

Photo-z-SQL: integrated, flexible photometric redshift computation in a database

Róbert Beck^{a,b,*}, László Dobos^a, Tamás Budavári^{c,d,b}, Alexander S. Szalay^b, István Csabai^a

^aDepartment of Physics of Complex Systems, Eötvös Loránd University, 1/A. Pázmány Péter sétány, 1117 Budapest, Hungary

^bDepartment of Physics and Astronomy, The Johns Hopkins University, 3400 N Charles Street, Baltimore, MD 21218, USA

^cDepartment of Applied Mathematics and Statistics, The Johns Hopkins University, 3400 N Charles Street, Baltimore, MD 21218, USA

^dDepartment of Computer Science, The Johns Hopkins University, 3400 N Charles Street, Baltimore, MD 21218, USA

Abstract

We present a flexible template-based photometric redshift estimation framework, implemented in C#, that can be seamlessly integrated into a SQL database (or DB) server and executed on-demand in SQL. The DB integration eliminates the need to move large photometric datasets outside a database for redshift estimation, and utilizes the computational capabilities of DB hardware. The code is able to perform both maximum likelihood and Bayesian estimation, and can handle inputs of variable photometric filter sets and corresponding broad-band magnitudes. It is possible to take into account the full covariance matrix between filters, and filter zero points can be empirically calibrated using measurements with given redshifts. The list of spectral templates and the prior can be specified flexibly, and the expensive synthetic magnitude computations are done via lazy evaluation, coupled with a caching of results. Parallel execution is fully supported. For large upcoming photometric surveys such as the LSST, the ability to perform in-place photo-z calculation would be a significant advantage. Also, the efficient handling of variable filter sets is a necessity for heterogeneous databases, for example the Hubble Source Catalog, and for cross-match services such as SkyQuery. We illustrate the performance of our code on two reference photo-z estimation testing datasets, and provide an analysis of execution time and scalability with respect to different configurations. The code is available for download at <https://github.com/beckrob/Photo-z-SQL>.

Keywords: galaxies: distances and redshifts, techniques: photometric, astronomical databases: miscellaneous

1. Introduction

In recent years, photometric redshift (often abbreviated as photo-z) estimation has become a vital and widespread tool in the astronomical repertoire. Large amounts of photometric data are produced by sky surveys such as the Sloan Digital Sky Survey (York et al., 2000; Eisenstein et al., 2011; Alam et al., 2015, SDSS), but spectroscopic measurements are more scarce due to the longer observing times required. Therefore it is important to accurately estimate the redshift – and thus the distance – of objects based just on their photometry.

There are two distinct approaches to photo-z estimation – the empirical approach starts from a training set with known redshifts and uses a machine learning method to perform the estimation (Connolly et al., 1995; Wang et al., 1998; Wadadekar, 2005; Csabai et al., 2007; Carliles et al., 2010; Gerdes et al., 2010; Brescia et al., 2014; Beck et al., 2016), while the template-based approach fits a spectral template using synthetic photometry, i.e. by computing the broad-band magnitudes corresponding to known filter transmission curves and a known spectrum at a given redshift. The former approach is generally more accurate

(Csabai et al., 2003), but the latter does not require an extensive training set, and is thus more flexible.

Within the family of template-based photometric redshift estimation methods, there are two main branches: those that search for the maximum likelihood, best-fitting spectral energy distribution (SED) template and record the redshift of this single SED (Bolzonella et al., 2000; Csabai et al., 2000; Arnouts et al., 2002; Ilbert et al., 2006), and Bayesian methods that aim to reproduce the full posterior redshift probability distribution based on the observations (Benítez, 2000; Coe et al., 2006; Brammer et al., 2008; Budavári, 2009). The latter method provides a more detailed answer, but it is computationally more expensive, and it does need more input information in the form of a prior (or, with a flat prior, the peak of the probability distribution gives the same result as the maximum likelihood method).

Large photometric catalogs are most often stored in relational databases (e.g. *SkyServer*¹ for SDSS), and cross-matching tools that connect such databases also use the relational model (Dobos et al., 2012; Budavari et al., 2013; Han et al., 2016). It would be advantageous to have a photo-z implementation that can work directly on data in

*Corresponding author

Email address: beckrob23@caesar.elte.hu (Róbert Beck)

¹<http://skyserver.sdss.org/>

the servers, and on results from cross-matches, by integrating into the well-established SQL language.

There are various software packages available that utilize either the template-based or the empirical approach to perform photometric redshift estimation². However, to the best of our knowledge, none of the public codes have the ability to perform the computations within the database itself. Thus, they require moving the photometric catalogs outside of the database, and the results also have to be loaded into the DB separately. This is a cumbersome process, especially when the amount of data is in the PB range, as will be the case with upcoming surveys such as the LSST (Ivezic et al., 2008) and Pan-STARRS (Tonry et al., 2012). The problem is also acute when on-demand photo-z would be needed to quickly process cross-match results.

Additionally, existing template-based photometric redshift estimation codes use a predefined set of broad-band photometric filters, and compute synthetic magnitudes in all of these bands before the actual photo-z calculations. This approach is not ideal for heterogeneous observations, such as the Hubble Source Catalog (Whitmore et al., 2016), where each object can have a different corresponding filter set, and the total number of filters is much larger than the number of available filters for an object. Also, this issue can become especially prominent if we consider that filter transmission curves can change over time – if this time-dependence is to be properly taken into account, precomputing synthetic magnitudes is simply not an option.

Even with all these difficulties, the template-based approach is at least feasible in such situations, while building a training set for empirical estimation is almost impossible when the input filter set is not static, and not defined beforehand. Moreover, cross-match results are generally not extensive enough to provide complete training set coverage in redshift and color-magnitude space, and empirical methods usually do not perform well when substantial extrapolation is necessary (Beck et al., 2017).

Thus, we decided to adopt the template-based photo-z approach, and sought to address the issues of existing codes with our own implementation, named Photo-z-SQL. We apply a computational innovation to deal with variable filter sets – we use lazy evaluation for synthetic magnitudes, computing them only when they are needed, and caching them to be reused.

The code has been developed in the C# language, which can be run stand-alone, but can also be integrated into Microsoft SQL Server as a set of user-defined functions. Microsoft SQL Server is a commercial DB server application, which has been chosen because it enables using a general-purpose computer language via CLR integration (see Dobos et al., 2011, for another astronomy-related use case), and also because it is used by e.g. the SDSS and Pan-STARRS collaborations.

²See <http://www.sedfitting.org/> for a collection of links to photo-z codes.

The database integration has the added benefit of utilizing the computational capabilities of DB servers – i.e. bringing the computations directly to the data –, which often have multiple processors specifically to handle computationally intensive tasks. Parallel execution is fully supported by our code, taking advantage of modern multi-core systems. The objective-oriented power of C# also leads to flexibility in defining and providing priors, templates and filters.

2. Photometric redshift estimation methods

Our code implements both the maximum likelihood and the Bayesian template-based photo-z method as separate functions, following the standard approaches in the literature (Ilbert et al., 2006; Coe et al., 2006). This section details the calculations performed by our algorithm.

2.1. Maximum likelihood method

The maximum likelihood method can be summarized with the formulae

$$\chi^2(z', t', m'_0) = \frac{1}{2} \sum_{ij} (m_i - (s_i(z', t') - m'_0)) C_{ij}^{-1} (m_j - (s_j(z', t') - m'_0)) \quad (1)$$

and

$$(z_{\text{phot}}, t, m_0) = \arg \min_{(z', t', m'_0)} \chi^2(z', t', m'_0), \quad (2)$$

where z' and t' are the redshift and the template SED, m'_0 is a factor that scales the total flux of the template spectrum, \mathbf{m} is the measured broad-band magnitude vector of an object, $\mathbf{s}(z', t')$ is the synthetic magnitude vector of template t' at redshift z' , and finally, \mathbf{C} is the covariance matrix of magnitude errors between filters.

The fitting process involves iterating over a set of redshifts and templates, but not over m'_0 , since for a given z' and t' the m'_0 value that yields the lowest χ^2 can be determined algebraically. $(z_{\text{phot}}, t, m_0)$ denote the set of best-fitting parameters, corresponding to the lowest χ^2 value found during the iteration. z_{phot} is the photometric redshift estimate.

Most applications assume that magnitude errors are uncorrelated, in addition to being Gaussian with a standard deviation matching the estimated measurement error for the given object. Under the uncorrelated assumption, the \mathbf{C} covariance matrix takes the form

$$C_{ij} = \sigma_i^2 \delta_{ij}, \quad (3)$$

where σ is the estimated magnitude error vector of the object, and δ is the Kroenecker delta. The expression for χ^2 can be simplified under this assumption to yield

$$\chi^2(z', t', m'_0) = \frac{1}{2} \sum_i \frac{(m_i - (s_i(z', t') - m'_0))^2}{\sigma_i^2}, \quad (4)$$

but in many cases the assumption may simply not hold (Scranton et al., 2005). For this reason, our code implements using the full covariance matrix \mathbf{C} , which, to our knowledge, is a missing feature in other public codes.

2.2. Bayesian method

Using the notation introduced in Sect. 2.1, the Bayesian approach starts from the expression

$$P(z, t, m_0 | \mathbf{m}, \mathbf{C}) = \frac{P(z, t, m_0) P(\mathbf{m}, \mathbf{C} | z, t, m_0)}{\int dz' \int dt \int dm_0 P(z', t, m_0) P(\mathbf{m}, \mathbf{C} | z', t, m_0)}, \quad (5)$$

which specifies the probability $P(z, t, m_0 | \mathbf{m}, \mathbf{C})$ of a redshift and a template SED – scaled to a flux – given the data, i.e. the measured magnitudes and magnitude errors (optionally taking the covariance between filters into account). $P(z, t, m_0)$ is the prior probability of the given template SED at the given redshift. $P(\mathbf{m}, \mathbf{C} | z, t, m_0)$ is the likelihood that the given template and redshift produce the data. The denominator contains a factor that normalizes the total probability to 1, but in practice this is not computed, and instead the final probability distribution is normalized. We note that we use the $\int dt$ term as shorthand for integrating over the model parameters that describe a template SED, t does not represent a real number.

