

# Large Synoptic Survey Telescope Data Products Definition Document (\*\*\* DRAFT \*\*\*)

Mario Jurić (mjuric@lsst.org)

*with input from*

T. Axelrod, A.C. Becker, J. Becla, J. Kantor, K-T Lim,  
R. Lupton, M. Strauss, *and* J.A. Tyson

*for LSST Data Management*

April 19, 2013

## Abstract

This document describes the plans for contents of Level 1 and 2 LSST data products, and the rationale behind various choices that were made. This is an **internal draft** and a work in progress. **It should not be circulated further until this notice is removed.**

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
1.1	Level 1 and 2 Data Products . . . . .	3
<b>2</b>	<b>Level 1 Data Products</b>	<b>4</b>
2.1	Overview . . . . .	4
2.2	Level 1 Data Processing . . . . .	5
2.2.1	Difference Image Analysis . . . . .	5
2.2.2	Solar System Object Processing . . . . .	7
2.3	The Level 1 database . . . . .	8

2.3.1	DIASource Table . . . . .	9
2.3.2	DIAObject Table . . . . .	13
2.3.3	SSObject Table . . . . .	14
2.3.4	Difference Images . . . . .	16
2.3.5	Image Differencing Templates . . . . .	16
2.3.6	Statistical Considerations . . . . .	16
2.3.7	Fluxes and Magnitudes . . . . .	17
2.3.8	Precovery . . . . .	17
2.3.9	Annual Reprocessings . . . . .	18
2.3.10	Repeatability of Queries . . . . .	20
2.3.11	Uniqueness of IDs across database versions . . . . .	20
2.4	Alerts to DIASources . . . . .	20
2.4.1	Information Contained in Each Alert . . . . .	20
2.4.2	Receiving and Filtering the Alerts . . . . .	21
2.5	Open Issues to be Closed before Baselineing this Document . .	22

## 1 Introduction

LSST will be a large, wide-field ground-based optical telescope system designed to obtain multiple images covering the sky that is visible from Cerro Pachón in Northern Chile. The current baseline design, with an 8.4m (6.7m effective) primary mirror, a 9.6 deg<sup>2</sup> field of view, and a 3.2 Gigapixel camera, will allow about 10,000 square degrees of sky to be covered using pairs of 15-second exposures twice per night every three nights on average, with typical 5 $\sigma$  depth for point sources of  $r \sim 24.5$  (AB). The system is designed to yield high image quality as well as superb astrometric and photometric accuracy. The total survey area will include 30,000 deg<sup>2</sup> with  $\delta < +34.5^\circ$ , and will be imaged multiple times in six bands, *ugrizy*, covering the wavelength range 320–1050 nm. The project is scheduled to begin the regular survey operations at the start of next decade. About 90% of the observing time will be devoted to a deep-wide-fast survey mode which will uniformly observe a 18,000 deg<sup>2</sup> region about 1000 times (summed over all six bands) during the anticipated 10 years of operations, and yield a coadded map to  $r \sim 27.5$ . These data will result in databases including 10 billion galaxies and a similar number of stars, and will serve the majority of the primary science programs. The remaining 10% of the observing time will be allocated to special projects such as a Very Deep and Fast time domain survey.

The LSST will be operated in fully automated survey mode. The images acquired by the LSST Camera will be processed by LSST Data Management software to a) detect and characterize imaged astrophysical sources and b) detect and characterize changes in time in LSST-observed universe. The results of that processing will be reduced images, catalogs of detected objects and the measurements of their properties, and prompt alerts to “events” – changes in astrophysical scenery discovered by differencing incoming images against older, deeper, images of the sky in the same direction (*templates*, see §2.3.5).

The *broad, high-level*, requirements for LSST Data Products are given by the LSST Science Requirements Document. This document lays out the *specifics* of what the data products will comprise of, how those data will be generated, and when. It informs the flow-down from the LSST Science Requirements Document and the LSST Observatory System Specifications document to the LSST Data Management System Requirements document, the UML model, and the database schema.

