution, defined as:

$$K = \frac{\mu_4}{\mu^2}, \quad (19)$$



**Figure 21.** Kurtosis of  $\Delta m$  of 7-DQ quasars and stars in the  $g$ ,  $r$  and  $i$  bands. Note that the size of the circular markers illustrates the relative number of points in that bin.

defined again in terms of the  $n^{\text{th}}$  central moment,

$$m_i = \frac{1}{N} \sum_{n=1}^N (x[n] - \bar{x})^i. \quad (20)$$

$K$  quantifies the presence of tails in a distribution. For a Gaussian distribution, we expect  $K = 3$ . Excess kurtosis is defined as  $K_{\text{ex}} = K - 3$ , which is used to show deviation from a normal process. Therefore, we expect  $K_{\text{ex}} = 0$  for a Gaussian process, whereas a distribution with heavy tails and many outliers will have  $K_{\text{ex}} \gg 0$ . As with skewness, the standard error on kurtosis depends only on sample size:

$$\sigma_{\text{kurt}} = \sqrt{\frac{6n(n-1)}{(n-2)(n+1)(n+3)}}, \quad (21)$$

We show the excess kurtosis for the 7-DQ quasars and stars in Figure 21. For the stars,  $K_{\text{ex}}$  seems to fluctuate around the average with no clear trend with  $\Delta t$ , which is reassuring. The average  $K_{\text{ex}}$  is significantly non-zero in all three bands, which implies that the combination of noise distributions results in a distribution with heavier tails than that of a Gaussian, on all timescales.

The excess kurtosis of the quasars shows a gradual decline with  $\Delta t$ , ranging from  $K_{\text{ex}} \approx 4$  on the shortest timescales to  $K_{\text{ex}} \approx 0$  on the longest timescales of  $\Delta t = 10^4$  days. Again, this is a dimensionless parameter and therefore the exact value of kurtosis is difficult to interpret physically.

## 9 INITIAL SCIENCE: ENSEMBLE STRUCTURE FUNCTION

The structure function - basically the variance as a function of lag - has been a popular way of statistically characterising quasar variability for many years. Here we take a first look at the structure function (SF) for the quasars in our sample. We will examine the SF more fully, and especially how it depends on quasar properties, in Paper II.

In this first look, we present only the  $r$ -band SF, and consider the *ensemble* structure function for all quasars combined, i.e. on the assumption that the variability behaviour of all quasars is drawn from a single universal SF. In Paper-II we divide the sample into sub-ensembles, based on luminosity, mass, and Eddington ratio.

### 9.1 Calculation of Structure Function

Like many other natural processes, measurements of quasar brightness taken a short time apart are typically very similar, whereas measurements taken with a longer time gap show a larger spread. The most popular method for quantifying this behaviour of quasar variability is the structure function - essentially the RMS of magnitude changes as a function of lag. It is a powerful tool which is robust against sparse sampling and may be applied to either a single object or an ensemble. It was first used by astronomers by [Simonetti et al. \(1985\)](#).

For a given quasar we can find all the magnitude differences  $\Delta m = m(t_j) - m(t_i)$  evaluated at (rest frame adjusted) times  $t_j$  and  $t_i$ , separated by a time-lag  $\Delta t = t_j - t_i$ . We can group these values in a range of  $\Delta t$  and take the RMS of  $\Delta m$  values as our estimate of  $SF(\Delta t)$  at the centre of the range. However the observed magnitude differences come from both intrinsic variability and photometric errors for each of the two points,  $\sigma_i$  and  $\sigma_j$ . We estimate the intrinsic  $SF_{int}$  by subtracting the photometric errors in quadrature.

However, if we group many quasars, assuming they have identical intrinsic structure functions, to form an ‘ensemble structure function’, this has the advantage of both improving signal noise and of sampling many different rest frame  $\Delta t$  values, given the range of redshifts. The ensemble SF is then calculated as follows:

$$SF_{int,ensemble}(\Delta t) = \sqrt{\frac{1}{N(\Delta t)} \sum_k \sum_{j < i} (\Delta m_{ij,k}^2 - \sigma_{i,k}^2 - \sigma_{j,k}^2)}, \quad (22)$$

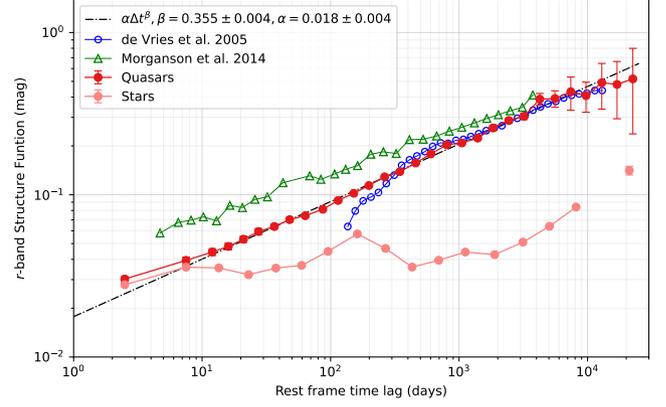
where  $k$  denotes the index of the quasar.