The likelihood term can be formulated as

$$P(\mathbf{m}, \mathbf{C} | z, t, m_0) = \exp(-\chi^2(z, t, m_0)), \quad (6)$$

where χ^2 is the expression defined in Eq. 1. Again, assuming uncorrelated magnitude errors, the simplified version in Eq. 4 can be applied, but our application does support using the full covariance matrix \mathbf{C} .

The method aims to extract the posterior redshift probability distribution, $P(z | \mathbf{m}, \mathbf{C})$, which can be obtained by integrating out the nuisance parameters in Eq. 5 such that

$$P(z | \mathbf{m}, \mathbf{C}) = \int dt \int dm_0 P(z, t, m_0 | \mathbf{m}, \mathbf{C}). \quad (7)$$

In the actual implementation, the integrals are discretized, replaced by a summation over a predetermined set of template parameters, a range in flux, and the points of a predetermined redshift grid, to finally yield the expression

$$P(z | \mathbf{m}, \mathbf{C}) = \frac{\sum_t \sum_{m_0} P(z, t, m_0) \exp(-\chi^2(z, t, m_0))}{\sum_{z'} \sum_t \sum_{m_0} P(z', t, m_0) \exp(-\chi^2(z', t, m_0))}. \quad (8)$$

The exact form of prior to use has to be chosen from the list currently available (see Sect. 3.2 and Sect. 5.1), but further priors could be easily added to the code. The output result is the $P(z | \mathbf{m}, \mathbf{C})$ posterior redshift probability distribution, normalized to an integral probability of 1 on the given redshift grid.

3. Database integration

Microsoft SQL Server provides a unique opportunity among relational database management systems in its Common Language Runtime (CLR) integration feature. Compiled assemblies produced by any Common Language Infrastructure (CLI) programming language, which include C++, C# and F#, can be loaded into a DB on a server, thus granting access to code with arbitrarily complex functionality.

The main difficulty in SQL-CLR integration pertains to the SQL interface of the code, which has to translate between storage types of the programming language and the types available in SQL, and additionally it has to conform to the limitations of user-defined functions and user-defined stored procedures in SQL. Conversions between basic types such as integers, floating-point numbers and strings are fairly trivial, but more complex containers do need to be implemented on a case-by-case basis, relying on the binary storage types of SQL, or potentially on user-defined types. The publicly available *SqlArray* library (Dobos et al., 2011; Dobos et al., 2012) is a good example of this process, providing variable-length array functionality in SQL.

Most algorithms need to store data, maintain a state between different execution steps. In a SQL interface, functions called in subsequent queries would normally only be able to share data by storing it on a disk, in DB tables using SQL types. However, in our implementation we bypass this limitation, and maintain a state in-memory without needing to convert to SQL types (see Sect. 3.1 below).

The Photo-z-SQL library has been coded in C#. The interface consists of a set of user-defined SQL functions which can be used to first configure the photo-z calculation engine, and then perform the computations themselves. Parallel execution is fully supported to utilize the multi-core processors of typical DB servers. We provide install and uninstall scripts to make the installation process straightforward.

The remainder of this section gives more details about our implementation and interface in relation to the database integration.

3.1. Implementation

The data needs of the algorithm are the following:

- observed fluxes/magnitudes and errors for every object (optionally with a Schlegel et al., 1998, extinction map value),
- the transmission curves corresponding to the broadband filters,
- the spectral templates that are considered in the fit,
- information about the prior,

- additional configuration details such as the resolution in redshift and luminosity space.

In addition to the fluxes/magnitudes, the filter set may also change from object to object, but the rest of the data are static in a single estimation run, and thus can be pre-loaded and stored in the DB server. Additionally, the synthetic magnitude cache corresponding to the given configuration also has to be stored.

There is a way to perform this storage in-memory, without moving data into temporary tables: a static C# class that uses the singleton pattern will retain its state in Microsoft SQL Server between function calls. Therefore we created a singleton wrapper class that stores the photo-z configuration and synthetic magnitudes. Currently only a single, global configuration can exist in a DB, so multiple users would have to share the set of spectral templates, redshift grid and prior when estimating (but not the set of photometric filters, that can vary). In the future this can easily be extended to allow separate configurations, e.g. linked to a globally unique identifier (GUID) that corresponds to a user.

The DB server needs access to filter and template curves. While this could be solved by locally loading the required data into tables, it would put an extra burden on users to create them, and to ensure that the photo-z code accesses the data correctly. Instead, we decided to use the Spectrum Services and Filter Profile Services of the Virtual Observatory (VO)³. Users can choose from the existing templates and filters, or they can upload their own. This solution has the added benefit that the identifiers used in the Virtual Observatory can uniquely link the fluxes/magnitudes of objects to corresponding filters. The filter curves are also cached, thus they have to be downloaded from the VO only once.

The per-object data is expected to come from tables in the local DB server, supplied to the SQL functions of the Photo-z-SQL library. In order to allow variable-length array inputs of fluxes/magnitudes (and their errors), we used the *SqlArray* library (Dobos et al., 2011; Dobos et al., 2012), which can also be installed into a DB server via a simple script.

The overall setup of Photo-z-SQL is illustrated in Fig. 1.

3.2. SQL interface

Interaction with the Photo-z-SQL code is achieved with a collection of user-defined SQL functions that are called from the DB. There are three types of functions, Config functions that set up the photo-z configuration, Compute functions that perform the photo-z estimation itself, and Util functions which e.g. help interfacing with *SqlArray*. When functions require a filter or template, they can be specified either with their VO integer identifier, or a complete URL address. These two options correspond to two

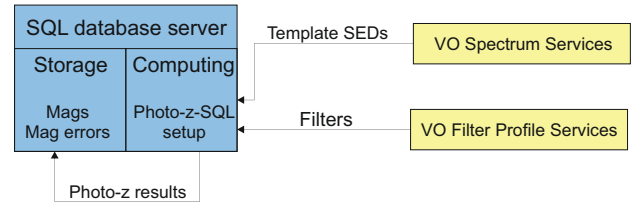


Figure 1: An illustration of the architecture of Photo-z-SQL. Computing refers to the dynamic memory and computing cores of the server, while storage refers to the disk-based data storage. Mag is shorthand for magnitude, but fluxes could also be used.

versions of the given function, with a `_ID` or `_URL` suffix in the function name, respectively. Here we briefly list the main functions and their role, and an example query of a complete setup is provided in Appendix A.

- **Config.SetupTemplateList** Specifies the list of SED templates to use in the photo-z estimation, along with the resolution in redshift and luminosity. Either the number of steps around the best-fitting luminosity can be given, or alternatively a range of physical luminosities (`_LuminositySpecified` suffix).
- **Config.SetupExtinctionLaw** Specifies the reference spectrum, R_V dust parameter and extinction law to use when correcting for Galactic extinction (see Sect. 4.4 for more details).
- **Config.RemoveExtinctionLaw** Removes the data associated with the extinction law – Galactic extinction correction is no longer applied.
- **Config.SetupFlatPrior** Specifies that a flat prior should be used (this is the default prior).
- **Config.SetupBenitezHDFPrior** Specifies that the HDF-N prior of Benítez (2000) should be used (see Sect. 5.1 for more details).
- **Config.SetupAbsoluteMagnitudeLimitPrior** Specifies a maximum limit in absolute magnitude as the prior, in the provided reference filter. Cosmological parameters are also needed to define the distance–redshift relation.
- **Config.SetupTemplateTypePrior** Specifies a prior that assigns a given probability to each of the used SED templates.
- **Config.AddMissingValueSpecifier** For convenience, specifies a value that denotes missing or otherwise invalid inputs. All such inputs are ignored.
- **Config.RemoveInitialization** Removes all data connected to the photo-z configuration, including cached filters, templates and synthetic magnitudes.

³<http://voservices.net/>

- **Compute.PhotoZMinChiSqr** Performs maximum likelihood photo-z estimation, returning a single scalar z_{phot} value (see Sect. 2.1). The prior is not used here, and luminosity is only evaluated at the best-fitting value. Magnitudes (or fluxes) and their errors, and corresponding filter identifiers are needed. An extinction map value can be specified, and the fit itself can be done in either magnitude or flux space. Also, an extra uncorrelated variance term can be added to the observational errors (see Sect. 5.1).
- **Compute.PhotoZBayesian** Performs Bayesian photo-z estimation, returning the entire posterior redshift probability distribution within the specified coverage as a SQL table (see Sect. 2.2 for more details). Input parameters are the same as in Compute.PhotoZMinChiSqr.

4. Implementation highlights

In this section, we detail some additional the features of the implementation that could be noteworthy to users.

4.1. Dynamic filter handling

The computation of synthetic magnitudes is relatively expensive, as it involves integrating a filter curve with an SED curve. This has to be done – ideally only once – for every filter and SED pairing, at every considered redshift. Available photo-z software generally solve this issue by precomputing the entire synthetic magnitude table before the actual photo-z computations. As we noted in Sect. 1, when the selection of filters changes between objects, or there is a large number of filters from which only a few are available at a time, or especially when filter time-dependence is to be taken into account, this choice is far from optimal.