## 1.1 Level 1 and 2 Data Products

LSST Data Management will perform two, somewhat overlapping in scientific intent, types of image analyses:

1. Analysis of difference images, with the goal of detecting and characterizing astrophysical phenomena revealed by their time-dependent nature. The detection of supernovae superimposed on bright extended galaxies is an example of this analysis. The processing is done on a nightly or daily basis and produces **Level 1** data products. They include the difference images, the sources detected in difference images (**DIASources**), astrophysical objects<sup>1</sup> these are associated to (**DIAObjects**), and Solar System objects (**SSObjects**<sup>2</sup>). These are added to the **Level 1 database** and made available in real time. Notifications (“alerts”)

---

<sup>1</sup>The LSST has adopted the terminology where single-epoch detections of astrophysical phenomena, of astrophysical *objects*, are called *sources*. Note that some other surveys call *detections* what we call *sources*, and use the name *sources* for what we call *objects*.

<sup>2</sup>**SSObject** used to be called call a “Moving Object”. The name is potentially confusing, as high-proper motion stars are moving objects as well. A more accurate distinction is the one between objects in an out of the Solar System.

about new **DIASources** are issued as **VOEvents** within 60 seconds of observation.

2. Analysis of science images, with the goal of detecting and characterizing astrophysical objects. Detection of faint galaxies on deep co-adds and their subsequent characterization is an example of this analysis. The results are **Level 2** data products. These products, released annually<sup>3</sup>, will include catalogs of **Objects** (detections on deep co-adds) and **Sources**<sup>4</sup> (measurements on individual science images), as well as fully reprocessed Level 1 data products (see §2.3.9). In contrast to the Level 1 database, which is updated in real-time, the Level 2 databases are static and will not change after release.

The two types of analyses have different requirements on timeliness. Changes in flux or position of objects may need to be immediately followed up, lest interesting information be lost. Thus the primary results of analysis of difference images – discovered and characterized **DIASources** – generally need to be broadcast as *event alerts* within 60 seconds of end of visit<sup>5</sup> acquisition. The analysis of science images is less time sensitive, and will be done as a part of annual data release process.

## 2 Level 1 Data Products

### 2.1 Overview

Level 1 data products are a result of difference image analysis (DIA; §2.2.1). They include the sources detected in difference images (**DIASources**), astrophysical objects that these are associated to (**DIAObjects**), identified Solar System objects<sup>6</sup> (**SSObject**), and related, broadly defined, metadata (includ-

---

<sup>3</sup>Except for the first two data releases, which will be created six months apart.

<sup>4</sup>When written in bold monospace type, **Objects** and **Sources** specifically refer to objects and sources detected and measured as a part of Level 2 processing.

<sup>5</sup>The LSST takes two (nominally 15 second) exposures per pointing, called *snaps*. That pair of exposures is called a *visit*.

<sup>6</sup>The LSST SRD considers Solar System object orbit catalog to be a Level 2 data product (LSSTSRD, Sec 3.5). Nevertheless, to successfully differentiate between apparitions of known Solar System objects and other types **DIASources** we consider it functionally a part of Level 1.

ing e.g., cut-outs<sup>7</sup>).

**DIASources** are sources detected on difference images (those above  $S/N = 5$  after correlation with an appropriate PSF profile). They represent changes in flux wrt. to the deep template. Physically, a **DIASource** may be an observation of new astrophysical object that was not present at that position in the template image (for example, an asteroid), or an observation of flux change in an existing source (for example, a variable star). Their flux can be negative (eg., if a source present in the template image reduced its brightness, or moved away).

**DIASources** detected on visits taken at different times are associated with **DIAObjects**. **DIAObjects** represent the underlying astrophysical phenomenon detected and measured by individual **DIASources**. The association can be done in two different ways: by assuming the underlying phenomenon is an object within the Solar System moving on an orbit around the Sun<sup>8</sup>, or by assuming the underlying phenomenon is distant enough to only exhibit small parallactic and proper motion<sup>9</sup>. The latter type of association is performed during difference image analysis right after the image has been acquired. The former is done at daytime by the Moving Objects Processing Software (MOPS), unless the **DIASource** is an apparition of an already known Solar System object (“**SSObjects**”) in which case it will be flagged as such during difference image analysis. All **DIASources** will be alerted on at the end of the difference image analysis<sup>10</sup>.