Sometimes the variance subtracted magnitude differences will be less than zero, and occasionally even the sum inside the square root in the above equation will be negative. We avoid this problem by using a variance weighted estimate, which can also be shown to be the optimal average. We therefore calculate

$$\widehat{SF}(\Delta t) = \sqrt{\frac{\sum_k SF_k^2(\Delta t) \cdot w_k}{\sum_k w_k}}, \quad (23)$$

where the weights,  $w_k$ , are defined as the inverse-variance of  $SF^2$

$$w_k = \text{Var}[SF^2(\Delta t)]^{-1} = \frac{1}{2}(\sigma_i^2 + \sigma_j^2)^{-2}. \quad (24)$$



**Figure 22.** This is a cut-out from previous figure - needs remaking properly. - DONE Ensemble structure function for the quasars and stars in the  $r$  band. The size of the points represents the relative number of points in the bin. The black dot-dashed line represents an SPL fit to the quasar structure function data points for  $\Delta t > 10$  days, with the slope shown in the legend. Comparison data from [de Vries et al. \(2005\)](#), and [Morganson et al. \(2014\)](#) are overplotted. A few points are omitted for the ensemble star structure function, as the photometric errors are greater than the observed variability.

In paper II we show a more complete derivation of this method, and show how it produces a smoother and better behaved estimate of the SF.

### 9.2 Ensemble Structure Function result

In Figure 22, we present the  $r$  band ensemble structure function of our 7-DQ quasar and star populations, stretching from a few days to many decades.

The structure function of the stars is flat, fluctuating around a median value of  $SF = 0.03$  mag. This can be interpreted as residual uncorrected variability due to photometric noise. For the quasars, SF values below this level can therefore be considered unreliable. This is only the case for the quasar structure function on timescales  $\Delta t < 10$  days, suggesting that there is minimal intrinsic variability exhibited by the ensemble quasar population on these timescales.

On timescales  $\Delta t > 10$  days, the quasar ensemble structure function appears to follow a single power law (SPL) up to the longest timescales of  $\Delta t \sim 10^4$  days. We quantified this by fitting the form  $\alpha\Delta t^\beta$  over these timescales, giving slope  $\beta_r = 0.355 \pm 0.004$  and normalisation  $\alpha_r = 0.018 \pm 0.004$ . This corresponds to the SF amplitude at 1 day. More helpfully, this corresponds to SF amplitude of 0.092 and 0.473 magnitudes at 100 days and 10,000 days respectively. The SF continues to follow an SPL (within error) even on the longest timescales, indicated by the right-most points. There is no sign of a long timescale turnover.

Ensemble structure functions calculated in the  $r$  band by [de Vries et al. \(2005\)](#) and [Morganson et al. \(2014\)](#) are overplotted. Our SF agrees well with that of [de Vries et al. \(2005\)](#) over the range of time lags in common, except that their SF dips at their shortest timescales. Given our improved calculation method, and the smoothness of our result, we are confident that our result is more reliable. Another comparable study is that of [Morganson et al. \(2014\)](#) who used a sample of 105,783

quasars observed over 10 years. They claim an SPL model fits the data well with a slope  $\beta = 0.2457 \pm 0.0025$ . This is qualitatively in agreement with our result in that a smooth power law is found over a large time range, but the slopes we find are not in agreement. We agree on the longest timescales, but the discrepancy becomes increasingly larger at shorter timescales. This suggests that the discrepancy might result from a different treatment of the photometric errors.

### 9.3 Discussion

Following Kelly et al. (2009), a popular assumption is that the variability of quasars is likely to be described by a Damped Random Walk (DRW), which is also equivalent to a first order autoregressive process, or to a linear system with random seed and exponential filter (see for example Lawrence (2012)). Indeed some studies have argued that the quasar ensemble SF does follow the expected DRW form, (e.g. MacLeod et al. (2012)), which is

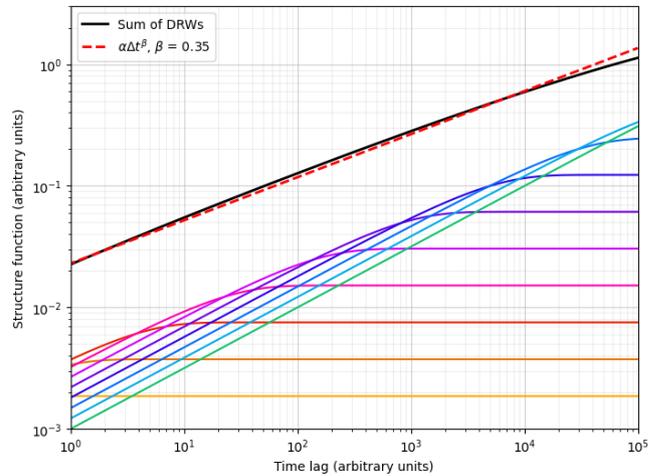
$$SF(\Delta t) = \sqrt{2}\sigma \left(1 - e^{-\Delta t/\tau}\right)^{1/2},$$

where  $\sigma$  is the overall variance of the process, and  $\tau$  is the damping timescale. On long timescales, the SF is flat, and on short timescales  $SF \propto \Delta t^{0.5}$