In our implementation, we instead chose to set up a synthetic magnitude cache. The memory addresses of the filters serve as keys in a hash table, while the corresponding values are arrays of synthetic magnitudes. The arrays are indexed by the given template parameters and redshift. Whenever a synthetic magnitude is needed, it is retrieved from the cache if available, or computed and then stored in the cache if not. Therefore values are computed only when necessary, and only once.

With our approach, there is no requirement to know beforehand which filters will be needed. While it does include an additional hash table lookup whenever a synthetic magnitude is accessed, that is an inexpensive operation. Additionally, all functions working on the cache have been programmed to support parallel access. Currently, the cache has to be cleared by the user via a function call, but it would be possible to implement e.g. the deletion of old synthetic magnitudes after a time or cache size limit has been passed.

4.2. Magnitudes and fluxes

All of our photo-z algorithms have been programmed to be able to perform computations using either magnitudes or fluxes. When fluxes are used, the additive m_0 flux scaling parameter is replaced by a multiplicative f_0 parameter, but the best-fitting value of f_0 can still be determined algebraically for a given template SED and redshift (see Eq. 1 and Sect. 2.1).

Currently, our code supports using the AB magnitude system (Oke and Gunn, 1983), and the SDSS asinh magnitude system (Lupton et al., 1999), but additional systems could be readily included. Functions are available for converting magnitudes and magnitude errors into fluxes and flux errors, and vice versa – conversions between magnitude systems also use this mechanism. The conversions assume a Gaussian distribution for errors, therefore users should take care not to perform such conversions when this assumption does not hold, e.g. in the case of faint objects where the flux errors are Gaussian, but the magnitude errors are not.

4.3. Zeropoint calibration and error scaling

Ideally, all photometric magnitude measurements should conform to the theoretical flux zeropoints of the given magnitude system, which, for example, is 3631 Jy for the AB magnitude system (Oke and Gunn, 1983). However, there are cases when this assumption is incorrect, and the actual zeropoint of a broad-band filter is slightly different (Doi et al., 2010).

To mitigate this issue, we implemented an optional zeropoint calibration algorithm which uses objects with known redshifts. The uncorrelated assumption is adopted here (see Eq. 3), and the best-fitting template is determined at the given redshift using Eq. 4 and Eq. 2 for every object j in the T training set. Then, for each of the broad-band filters, indexed by i , a similar algebraic minimization is performed to find the filter-dependent zeropoint shift in magnitude, expressed by

$$m_i = \arg \min_{(m'_i)} \sum_{j \in T} \frac{(m_{i,j} - m'_i - (s_i(z_j, t_j) - m_{0,j}))^2}{\sigma_{i,j}^2}. \quad (9)$$

Since the zeropoint shift can change the best-fitting template t_j and the corresponding $m_{0,j}$, the two minimization steps are repeated iteratively – with the latest m_i shift applied to the $m_{i,j}$ measured magnitudes – until a chosen zeropoint precision is achieved.

There is an additional, optional refinement to the calibration that can be applied after the iteration. If our assumptions are correct, and the data are well-described by the template spectra, the residuals scaled by the estimated photometric errors should follow the standard normal distribution. This will not be the case if the photometric errors are underestimated in a band, but the errors could

be scaled by a factor

$$\alpha_i = \sqrt{\frac{1}{N_T - 1} \sum_{j \in T} \frac{(m_{i,j} - m_i - (s_i(z_j, t_j) - m_{0,j}))^2}{\sigma_{i,j}^2}} \quad (10)$$

to ensure that the scaled residual distribution

$$\sigma_{i,j}^{\text{scaled}} = \alpha_i \sigma_{i,j} \quad (11)$$

has a standard deviation of 1. Above, N_T is the number of objects in the training set, and the error scaling factor is linked to a given filter i . Thus, problematic passbands with ill-estimated errors or a higher proportion of badly matching templates can be downweighted during the actual photo- z estimation.

We note that the same calibration algorithm can also be performed using fluxes, in that case the f_i zeropoint flux multiplier is estimated for each filter.

4.4. Galactic extinction

We provide a built-in implementation for correcting broad-band magnitudes for Galactic extinction using the IR dust map of Schlegel et al. (1998).

At present, two different extinction models are supported by the code. The first computes Eq. B2 of Schlegel et al. (1998) and assumes the O’Donnell (1994) extinction law in the optical and the Cardelli et al. (1989) law in the UV and IR wavelength ranges. The second calculates Eq. A1 of Schlafly and Finkbeiner (2011) while assuming the Fitzpatrick (1999) extinction law.

The reference galaxy spectrum used in the computation can be specified freely, but it should have coverage in the entire wavelength range of the broad-band filters for which the correction is applied.

5. Photo- z results

There have been two recent large publications that allowed the comparison of photometric redshift estimation methods by publishing blind datasets for testing purposes. These two are “Photo- z Accuracy Testing” or *PHAT* by Hildebrandt et al. (2010), and “A Critical Assessment of Photometric Redshift Methods: A CANDELS Investigation” by Dahlen et al. (2013), which we will refer to as *CAPR*. We use the public datasets of these articles to demonstrate the performance of our code.

5.1. Configurations

We adopted some of the more successful approaches in the literature in our setups (Hildebrandt et al., 2010; Dahlen et al., 2013). We use two different sets of galaxy template SEDs, the Hubble UDF set of Coe et al. (2006), denoted by *BPZ*, and the COSMOS set of Ilbert et al. (2009), indicated by *LP*. In the case of the *BPZ* set, following Coe et al. (2006) we linearly interpolated 9 galaxies between each of the 8 neighboring templates to generate

a total of 71 SEDs. The wavelength coverage of these templates ends at 25600Å, after which they were linearly extrapolated. For the *LP* set, adopting the choices of Ilbert et al. (2009) we used the *Le PHARE* code (Arnouts et al., 2002; Ilbert et al., 2006) to add emission lines of different fluxes, and to apply a selection of extinction laws with different parameters to the templates, yielding a total of 641 SEDs.

The Bayesian estimation method was selected. The resolution of the redshift grid was taken to be 0.01. We use two different priors, a simple flat prior (indicated by *Flat*), and the apparent *I*-band magnitude, galaxy type and redshift prior of Benítez (2000), empirically calibrated on HDF-N data (denoted by *HDF*). The latter prior has been adapted to the *LP* template set based on galaxy type, distributing the total probability of a given type evenly among its instances. Since a measured *I*-band magnitude may not be available, the synthetic *I*-band magnitude corresponding to the given parameters is used as a proxy to allow applying this prior in all situations.

We found that adding an independent magnitude variance term to the measured magnitude variance can improve the estimation results. This extra error term can represent uncertainties in Galactic extinction, or in the template SEDs themselves, and it prevents single filters with very low errors from placing too stringent limits on the fitted templates. Thus, as an additional refinement, we may add an extra 0.01 mag of error when using the *LP* set, or 0.02 mag in the case of the *BPZ* set (shown by the tag *Err*). The exact values were chosen based on the redshift estimation results of the *PHAT* dataset, but have not been fine-tuned.

The combination of the two template sets, two priors, and whether or not extra error is added yields a total of 8 different configurations that we present.

5.2. PHAT

The PHAT1 dataset (Hildebrandt et al., 2010) contains 515 galaxies that have a spectroscopic redshift, with magnitude measurements in 14 different broad-band filters that span the wavelength range between 3000Å and 25000Å (here we do not use the extra 4 IRAC filters). The measures that we report are the $\overline{\Delta z_{\text{norm}}}$ bias and $\sigma(\Delta z_{\text{norm}})$ standard deviation – excluding outliers – of the normalized redshift error, $\Delta z_{\text{norm}} = \frac{z_{\text{spec}} - z_{\text{phot}}}{1 + z_{\text{spec}}}$, and also the P_o percentage of outliers. Outliers are defined as having $|\Delta z_{\text{norm}}| > 0.15$. z_{phot} is chosen to be the highest-probability redshift in the posterior redshift distribution defined in Eq. 8.

The results are presented in Tab. 1, while the spectroscopic-photometric redshift scatterplots are shown in Fig. 2. Applying the *HDF* prior is slightly beneficial for the smaller *BPZ* template set in that it reduces scatter while only marginally changing other measures, and it is somewhat detrimental in the case of the detailed *LP* SED set, mainly because of an increased bias. Adding the extra error term

Configuration	$\overline{\Delta z_{\text{norm}}}$	$\sigma(\Delta z_{\text{norm}})$	P_o
<i>BPZ Flat</i>	0.0114	0.0494	9.91%
<i>BPZ HDF</i>	0.0119	0.0484	9.94%
<i>BPZ Flat Err</i>	0.0029	0.0486	10.34%
<i>BPZ HDF Err</i>	0.0026	0.0462	10.77%
<i>LP Flat</i>	0.0049	0.0466	9.36%
<i>LP HDF</i>	0.0055	0.0467	9.38%
<i>LP Flat Err</i>	0.0013	0.0445	9.88%
<i>LP HDF Err</i>	0.0025	0.0448	10.08%

Table 1: Numerical results for the *PHAT* dataset. The table lists the Δz_{norm} average bias and $\sigma(\Delta z_{\text{norm}})$ standard deviation of the $\Delta z_{\text{norm}} = \frac{z_{\text{spec}} - z_{\text{phot}}}{1 + z_{\text{spec}}}$ normalized redshift estimation error, and the P_o outlier rate with $|\Delta z_{\text{norm}}| > 0.15$ outliers, for each configuration we ran.

reduces the estimation bias considerably while also slightly reducing the scatter, however, the outlier rate is marginally increased. All in all, our results are comparable to those of the better-performing methods in Table 5 of Hildebrandt et al. (2010), whose typical values were: $\overline{\Delta z_{\text{norm}}} \approx 0.004 - 0.0011$, $\sigma(\Delta z_{\text{norm}}) \approx 0.038 - 0.048$ and $P_o \approx 9.2\% - 13.5\%$.