## 2.2 Level 1 Data Processing

### 2.2.1 Difference Image Analysis

The following is a high-level description of steps which will occur during normal difference image analysis:

---

<sup>7</sup>Small,  $30 \times 30$ , sub-images at the position of a detected source. Also known as *postage stamps*.

<sup>8</sup>We don’t plan to fit for motion around other Solar System bodies; eg., identifying new satellites of Jupiter is left to the community.

<sup>9</sup>Where ‘small’ is small enough to unambiguously positionally associate together individual apparitions of the object.

<sup>10</sup>For observations on the ecliptic near the opposition, Solar System objects will dominate the **DIASource** counts, and (until they’re recognized as such) overwhelm the explosive transient signal. It will therefore be advantageous to quickly identify the majority of Solar System objects early in the survey.

1. A visit is acquired and the images reduced and combined to a single science image (cosmic ray rejection, instrumental signature removal<sup>11</sup>, combining of snaps, etc.).
2. The visit image is differenced against the appropriate template and **DIASources** are detected.
3. The flux and shape<sup>12</sup> of the **DIASource** are measured on the difference image. The science image is force-photometered at the position of the **DIASource** to obtain a measure of the absolute flux. No deblending will be attempted.
4. The Level 1 database (see §2.3) is searched for a **DIAObject** or **SSObject** with which to positionally associate the observed **DIASource**<sup>13</sup>. If no match is found, a new **DIAObject** is created and the observed **DIASource** is associated to it.
5. If the **DIASource** has been associated with an **SSObject** (a known moving object), it will be flagged as such and an alert issued. Further processing will occur in daytime (see section 2.2.2).
6. Otherwise, the associated **DIAObject** measurements are updated with new data. All affected columns are recomputed, including proper motions, centroids, light curves, etc.
7. The Level 2 database<sup>14</sup> is searched for one or more **Objects** positionally close to the **DIAObject**, out to some maximum radius<sup>15</sup>. The IDs of these **Objects** are recorded in the **DIAObject** record and provided in the event alert.

---

<sup>11</sup>E.g., subtraction of bias and dark frames, flat fielding, bad pixel/column interpolation, etc.

<sup>12</sup>The “shape” in this context are weighted 2<sup>nd</sup> moments, as well as a fit to a trailed source model.

<sup>13</sup>The association algorithm will guarantee that a **DIASource** is associated with not more than one **DIAObject** or **SSObject**. The algorithm will take into account the parallax and proper or Keplerian motions, as well as the errors in estimated positions of **DIAObject**, **SSObject**, and **DIASource** to find the maximally likely match. Multiple **DIASources** in the same visit will not be matched to the same **DIAObject**.

<sup>14</sup>Level 2 database is a database resulting from annual data release processing.

<sup>15</sup>E.g., a few arcseconds.

8. An alert is issued that includes: the name of the Level 1 database, the timestamp of when this database has been queried to issue this alert, the **DIASource** ID, the **DIAObject** ID<sup>16</sup>, name of the Level 2 database and the IDs of nearby **Objects**, and the associated science content (centroid, fluxes, low-order lightcurve moments, periods, etc.), *including the full light curves*. See Section 2.4 for a more complete enumeration.
9. For all **DIAObjects** overlapping the field of view, to which a **DIASource** from this visit has not been associated, perform forced photometry (point source photometry only). Store those measurements as appropriately flagged **DIASources**<sup>17</sup>. No alerts will be issued for these **DIASources**.
10. Precoverage PSF forced photometry is performed on any difference image overlapping the position of new **DIAObjects** taken within the past 30 days, and added to the database within 24 hours. No additional alerts are issued with the precovery photometry.

### 2.2.2 Solar System Object Processing

The following will occur during normal Solar System object processing (in daytime<sup>18</sup> after a night of observing):

1. The orbits/physical properties of **SSObjects** that were re-observed on the previous night are recomputed. Updated data are entered to the **SSObjects** table.
2. All **DIASources** detected on the previous night, that have not been matched with high probability to a known **Object**, **SSObject**, or an

---

<sup>16</sup>We guarantee that a receiver will always be able to regenerate the alert contents at any later date using the included timestamps and metadata (IDs and database names).