The observed SF shows a smooth power law over many decades with slope  $\beta_r = 0.355 \pm 0.004$ , strongly inconsistent with the DRW expectation of  $\beta = 0.5$ . Furthermore we see no sign of a long timescale turnover. If quasar variability is a stationary random process, there must eventually be a long timescale turnover, or the process will have infinite variance. Of course, it may be that the process is *not* stationary. Assuming that it is stationary, perhaps the turnover is just beyond the timescales we have sampled, and the reason for the flat slope is that we are measuring within the turnover region. However, this is not plausible, because the gradient of the DRW form is only consistent with our observed value of  $0.35 \pm 0.004$  over a very small range, whereas we see a consistent slope over many decades of timescale.

The ensemble SF assumes that every quasar is the same, at least in some statistical sense. It may be that individual quasars do have characteristic timescales, but there exists a wide range of these timescales, with the resultant SF being the envelope of the combined SFs. A little numerical experimentation shows that this is indeed a possible explanation. Fig. 23 illustrates how this could be done, combining 10 DRWs with logarithmically spaced timescales. Note that to mimic the 0.35 slope, there needs to be a range of amplitudes as well as a range of timescales, and the objects with shorter timescales need to have larger overall variability amplitude. If this the origin of the 0.35 slope, the implication is that around 30% of quasars have characteristic timescales of less than 5 days. Overall, this does not seem a reasonable explanation of the flat SF slope. However, we will explore this option further in Paper II.

If there is a distribution of amplitudes and/or timescales, it would seem reasonable to expect that these depend on the luminosity of the quasar, or the mass of the central black hole, or both. In paper II we divide our sample into sub-ensembles based on luminosity, estimated black hole mass, and Eddington ratio. We will report on the results in Paper II, but the brief answer is that while the shape of the SF



**Figure 23.** An analytical demonstration showing that a combination of DRW structure functions can approximate a single power law (SPL) over a range of timescales. Ten DRWs were used with a set of timescales and amplitudes that were logarithmically spaced and proportional to each other. The envelope reasonably reproduces a typical slope of  $\beta \sim 0.35$

does depend on mass and luminosity, the sub-ensembles still do not show a DRW form.

Meanwhile, Burke et al. (2021) have taken a different approach to studying the dependence of variability on black hole mass. Rather than ensembles, they have fitted the DR model to a range of individual quasars, deriving  $\tau$  for each. They find evidence that  $\tau$  depends on the black hole mass. In a later paper, we will attempt to duplicate this analysis with our much larger sample.

## 10 DATA AVAILABILITY

Time to fill this in...

## 11 CONCLUSIONS

To come

## ACKNOWLEDGEMENTS

### References

- Bauer A., Baltay C., Coppi P., Ellman N., Jerke J., Rabinowitz D., Scalzo R., 2009, *ApJ*, **696**, 1241
- Burke C. J., et al., 2021, *Science*, **373**, 789
- Caplar N., Pena T., Johnson S. D., Greene J. E., 2020, *ApJ*, **889**, L29
- Hook I. M., McMahon R. G., Boyle B. J., Irwin M. J., 1994, *MNRAS*, **268**, 305
- Ivezić Ž., et al., 2007, *AJ*, **134**, 973
- Kelly B. C., Bechtold J., Siemiginowska A., 2009, *ApJ*, **698**, 895
- Lawrence A., 2016, Clues to the Structure of AGN Through Massive Variability Surveys. p. 107

- Li Z., McGreer I. D., Wu X.-B., Fan X., Yang Q., 2018, *ApJ*, **861**, 6
- MacLeod C. L., et al., 2012, *ApJ*, **753**, 106
- Masci F. J., et al., 2019, *PASP*, **131**, 018003
- Morganson E., et al., 2014, *ApJ*, **784**, 92
- Pâris I., et al., 2018, *A&A*, **613**, A51
- Peacock J. A., Hambly N. C., Bilicki M., MacGillivray H. T., Miller L., Read M. A., Tritton S. B., 2016, *MNRAS*, **462**, 2085
- Rumbaugh N., et al., 2018, *ApJ*, **854**, 160
- Sartori L. F., et al., 2016, *MNRAS*, **457**, 3629
- Shen Y., Burke C. J., 2021, *ApJ*, **918**, L19
- Simonetti J. H., Cordes J. M., Heeschen D. S., 1985, *ApJ*, **296**, 46
- Tonry J. L., et al., 2012, *ApJ*, **750**, 99
- Voevodkin A., 2011, *arXiv e-prints*, p. arXiv:1107.4244
- Wu Q., Shen Y., 2022, *ApJS*, **263**, 42
- de Vries W. H., Becker R. H., White R. L., 2003, *AJ*, **126**, 1217
- de Vries W. H., Becker R. H., White R. L., Loomis C., 2005, *AJ*, **129**, 615
- No appendices? Only in Paper II