We note that the *PHAT* dataset is not large enough to warrant performing the calibration detailed in Sect. 4.3 – with such a small sample, the calibration generally converges to a local minimum in filter zeropoints (as opposed to the global one), improving performance on the calibration set, but actually decreasing it on the validation set. Additionally, the Schlegel et al. (1998) dust map value corresponding to the galaxies was not published, and neither were the coordinates on the sky, therefore Galactic extinction has not been taken into account.

5.3. CAPR

The *CAPR* dataset (Dahlen et al., 2013) contains 589 galaxies intended for redshift estimation, along with a training set of 580 galaxies intended for calibration. There are flux measurements in 14 different passbands, of which we do not use the the IRAC $5.8\mu\text{m}$ and $8.0\mu\text{m}$ channels, similarly to code H (*Le PHARE*) in Dahlen et al. (2013). This way, the wavelength coverage is between 3000\AA and 50000\AA . The IRAC passbands are problematic because they probe wavelength ranges where the spectral templates at low redshifts are not as reliable, and where Galactic extinction is a more significant factor (Dahlen et al., 2013). As with the *PHAT* dataset, there were no published IR dust map values, therefore we did not take Galactic extinction into account.

First, our code was executed without performing any calibration. Again, we report the same numeric measures as in Sect. 5.2, collated in Tab. 2. Refer to Fig. 3 for the redshift estimation scatterplots. Whether or not the *HDF* prior was applied does not significantly impact the results, and the additional error term leads to a small improvement in the case of the *LP* template set, but slightly worse bias and scatter for the *BPZ* template set. When compared with the results in Table 2 of Dahlen et al. (2013) (bias, σ_o

Configuration	$\overline{\Delta z_{\text{norm}}}$	$\sigma(\Delta z_{\text{norm}})$	P_o
<i>BPZ Flat</i>	-0.0359	0.0732	19.84%
<i>BPZ HDF</i>	-0.0352	0.0727	20.49%
<i>BPZ Flat Err</i>	-0.0428	0.0747	17.76%
<i>BPZ HDF Err</i>	-0.0421	0.0742	18.45%
<i>LP Flat</i>	-0.0083	0.0514	8.54%
<i>LP HDF</i>	-0.0083	0.0514	8.54%
<i>LP Flat Err</i>	-0.0083	0.0492	7.65%
<i>LP HDF Err</i>	-0.0081	0.0490	7.65%

Table 2: Numerical results for the *CAPR* dataset, without calibration. The table lists the Δz_{norm} average bias and $\sigma(\Delta z_{\text{norm}})$ standard deviation of the $\Delta z_{\text{norm}} = \frac{z_{\text{spec}} - z_{\text{phot}}}{1 + z_{\text{spec}}}$ normalized redshift estimation error, and the P_o outlier rate with $|\Delta z_{\text{norm}}| > 0.15$ outliers, for each configuration we ran.

and OLF columns), where typical values were: $\overline{\Delta z_{\text{norm}}} \approx 0.005 - 0.023$, $\sigma(\Delta z_{\text{norm}}) \approx 0.034 - 0.064$ and $P_o \approx 3.9\% - 9.3\%$, the *BPZ* setups are among the worse performers, while the *LP* configurations are in the middle of the pack. However, it should be noted that the better performers all included some sort of training using the spectroscopic sample.

For the sake of a better comparison, we tested our calibration algorithm, utilizing the published training set. We selected the *BPZ HDF Err* and *LP Flat Err* configurations to illustrate the effects of calibration – the prior does not have a significant influence on results, and the added error was necessary because a few galaxies with very low estimated errors can adversely impact the calibration of a filter. As discussed in Sect. 4.3, we can either only perform zeropoint calibration (*ZP*), or both zeropoint calibration and photometric error scaling (*ZP+E*). The results are presented in Tab. 3, and the estimation scatterplots are shown in Fig. 4. Even with the calibration, the *BPZ* template set appears inadequate in describing the data, probably because its ”proper” wavelength coverage ends at 25600\AA . However, after calibration the results of the *LP* set are similar to what the better performers in Dahlen et al. (2013) can produce – the only exception is the higher outlier fraction, which is dominated by high-redshift galaxies: 49% of outliers have $z > 3$, as opposed to only 8% of the whole sample. In fact, the fraction of $z > 3$ outliers jumps from 39% to 52% because of the zeropoint calibration, and to 76% with the error scaling applied. This is because the training set is also mainly made up of low-redshift galaxies, and at different redshifts different passbands become important. Additionally, with low-redshift galaxies dominating the calibration, the less well-modeled higher wavelength ranges in the templates will lead to a larger than needed assigned uncertainty in the case of higher wavelength filters, downweighting filters that should constrain the photo- z estimation at high z . This issue could be solved by upweighting less well populated redshift bins in the calibration, and by applying a wavelength-dependent template error correction e.g. following Brammer et al. (2008).

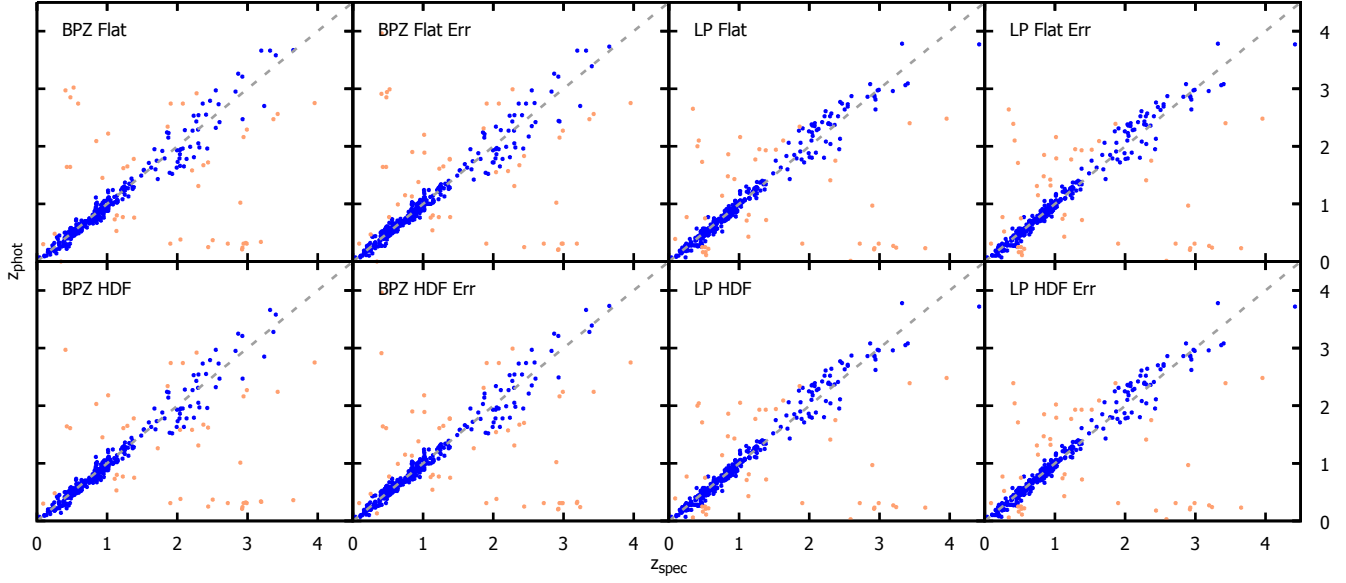


Figure 2: The z_{phot} photometric redshift as a function of the z_{spec} spectroscopic redshift, for all the configurations we ran when estimating the *PHAT* dataset. The text in the top left corner of each panel indicates the given setup. Outlying galaxies with $|\Delta z_{\text{norm}}| > 0.15$ are shown in light red, non-outlying galaxies in blue.

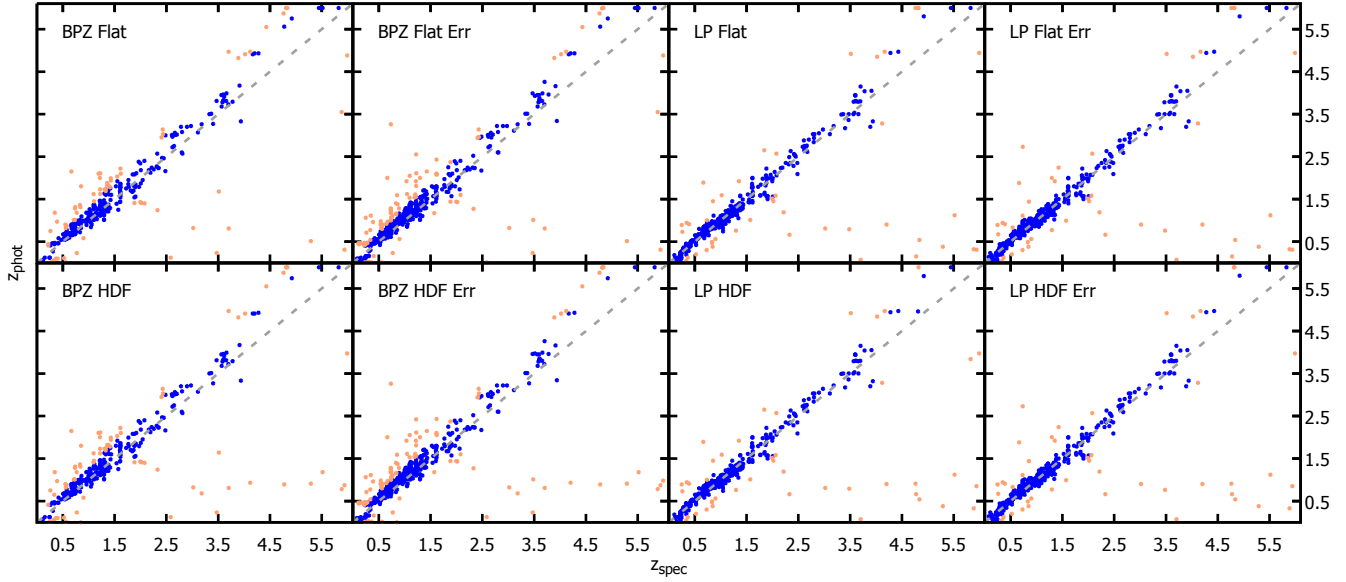


Figure 3: The z_{phot} photometric redshift as a function of the z_{spec} spectroscopic redshift, for all the configurations we ran when estimating the *CAPR* dataset, without calibration. The text in the top left corner of each panel indicates the given setup. Outlying galaxies with $|\Delta z_{\text{norm}}| > 0.15$ are shown in light red, non-outlying galaxies in blue.