<sup>17</sup>For the purposes of this document, we're treating the **DIASources** generated by precovery measurements to be the same as **DIASources** detected in difference images (but flagged appropriately). In the logical schema, these may be separated into two different tables.

<sup>18</sup>Note that there *is no guarantee on when daytime Solar System processing must finish*, just that, averaged over some reasonable timescale (eg., a month), a night's worth of observing is processed within 24 hours. Nights rich in moving objects may take longer to process, while nights with less will finish more quickly. In other words, the requirement is on *throughput*, not latency.

artifact, are analyzed for potential pairs, forming *tracklets*.

3. The collection of tracklets collected over the past 30 days is analyzed for those *tracks* consistent with being on the same Keplerian orbit around the Sun.
4. For those that are, an orbit is fitted and a new **SSObject** table entry created. **DIASource** records are updated to point to the new **SSObject** record. **DIAObjects** “orphaned” by this unlinking are deleted.<sup>19</sup>.
5. Preccovery linking is attempted for all **SSObjects** whose orbits were updated in this process. Where successful, **SSObjects** (orbits) are updated as needed.

### 2.3 The Level 1 database

The described alert processing design presupposes the existence of an Level 1 database that contains the objects and sources observed on difference images since the beginning of the survey. At the very least<sup>20</sup>, this database will have tables of **DIASources**, **DIAObjects**, and **SSObjects**. They are populated in the course of difference image and Solar System object processing<sup>21</sup>. As these get updated and added to, their updated contents becomes visible (queryable) immediately<sup>22</sup>.

Note that *this database is only loosely coupled to the Level 2 database*. All of the coupling is through providing positional matches between the **DIAObjects** table in the Level 1 database and the **Objects** in the Level 2 database database. There is no direct **DIASource**-to-**Object** match. The adopted data model emphasizes that *having a **DIASource** be positionally coincident with an **Object** does not imply it is physically related to it*. Absent other information, the least presumptuous data model relationship is one of *positional association*, not *physical identity*.

---

<sup>19</sup>Some **DIAObjects** may only be left with forced photometry measurements at their location (since all **DIAObjects** are force-photometered on previous and subsequent visits); these will be kept but flagged as such.

<sup>20</sup>It will also contain exposure and visit metadata, MOPS-specific tables, etc. These are either standard/uncontroversial, or implementation-dependent, irrelevant for science, and therefore not discussed here.

<sup>21</sup>The latter is also colloquially known as *DayMOPS*.

<sup>22</sup>No later than the moment of issuance of any event alert that may refer to it.



This may seem odd at first: for example, in a simple case of a variable star, matching individual **DIASources** to **Objects** is exactly what an astronomer would want. That approach, however, fails in the following scenarios:

- *A supernova in a galaxy.* The matched object in the **Object** table will be the galaxy, which is a distinct astrophysical object. We want to keep the information related to the supernova (e.g., colors, the light curve) separate from those measurements for the galaxy.
- *An asteroid occulting a star.* If associated with the star on first apparition, the association would need to be dissolved when the source is recognized as an asteroid (perhaps even as early as a day later).
- *A supernova on top of a pair of blended galaxies.* It is not clear in general to which galaxy this **DIASource** would belong. That in itself is a research question.

**DIASource-to-Object** matches can still be emulated via a three-step link (**DIASource-DIAObject-Object**). For ease of use, views or pre-built table with these will be offered to end-users.

In the sections to follow, we present the *conceptual schemas* for the most important Level 1 database tables. These convey *what* data will be recorded in each table, rather than the details of *how*. For example, columns whose type is an array (eg., **radec**) may be expanded to one table column per element of the array (eg., **ra**, **decl**) once this schema is translated to SQL. Secondly, the tables to be presented are normalized (i.e., contain no redundant information). For example, since the band of observation can be found by joining a **DIASource** table to the table with exposure metadata, there's no column for 'band' in the **DIASource** table. In the as-built database, the views presented to the users will be appropriately denormalized for ease of use.

### 2.3.1 DIASource Table

This is a table of sources detected at  $SNR \geq 5$  on difference images (**DIASources**). On average, we expect  $\sim 2000$  **DIASources** per visit ( $\sim 2$ M per night; 20,000 per deg<sup>2</sup> per hour).

Some  $SNR \geq 5$  sources will not be caused by observed astrophysical phenomena, but by artifacts (bad columns, diffraction spikes, etc.). The

difference image analysis software will attempt to identify and flag these as such.

Unless noted otherwise, all **DIASource** quantities (fluxes, centroids, etc.) are measured on the difference image.

Table 1: DIASource Table

Name	Type	Unit	Description
diaSourceId	uint128		Unique source identifier
ccdVisitId	uint64		Id. of CCD and visit where this source was measured
diaObjectId	uint128		Id. of the <b>DIAObject</b> this source was associated with <sup>23</sup>
ssObjectId	uint64		Id. of the <b>SSObject</b> this source has been linked to <sup>24</sup>
midPointTai	double	time	Time of mid-exposure for this DIASource.
radec	double[2]	degrees	$(\alpha, \delta)$ <sup>25</sup>
radecCov	float[3]	various	<b>radec</b> covariance matrix
xy	float[2]	pixels	Column and row of the centroid.
xyCov	float[3]	various	Centroid covariance matrix
SNR	float		The signal-to-noise ratio at which this source was detected in the difference image. <sup>26</sup>
psFlux	float	nmgy <sup>27</sup>	Calibrated flux for point source model. Note this actually measures the flux <i>difference</i> between the template and the science image.

*Continued on next page*

<sup>23</sup>diaObjectId will be NULL if ssObjectId is not NULL

<sup>24</sup>ssObjectId will be NULL if diaObjectId is not NULL

<sup>25</sup>The astrometric reference frame will be chosen closer to start of operations.

<sup>26</sup>This is not necessarily the same as psFlux/psFluxSigma, as the flux measurement algorithm may be more accurate than the detection algorithm.

<sup>27</sup>A “maggie”, as introduced by SDSS, is a linear measure of flux; one maggie has an

Table 1: DIASource Table

Name	Type	Unit	Description
psFluxSigma	float	nmgy	Estimated uncertainty of <b>psFlux</b> .
psLnL	float		Natural <i>log</i> likelihood of the observed data given the point source model.
trailFlux	float	nmgy	Calibrated flux for a trailed source model <sup>28,29</sup> . Note this actually measures the flux <i>difference</i> between the template and the science image.
trailLength	float	arcsec	Maximum likelihood fit of trail length <sup>30,31</sup> .
trailAngle	float	degrees	Maximum likelihood fit of the angle between the meridian through the centroid and the trail direction (bearing).
trailLnL	float		Natural <i>log</i> likelihood of the observed data given the trailed source model.

*Continued on next page*

AB magnitude of 0. “nmgy” is short for a nanomaggie. Flux of 0.063 nmgy corresponds to a 24.5<sup>th</sup> magnitude star. See §2.3.7 for details.

<sup>28</sup>A *Trailed Source Model* attempts to fit a (PSF-convolved) model of a point source that was trailed by a certain amount in some direction (taking into account the two-snap nature of the visit, which may lead to a dip in flux around the mid-point of the trail). Roughly, it’s a fit to a PSF-convolved line. The primary use case is to characterize fast-moving Solar System objects.

<sup>29</sup>This model does not fit for the *direction* of motion; to recover it, we would need to fit the model to separately to individual snaps of a visit. This adds to system complexity, and is not clearly justified by increased MOPS performance given the added information.

<sup>30</sup>Note that we’ll likely measure trailRow and trailCol, and transform to trailLength/trailAngle (or trailRa/trailDec) for storage in the database. A stretch goal is to retain both.

<sup>31</sup>TBD: Do we need a separate trailCentroid? It’s unlikely that we do, but one may wish to prove it.