Configuration	$\overline{\Delta z_{\text{norm}}}$	$\sigma(\Delta z_{\text{norm}})$	P_o
<i>BPZ HDF Err</i>	-0.0421	0.0742	18.45%
<i>BPZ HDF Err ZP</i>	0.0117	0.0615	15.45%
<i>BPZ HDF Err ZP+E</i>	0.0200	0.0631	15.62%
<i>LP Flat Err</i>	-0.0083	0.0492	7.65%
<i>LP Flat Err ZP</i>	0.0076	0.0445	9.85%
<i>LP Flat Err ZP+E</i>	0.0123	0.0402	12.05%

Table 3: Numerical results for the *CAPR* dataset, including calibration. The table lists the $\overline{\Delta z_{\text{norm}}}$ average bias and $\sigma(\Delta z_{\text{norm}})$ standard deviation of the $\Delta z_{\text{norm}} = \frac{z_{\text{spec}} - z_{\text{phot}}}{1 + z_{\text{spec}}}$ normalized redshift estimation error, and the P_o outlier rate with $|\Delta z_{\text{norm}}| > 0.15$ outliers. *ZP* denotes only zeropoint calibration, while *ZP+E* indicates zeropoint calibration and error scaling.

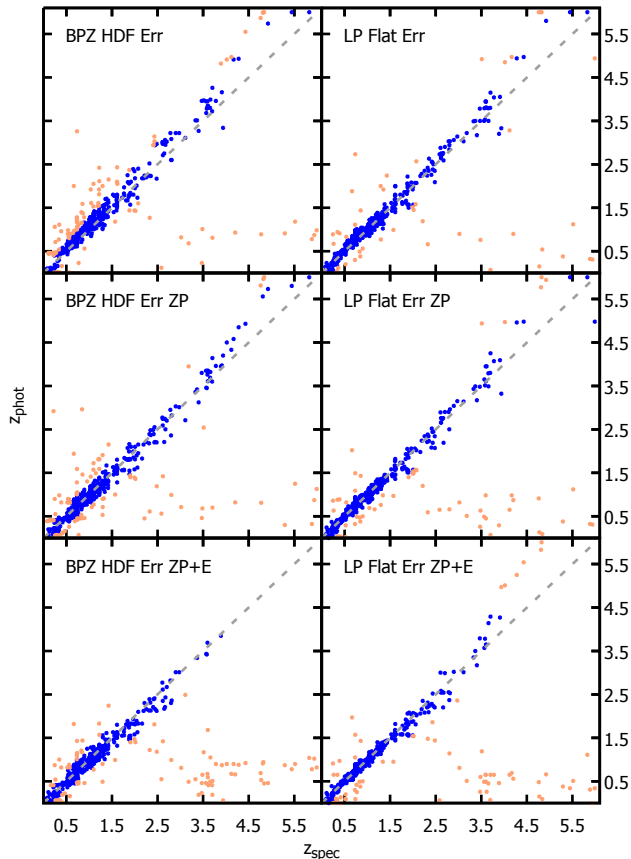


Figure 4: The z_{phot} photometric redshift as a function of the z_{spec} spectroscopic redshift, for the calibration tests we ran on the *CAPR* dataset. The text in the top left corner of each panel indicates the given setup. Outlying galaxies with $|\Delta z_{\text{norm}}| > 0.15$ are shown in light red, non-outlying galaxies in blue.

6. Performance analysis

In this section, we provide an analysis of the computational performance and memory usage of our code. We test how these factors scale when varying the details of the adopted configuration.

The test system was a slightly outdated server computer with a 16-core, 2.2 GHz Intel Xeon processor running Microsoft SQL Server 2012. With hyper-threading, this system can support 32 parallel threads.

The server could and was used by multiple users at the time, but no other queries were running that were as computationally intensive or lengthy as the Photo-z-SQL testing queries, leading to relatively stable benchmark results.

When first specifying the photo-z configuration, the synthetic magnitude cache is empty, and it will be filled during the starting phase of the first estimation query. Any subsequent queries that use the same filters will then work with the existing magnitude cache. For this reason, for each configuration we execute the same query twice, recording two times, T_{first} and T_{second} . T_{first} includes the synthetic magnitude computation for the first few galaxies, the estimation time for the rest of the galaxies, and any overhead from the initial setup and the transmission of results. T_{second} includes the estimation time for all the galaxies, and again any overhead. Therefore $T_{\text{cache}} \equiv T_{\text{first}} - T_{\text{second}}$ is a good measure of the time required to fill the magnitude cache, and using the number of galaxies N_{gal} , $T_{\text{gal}} \equiv T_{\text{second}}/N_{\text{gal}}$ is a reasonable estimate of the time required to perform photo-z on a single galaxy.

One of our goals is to quantify memory usage, however, C# is a garbage-collected language with automatically managed memory, and it would be very difficult to track the memory used by different objects in a detailed way. Still, we can measure the total CLR memory usage $\text{MEM}_{\text{total}}$, which is the final metric in our analysis.

We use the *PHAT* dataset in our tests, with $N_{\text{gal}} = 515$ (see Sect. 5.2). Each test run was performed five times, and we provide the average and standard deviation of the resulting values. In Tab. 4, we report the metrics defined above for the different configurations introduced in Sect. 5.1. We only mention the four *Err* configurations, since the increased error values do not impact performance. As expected, having ≈ 9 times as many template SEDs in the *LP* configurations scales up runtime and memory usage by a correspondingly large factor of ≈ 8 . Additionally, applying the *HDF* prior, which has to be evaluated in the innermost loop, more than doubles the per-galaxy estimation times. The synthetic magnitude computation time is not that much affected by the prior, since in that case the filter-spectrum integration is the most expensive operation.

We follow up our tests of the predetermined configurations with an analysis of how the metrics scale when different parameters are varied. Starting from the *LP Flat*

Configuration	$T_{\text{first}} (s)$	$T_{\text{second}} (s)$	$T_{\text{cache}} (s)$	$T_{\text{gal}} (s)$	$\text{MEM}_{\text{total}} (\text{MB})$
<i>BPZ Flat Err</i>	49.7 ± 3.5	18.1 ± 1.1	31.6 ± 3.7	0.035 ± 0.002	101 ± 3
<i>BPZ HDF Err</i>	77.4 ± 1.7	41.4 ± 2.8	36.0 ± 3.3	0.080 ± 0.005	107 ± 3
<i>LP Flat Err</i>	296 ± 29	147 ± 10	148 ± 31	0.286 ± 0.020	760 ± 19
<i>LP HDF Err</i>	502 ± 20	348 ± 5	154 ± 21	0.676 ± 0.010	823 ± 18

Table 4: Performance results for four configurations on the *PHAT* dataset. The table lists the execution time T_{first} of the Bayesian query for the first time (which includes synthetic magnitude computation), the execution time T_{second} of the same query for the second time (with synthetic magnitudes already cached), the cache filling time T_{cache} , the photo-z estimation time T_{gal} for a single galaxy, and the total CLR memory usage $\text{MEM}_{\text{total}}$. See the text for a discussion.

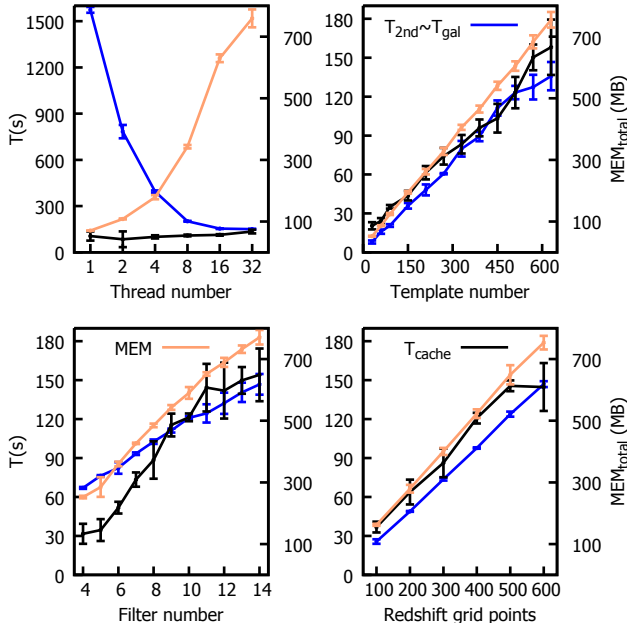


Figure 5: The scaling of the execution time T_{second} (in blue, left axes), the cache filling time T_{cache} (black, left axes), and the total memory usage $\text{MEM}_{\text{total}}$ (light red, right axes) with respect to photo-z configuration parameters. To allow the information in T_{cache} and T_{gal} to share the same y axis, we show T_{second} in place of T_{gal} , since T_{second} can be divided by 515 to yield T_{gal} . See the text for a discussion.

Err configuration, again on the *PHAT* dataset, we modify the number of templates used between 30 and 630, the number of broad-band filters considered between 4 and the original 14, the number of redshift grid points between 100 and the original 600, and finally the number of allowed parallel threads between 1 and the original 32.

Results from these scaling tests are presented in Fig. 5. The most intriguing behavior can be observed in case of the thread number: the per-galaxy execution time T_{gal} drops sharply up to 16 threads as expected, but remains roughly the same for 32 threads. Since the number of cores in the server is 16, this shows that hyper-threading is not an advantage for this algorithm. Also, the cache filling time T_{cache} remains roughly the same with an increasing thread number (in fact, it gets slightly worse due to more collisions), which is caused by the fact that all galaxies in our sample have the same filter list, filling the cache in the same order. However, this is a one-time cost, and could

be parallelized if it ever becomes an issue. The allocated memory also sharply increases with more threads, which is expected since every thread requires resources to work on.

The execution time and memory usage functions with respect to redshift grid points, template number and filter number all follow the same, linearly growing trend. This also matches our expectations, since the storage unit sizes and number of required function evaluations scale linearly with these parameters. While the T_{cache} curve apparently breaks from the trend at filter numbers above 10, this is due to the fact that those are infrared filters with a smaller resolution, hence the synthetic magnitudes are computed more quickly.

7. Summary

In this article, we presented a new C# photometric redshift estimation code, Photo-z-SQL. We detailed the photo-z approaches implemented by the code, and listed the most important available features. We demonstrated the performance of our code on two public datasets that had been used for photo-z method comparison, *PHAT* and *CAPR* (Hildebrandt et al., 2010; Dahlen et al., 2013). We also provided an analysis of how execution time scales with various configuration parameter choices.

Our photo-z estimation results are on par with those of the better performers in the literature, as expected from our adopted configurations, e.g. choice of template SEDs matching Ilbert et al. (2009). Interested readers are referred to Beck et al. (2017) for further method comparisons. While we do not introduce anything inherently new from a photo-z estimation point of view, merely adopt some of the more successful methods in the literature, our implementation does possess the following important technical advantages:

- Capability to be directly integrated into Microsoft SQL Server, eliminating the need to move photometric data outside the DB, and giving users and administrators an integrated platform for photo-z computations.
- Dynamic, on-demand handling of the set of broad-band photometric filters, which is a requirement for heterogeneous databases such as the Hubble Source Catalog (Whitmore et al., 2016).

- Ability to utilize the full covariance matrix between filters.

We should note that, while from a modeling standpoint it is better not to neglect the covariance between filters, in practice it may be difficult to appropriately determine the covariance matrix when no direct measurements are available. An empirical approach similar to the one in Sect. 4.3 might be adopted, but that solution will depend on the adopted template set and calibration sample.

Our stand-alone C# code can be compiled and executed on Windows, Unix and Mac OS X systems, but the SQL-CLR integration is strictly tied to Microsoft SQL Server – in this sense, the implementation is not portable. However, SDSS, Pan-STARRS and the Hubble Source Catalog all use Microsoft SQL Server, therefore there are plenty of useful applications even with this limitation.

The implementation is actively in development, especially regarding priors, calibration techniques and the SQL interface. The planned applications include the assembly of a photo-z table for the Hubble Source Catalog, and integration into the SDSS *SkyServer*, or the database cross-matching platform *SkyQuery*⁴ (Dobos et al., 2012; Budavari et al., 2013).

The main unsolved challenge in running photo-z on heterogeneous, cross-matched data is the handling of different aperture sizes used in different catalogs. The differing apertures will result in a per-object magnitude shift between photometric bands, dependent on the apparent size of the object on the sky, and on what portion of it is covered by the apertures. Without going to the source images, these magnitude shifts cannot be transformed out – unless there are multiple available aperture sizes for the same band, and some form of profile fitting can be utilized. More research will be needed on how strongly this issue affects photo-z results, and on how to mitigate the effects.

Another avenue for continuing this research could be the implementation of an empirical photo-z technique in C#, to be integrated into Microsoft SQL Server. While the interface for specifying the training set and the learning step in SQL would presumably be more difficult to implement and to use than the interface of the template-based method, in principle there are no technical obstacles since a pre-loaded state can be stored in a static C# class. Furthermore, the 2016 version of Microsoft SQL Server will provide support for R, potentially making it straightforward to perform empirical photo-z calculations directly in a DB.

The Photo-z-SQL code is available for download at <https://github.com/beckrob/Photo-z-SQL>.

⁴<http://www.sciserver.org/tools/skyquery/>

Acknowledgments

The realization of this work was supported by the Hungarian NKFI NN grant 114560. RB was supported through the New National Excellence Program of the Ministry of Human Capacities, Hungary.

References

- Alam, S., Albareti, F.D., Allende Prieto, C., Anders, F., Anderson, S.F., Anderton, T., Andrews, B.H., Armengaud, E., Aubourg, É., Bailey, S., et al., 2015. The Eleventh and Twelfth Data Releases of the Sloan Digital Sky Survey: Final Data from SDSS-III. *ApJS* 219, 12. doi:10.1088/0067-0049/219/1/12, arXiv:1501.00963.
- Arnouts, S., Moscardini, L., Vanzella, E., Colombi, S., Cristiani, S., Fontana, A., Giallongo, E., Matarrese, S., Saracco, P., 2002. Measuring the redshift evolution of clustering: the Hubble Deep Field South. *MNRAS* 329, 355–366. doi:10.1046/j.1365-8711.2002.04988.x, arXiv:astro-ph/0109453.
- Beck, R., Dobos, L., Budavári, T., Szalay, A.S., Csabai, I., 2016. Photometric redshifts for the SDSS Data Release 12. *MNRAS* 460, 1371–1381. doi:10.1093/mnras/stw1009, arXiv:1603.09708.
- Beck, R., Lin, C.A., Ishida, E.E.O., Gieseke, F., de Souza, R.S., Costa-Duarte, M.V., Hattab, M.W., Krone-Martins, A., for the COIN Collaboration, 2017. On the realistic validation of photometric redshifts, or why Teddy will never be Happy. *ArXiv e-prints* arXiv:1701.08748.
- Benítez, N., 2000. Bayesian Photometric Redshift Estimation. *ApJ* 536, 571–583. doi:10.1086/308947, arXiv:astro-ph/9811189.
- Bolzonella, M., Miralles, J.M., Pelló, R., 2000. Photometric redshifts based on standard SED fitting procedures. *A&A* 363, 476–492. arXiv:astro-ph/0003380.
- Brammer, G.B., van Dokkum, P.G., Coppi, P., 2008. EAZY: A Fast, Public Photometric Redshift Code. *ApJ* 686, 1503–1513. doi:10.1086/591786, arXiv:0807.1533.
- Brescia, M., Cavuoti, S., Longo, G., De Stefano, V., 2014. A catalogue of photometric redshifts for the SDSS-DR9 galaxies. *A&A* 568, A126. doi:10.1051/0004-6361/201424383, arXiv:1407.2527.
- Budavári, T., 2009. A Unified Framework for Photometric Redshifts. *ApJ* 695, 747–754. doi:10.1088/0004-637X/695/1/747, arXiv:0811.2600.
- Budavari, T., Dobos, L., Szalay, A.S., 2013. Skyquery: Federating astronomy archives. *Computing in Science and Engineering* 15, 12–20. URL: <http://dx.doi.org/10.1109/MCSE.2013.41>, doi:10.1109/MCSE.2013.41.
- Cardelli, J.A., Clayton, G.C., Mathis, J.S., 1989. The relationship between infrared, optical, and ultraviolet extinction. *ApJ* 345, 245–256. doi:10.1086/167900.
- Carliles, S., Budavári, T., Heinis, S., Priebe, C., Szalay, A.S., 2010. Random Forests for Photometric Redshifts. *ApJ* 712, 511–515. doi:10.1088/0004-637X/712/1/511.
- Coe, D., Benítez, N., Sánchez, S.F., Jee, M., Bouwens, R., Ford, H., 2006. Galaxies in the Hubble Ultra Deep Field. I. Detection, Multiband Photometry, Photometric Redshifts, and Morphology. *AJ* 132, 926–959. doi:10.1086/505530, arXiv:astro-ph/0605262.
- Connolly, A.J., Csabai, I., Szalay, A.S., Koo, D.C., Kron, R.G., Munn, J.A., 1995. Slicing Through Multicolor Space: Galaxy Redshifts from Broadband Photometry. *AJ* 110, 2655. doi:10.1086/117720, arXiv:astro-ph/9508100.
- Csabai, I., Budavári, T., Connolly, A.J., Szalay, A.S., Gyóry, Z., Benítez, N., Annis, J., Brinkmann, J., Eisenstein, D., Fukugita, M., Gunn, J., Kent, S., Lupton, R., Nichol, R.C., Stoughton, C., 2003. The Application of Photometric Redshifts to the SDSS Early Data Release. *AJ* 125, 580–592. doi:10.1086/345883, arXiv:astro-ph/0211080.
- Csabai, I., Connolly, A.J., Szalay, A.S., Budavári, T., 2000. Reconstructing Galaxy Spectral Energy Distributions from Broadband Photometry. *AJ* 119, 69–78. doi:10.1086/301159, arXiv:astro-ph/9910389.

- Csabai, I., Dobos, L., Trencsényi, M., Herczegh, G., Józsa, P., Purger, N., Budavári, T., Szalay, A.S., 2007. Multidimensional indexing tools for the virtual observatory. *Astronomische Nachrichten* 328, 852. doi:10.1002/asna.200710817.
- Dahlen, T., Mobasher, B., Faber, S.M., Ferguson, H.C., Barro, G., Finkelstein, S.L., Finlator, K., Fontana, A., Gruetzbauch, R., Johnson, S., Pforr, J., Salvato, M., Wiklind, T., Wuyts, S., Acquaviva, V., Dickinson, M.E., Guo, Y., Huang, J., Huang, K.H., Newman, J.A., Bell, E.F., Conselice, C.J., Galametz, A., Gawiser, E., Giavalisco, M., Grogin, N.A., Hathi, N., Kocevski, D., Koekemoer, A.M., Koo, D.C., Lee, K.S., McGrath, E.J., Papovich, C., Peth, M., Ryan, R., Somerville, R., Weiner, B., Wilson, G., 2013. A Critical Assessment of Photometric Redshift Methods: A CANDELS Investigation. *ApJ* 775, 93. doi:10.1088/0004-637X/775/2/93, arXiv:1308.5353.
- Dobos, L., Budavari, T., Li, N., Szalay, A.S., Csabai, I., 2012. Skyquery: An implementation of a parallel probabilistic join engine for cross-identification of multiple astronomical databases, in: *Scientific and Statistical Database Management - 24th International Conference, SSDBM 2012, Chania, Crete, Greece, June 25-27, 2012. Proceedings*, pp. 159–167. URL: http://dx.doi.org/10.1007/978-3-642-31235-9_10, doi:10.1007/978-3-642-31235-9_10.
- Dobos, L., Szalay, A.S., Blakeley, J., Falck, B., Budavári, T., Csabai, I., 2012. An Array Library for Microsoft SQL Server with Astrophysical Applications, in: *Ballester, P., Egret, D., Lorente, N.P.F. (Eds.), Astronomical Data Analysis Software and Systems XXI*, p. 323.
- Dobos, L., Szalay, A.S., Blakeley, J.A., Budavari, T., Csabai, I., Tomic, D., Milovanovic, M., Tintor, M., Jovanovic, A., 2011. Array requirements for scientific applications and an implementation for microsoft SQL server, in: *Proceedings of the 2011 EDBT/ICDT Workshop on Array Databases*, Uppsala, Sweden, March 25, 2011, pp. 13–19. URL: <http://doi.acm.org/10.1145/1966895.1966897>, doi:10.1145/1966895.1966897.
- Doi, M., Tanaka, M., Fukugita, M., Gunn, J.E., Yasuda, N., Ivezić, Z., Brinkmann, J., de Haars, E., Kleinman, S.J., Krzesinski, J., French Leger, R., 2010. Photometric Response Functions of the Sloan Digital Sky Survey Imager. *AJ* 139, 1628–1648. doi:10.1088/0004-6256/139/4/1628, arXiv:1002.3701.
- Eisenstein, D.J., Weinberg, D.H., Agol, E., Aihara, H., Allende Prieto, C., Anderson, S.F., Arns, J.A., Aubourg, É., Bailey, S., Balbinot, E., et al., 2011. SDSS-III: Massive Spectroscopic Surveys of the Distant Universe, the Milky Way, and Extra-Solar Planetary Systems. *AJ* 142, 72. doi:10.1088/0004-6256/142/3/72, arXiv:1101.1529.
- Fitzpatrick, E.L., 1999. Correcting for the Effects of Interstellar Extinction. *PASP* 111, 63–75. doi:10.1086/316293, arXiv:astro-ph/9809387.
- Gerdes, D.W., Sypniewski, A.J., McKay, T.A., Hao, J., Weis, M.R., Wechsler, R.H., Busha, M.T., 2010. ArborZ: Photometric Redshifts Using Boosted Decision Trees. *ApJ* 715, 823–832. doi:10.1088/0004-637X/715/2/823, arXiv:0908.4085.
- Han, B., Zhang, Y., Zhong, S., Zhao, Y., 2016. Astronomical Data Fusion Tool Based on PostgreSQL. *ArXiv e-prints* arXiv:1609.01079.
- Hildebrandt, H., Arnouts, S., Capak, P., Moustakas, L.A., Wolf, C., Abdalla, F.B., Assef, R.J., Banerji, M., Benítez, N., Brammer, G.B., Budavári, T., Carliles, S., Coe, D., Dahlen, T., Feldmann, R., Gerdes, D., Gillis, B., Ilbert, O., Kotulla, R., Lahav, O., Li, I.H., Miralles, J.M., Purger, N., Schmidt, S., Singal, J., 2010. PHAT: PHoto-z Accuracy Testing. *A&A* 523, A31. doi:10.1051/0004-6361/201014885, arXiv:1008.0658.
- Ilbert, O., Arnouts, S., McCracken, H.J., Bolzonella, M., Bertin, E., Le Fèvre, O., Mellier, Y., Zamorani, G., Pellò, R., Iovino, A., Tresse, L., Le Brun, V., Bottini, D., Garilli, B., Maccagni, D., Picat, J.P., Scaramella, R., Scodreggio, M., Vettolani, G., Zanicchelli, A., Adami, C., Bardelli, S., Cappi, A., Charlot, S., Ciliegli, P., Contini, T., Cucciati, O., Foucaud, S., Franzetti, P., Gavignaud, I., Guzzo, L., Marano, B., Marinoni, C., Mazure, A., Meneux, B., Merighi, R., Paltani, S., Pollo, A., Pozzetti, L., Radovich, M., Zucca, E., Bondi, M., Bongiorno, A., Busarello, G., de La Torre, S., Gregorini, L., Lamareille, F., Mathez, G., Merluzzi, P., Ripepi, V., Rizzo, D., Vergani, D., 2006. Accurate photometric redshifts for the CFHT legacy survey calibrated using the VIMOS VLT deep survey. *A&A* 457, 841–856. doi:10.1051/0004-6361:20065138, arXiv:astro-ph/0603217.
- Ilbert, O., Capak, P., Salvato, M., Aussel, H., McCracken, H.J., Sanders, D.B., Scoville, N., Kartaltepe, J., Arnouts, S., Le Floc'h, E., Mobasher, B., Taniguchi, Y., Lamareille, F., Leauthaud, A., Sasaki, S., Thompson, D., Zamojski, M., Zamorani, G., Bardelli, S., Bolzonella, M., Bongiorno, A., Brusa, M., Caputi, K.I., Carollo, C.M., Contini, T., Cook, R., Coppa, G., Cucciati, O., de la Torre, S., de Ravel, L., Franzetti, P., Garilli, B., Hasinger, G., Iovino, A., Kampczyk, P., Kneib, J.P., Knobel, C., Kovac, K., Le Borgne, J.F., Le Brun, V., Fèvre, O.L., Lilly, S., Looper, D., Maier, C., Mainieri, V., Mellier, Y., Mignoli, M., Murayama, T., Pellò, R., Peng, Y., Pérez-Montero, E., Renzi, A., Ricciardelli, E., Schiminovich, D., Scodreggio, M., Shioya, Y., Silverman, J., Surace, J., Tanaka, M., Tasca, L., Tresse, L., Vergani, D., Zucca, E., 2009. Cosmos Photometric Redshifts with 30-Bands for 2-deg². *ApJ* 690, 1236–1249. doi:10.1088/0004-637X/690/2/1236, arXiv:0809.2101.
- Ivezić, Z., Tyson, J.A., Abel, B., Acosta, E., Allsman, R., AlSayyad, Y., Anderson, S.F., Andrew, J., Angel, R., Angeli, G., Ansari, R., Antilogus, P., Arndt, K.T., Astier, P., Aubourg, E., Axelrod, T., Bard, D.J., Barr, J.D., Barrau, A., Bartlett, J.G., Bauman, B.J., Beaumont, S., Becker, A.C., Becla, J., Beldica, C., Bellavia, S., Blanc, G., Blandford, R.D., Bloom, J.S., Bogart, J., Borne, K., Bosch, J.F., Boutigny, D., Brandt, W.N., Brown, M.E., Bullock, J.S., Burchat, P., Burke, D.L., Cagnoli, G., Calabrese, D., Chandrasekharan, S., Chesley, S., Cheu, E.C., Chiang, J., Claver, C.F., Connolly, A.J., Cook, K.H., Cooray, A., Covey, K.R., Cribbs, C., Cui, W., Cutri, R., Daubard, G., Daues, G., Delgado, F., Digel, S., Doherty, P., Dubois, R., Dubois-Felsmann, G.P., Durech, J., Eracleous, M., Ferguson, H., Frank, J., Freemon, M., Gangler, E., Gawiser, E., Geary, J.C., Gee, P., Geha, M., Gibson, R.R., Gilmore, D.K., Glanzman, T., Goodenow, I., Gressler, W.J., Gris, P., Guyonnet, A., Hascall, P.A., Haupt, J., Hernandez, F., Hogan, C., Huang, D., Huffer, M.E., Innes, W.R., Jacoby, S.H., Jain, B., Jee, J., Jernigan, J.G., Jevremovic, D., Johns, K., Jones, R.L., Juramy-Gilles, C., Juric, M., Kahn, S.M., Kalirai, J.S., Kallivayalil, N., Kalmbach, B., Kantor, J.P., Kasliwal, M.M., Kessler, R., Kirkby, D., Knox, L., Kotov, I., Krabbenham, V.L., Krughoff, S., Kubanek, P., Kuczewski, J., Kulkarni, S., Lambert, R., Le Guillou, L., Levine, D., Liang, M., Lim, K., Lintott, C., Lupton, R.H., Mahabal, A., Marshall, P., Marshall, S., May, M., McKecher, R., Migliore, M., Miller, M., Mills, D.J., Monet, D.G., Moniez, M., Neill, D.R., Nief, J., Nomerotski, A., Nordby, M., O'Connor, P., Oliver, J., Olivier, S.S., Olsen, K., Ortiz, S., Owen, R.E., Pain, R., Peterson, J.R., Petry, C.E., Pierfederici, F., Pietrowicz, S., Pike, R., Pinto, P.A., Plante, R., Plate, S., Price, P.A., Prouza, M., Radeka, V., Rajagopal, J., Rasmussen, A., Regnault, N., Ridgway, S.T., Ritz, S., 2008. LSST: from Science Drivers to Reference Design and Anticipated Data Products. *ArXiv e-prints* arXiv:0805.2366.
- Lupton, R.H., Gunn, J.E., Szalay, A.S., 1999. A Modified Magnitude System that Produces Well-Behaved Magnitudes, Colors, and Errors Even for Low Signal-to-Noise Ratio Measurements. *AJ* 118, 1406–1410. doi:10.1086/301004, arXiv:astro-ph/9903081.
- O'Donnell, J.E., 1994. R_{nu} -dependent optical and near-ultraviolet extinction. *ApJ* 422, 158–163. doi:10.1086/173713.
- Oke, J.B., Gunn, J.E., 1983. Secondary standard stars for absolute spectrophotometry. *ApJ* 266, 713–717. doi:10.1086/160817.
- Schlafly, E.F., Finkbeiner, D.P., 2011. Measuring Reddening with Sloan Digital Sky Survey Stellar Spectra and Recalibrating SFD. *ApJ* 737, 103. doi:10.1088/0004-637X/737/2/103, arXiv:1012.4804.
- Schlegel, D.J., Finkbeiner, D.P., Davis, M., 1998. Maps of Dust Infrared Emission for Use in Estimation of Reddening and Cosmic Microwave Background Radiation Foregrounds. *ApJ* 500, 525–553. doi:10.1086/305772, arXiv:astro-ph/9710327.

Scranton, R., Connolly, A.J., Szalay, A.S., Lupton, R.H., Johnston, D., Budavari, T., Brinkman, J., Fukugita, M., 2005. Photometric Covariance in Multi-Band Surveys: Understanding the Photometric Error in the SDSS. ArXiv Astrophysics e-prints arXiv:astro-ph/0508564.

Tonry, J.L., Stubbs, C.W., Lykke, K.R., Doherty, P., Shivvers, I.S., Burgett, W.S., Chambers, K.C., Hodapp, K.W., Kaiser, N., Kudritzki, R.P., Magnier, E.A., Morgan, J.S., Price, P.A., Wainscoat, R.J., 2012. The Pan-STARRS1 Photometric System. ApJ 750, 99. doi:10.1088/0004-637X/750/2/99, arXiv:1203.0297.

Wadadekar, Y., 2005. Estimating Photometric Redshifts Using Support Vector Machines. PASP 117, 79–85. doi:10.1086/427710, arXiv:astro-ph/0412005.

Wang, Y., Bahcall, N., Turner, E.L., 1998. A Catalog of Color-based Redshift Estimates for $Z < 4$ Galaxies in the Hubble Deep Field. AJ 116, 2081–2085. doi:10.1086/300592, arXiv:astro-ph/9804195.

Whitmore, B.C., Allam, S.S., Budavári, T., Casertano, S., Downes, R.A., Donaldson, T., Fall, S.M., Lubow, S.H., Quick, L., Strolger, L.G., Wallace, G., White, R.L., 2016. Version 1 of the Hubble Source Catalog. AJ 151, 134. doi:10.3847/0004-6256/151/6/134, arXiv:1602.04861.

York, D.G., Adelman, J., Anderson, Jr., J.E., Anderson, S.F., Annis, J., Bahcall, N.A., Bakken, J.A., Barkhouser, R., Bastian, S., Berman, E., Boroski, W.N., Bracker, S., Briegel, C., Briggs, J.W., Brinkmann, J., Brunner, R., Burles, S., Carey, L., Carr, M.A., Castander, F.J., Chen, B., Colestock, P.L., Connolly, A.J., Crocker, J.H., Csabai, I., Czarapata, P.C., Davis, J.E., Doi, M., Dombeck, T., Eisenstein, D., Ellman, N., Elms, B.R., Evans, M.L., Fan, X., Federwitz, G.R., Fiscelli, L., Friedman, S., Frieman, J.A., Fukugita, M., Gillespie, B., Gunn, J.E., Gurbani, V.K., de Haas, E., Haldeman, M., Harris, F.H., Hayes, J., Heckman, T.M., Hennessy, G.S., Hindsley, R.B., Holm, S., Holmgren, D.J., Huang, C.h., Hull, C., Husby, D., Ichikawa, S.I., Ichikawa, T., Ivezić, Ž., Kent, S., Kim, R.S.J., Kinney, E., Klaene, M., Kleinman, A.N., Kleinman, S., Knapp, G.R., Korienek, J., Kron, R.G., Kunszt, P.Z., Lamb, D.Q., Lee, B., Leger, R.F., Lim-mongkol, S., Lindenmeyer, C., Long, D.C., Loomis, C., Loveday, J., Lucinio, R., Lupton, R.H., MacKinnon, B., Mannery, E.J., Mantsch, P.M., Margon, B., McGehee, P., McKay, T.A., Meiksin, A., Merelli, A., Monet, D.G., Munn, J.A., Narayanan, V.K., Nash, T., Neilsen, E., Neswold, R., Newberg, H.J., Nichol, R.C., Nicinski, T., Nonino, M., Okada, N., Okamura, S., Ostriker, J.P., Owen, R., Pauls, A.G., Peoples, J., Peterson, R.L., Petravick, D., Pier, J.R., Pope, A., Pordes, R., Prosapio, A., Rechenmacher, R., Quinn, T.R., Richards, G.T., Richmond, M.W., Rivetta, C.H., Rockosi, C.M., Ruthmansdorfer, K., Sandford, D., Schlegel, D.J., Schneider, D.P., Sekiguchi, M., Sergey, G., Shimasaku, K., Siegmund, W.A., Smeed, S., Smith, J.A., Snedden, S., Stone, R., Stoughton, C., Strauss, M.A., Stubbs, C., SubbaRao, M., Szalay, A.S., Szapudi, I., Szokoly, G.P., Thakar, A.R., Tremonti, C., Tucker, D.L., Uomoto, A., Vanden Berk, D., Vo-geley, M.S., Waddell, P., Wang, S.i., Watanabe, M., Weinberg, D.H., Yanny, B., Yasuda, N., SDSS Collaboration, 2000. The Sloan Digital Sky Survey: Technical Summary. AJ 120, 1579–1587. doi:10.1086/301513, arXiv:astro-ph/0006396.

Appendix A. Example SQL queries

The first example sets up a photo-z configuration. First it clears the potentially existing photo-z setup, specifies that -9999 in the input data denotes a missing value, and chooses a flat prior. Then, it creates a string with VO template identifiers corresponding to the *LP* template set (see Sect. 5.1), and parses it into a *SqlArray* object. Finally, this template list is given to the photo-z code, with a redshift coverage between $0.001 - 1.001$ and a linear step size of 0.01 . 11 steps will be taken in luminosity space around the best-fitting luminosity.

```
SELECT PhotoZSQL.Config.RemoveInitialization()

SELECT PhotoZSQL.Config.AddMissingValueSpecifier(-9999)

SELECT PhotoZSQL.Config.SetupFlatPrior()

--To get a SqlArray int array of template IDs, a string
--has to be set up with this format: '[511,512,...],1151]'
DECLARE @IDString varchar(max) = '['
DECLARE @id INT = 511
WHILE @id < 1151
BEGIN
    SET @IDString = @IDString + CAST(@id AS varchar(max)) + ', '
    SET @id = @id + 1
END
SET @IDString = @IDString + CAST(@id AS varchar(max)) + ']'

--Parsing the string into an integer array
DECLARE @TemplateIDArray varbinary(max) =
    SqlArray.IntArrayMax.ParseInvariant(@IDString)

SELECT PhotoZSQL.Config.SetupTemplateList_ID(@TemplateIDArray,
    1.0,0.001,1.001,0.01,0,11)
```

The second example performs maximum likelihood photo-z estimation on a sample of SDSS data, using the 5 SDSS filters. The fit is performed in magnitude space, with no additional extinction correction, and with a 0.01 mag independent error term added.

```
DECLARE @FilterIDArray varbinary(max) =
    SqlArray.IntArrayMax.Vector_5(14,15,16,17,18)

SELECT TOP 100 gal.objID,
    PhotoZSQL.[Compute].PhotoZMinChiSqr_ID(
        SqlArray.FloatArrayMax.Vector_5(gal.dered_u,
            gal.dered_g,
            gal.dered_r,
            gal.dered_i,
            gal.dered_z),
        SqlArray.FloatArrayMax.Vector_5(gal.modelMagErr_u,
            gal.modelMagErr_g,
            gal.modelMagErr_r,
            gal.modelMagErr_i,
            gal.modelMagErr_z),
        1,@FilterIDArray,0.0,0,0.01) AS zphot
FROM DR12.Galaxy AS gal
```

The third example performs Bayesian photo-z estimation, otherwise with the same particulars as the previous query.

```
DECLARE @FilterIDArray varbinary(max) =
    SqlArray.IntArrayMax.Vector_5(14,15,16,17,18)
```

```
SELECT TOP 100 gal.objID,pz.*
FROM DR12.Galaxy AS gal
CROSS APPLY PhotoZSQL.[Compute].PhotoZBayesian_ID(
    SqlArray.FloatArrayMax.Vector_5(gal.dered_u,
        gal.dered_g,
        gal.dered_r,
        gal.dered_i,
        gal.dered_z),
    SqlArray.FloatArrayMax.Vector_5(gal.modelMagErr_u,
        gal.modelMagErr_g,
        gal.modelMagErr_r,
        gal.modelMagErr_i,
        gal.modelMagErr_z),
    1,@FilterIDArray,0.0,0,0.01) AS pz
```