Table 1: DIASource Table

<b>Name</b>	<b>Type</b>	<b>Unit</b>	<b>Description</b>
trailCov	float[6]	various	Covariance matrix of trailed source model parameters.
fpFlux	float	nmgy	Calibrated flux for point source model measured on the science image centered at the centroid measured on the difference image (forced photometry flux)
fpFluxSigma	float	nmgy	Estimated uncertainty of <b>fpFlux</b> .
fpSky	float	nmgy/asec <sup>2</sup>	Estimated sky background at the position (centroid) of the object.
fpSkySigma	float	nmgy/asec <sup>2</sup>	Estimated uncertainty of <b>fpSky</b> .
moments	float[5]	various	Adaptive first and second moments ( $I_x, I_y, I_{xx}, I_{yy}, I_{xy}$ ), measured on the difference image.
momentsSigma	float[5]	various	Estimated uncertainty for each entry in <b>moments</b> .

*Continued on next page*

Table 1: DIASource Table

Name	Type	Unit	Description
extendedness	float		A measure of extendedness, computed using a combination of available moments and model fluxes or from a likelihood ratio of point/trailed source models (exact algorithm TBD). <i>extendedness</i> = 1 implies a high degree of confidence that the source is extended. <i>extendedness</i> = 0 implies a high degree of confidence that the source is point-like.
flags	bit[64]	bit	Flags

### 2.3.2 DIAObject Table

Table 2: DIAObject Table

Name	Type	Unit	Description
diaObjectId	uint128		Unique identifier
radec	double[2]	degrees	$(\alpha, \delta)$ position of the object at time <b>radecTai</b>
radecCov	float[3]	various	<b>radec</b> covariance matrix
radecTai	double	time	Time at which the object was at a position <b>radec</b> .
pm	float[2]	mas/yr	Proper motion vector <sup>32</sup>
parallax	float	mas	Parallax
pmParallaxCov	float[6]	various	Proper motion - parallax covariances.

*Continued on next page*

<sup>32</sup>High proper-motion or parallax objects will appear as “dipoles” in difference images. Great care will have to be taken not to misidentify these as subtraction artifacts.

Table 2: DIAObject Table

Name	Type	Unit	Description
psFlux	float[ugrizy]	nmgy	Weighted mean point-source model magnitude.
psFluxErr	float[ugrizy]	nmgy	Standard error of <b>psFlux</b>
psFluxSigma	float[ugrizy]	nmgy	Standard deviation of the distribution of <b>psFlux</b> .
fpFlux	float[ugruzy]	nmgy	Weighted mean forced photometry flux.
fpFluxErr	float[ugrizy]	nmgy	Standard error of <b>fpFlux</b>
fpFluxSigma	float[ugrizy]	nmgy	Standard deviation of the distribution of <b>fpFlux</b> .
lsPeriod	float[ugrizy]	day	Period (the coordinate of the highest peak in Lomb-Scargle periodogram)
lsSigma	float[ugrizy]	day	Width of the peak at <b>lsPeriod</b> .
lsPower	float[ugrizy]		Power associated with <b>lsPeriod</b> peak.
lcChar	float[6 × M]		Light-curve characterization summary statistics (e.g., 2nd moments, etc.). The exact contents, and an appropriate value of M, are to be determined in consultation with time-domain experts.
nearbyObj	uint128[3]		Closest <b>Objects</b> in Level 2 database.
nearbyObjDist	float[3]	arcsec	Distances to <b>nearbyObj</b> .
flags	bit[64]	bit	Flags

### 2.3.3 SSObject Table

Table 3: SSObject Table

Name	Type	Unit	Description
ssObjectId	uint64		Unique identifier
oe	double[7]	various	Osculating orbital elements at epoch ( $q$ , $e$ , $i$ , $\Omega$ , $\omega$ , $M_0$ , epoch)
oeCov	double[21]	various	Covariance matrix for <b>oe</b>
arc	float	days	Arc of observation.
orbFitLnL	float		Natural log of the likelihood of the orbital elements fit.
nOrbFit	int16		Number of observations used in the fit.
MOID	float[2]	AU	Minimum orbit intersection distances <sup>33</sup>
moidLon	double[2]	degrees	MOID longitudes.
H	float[6]	mag	Mean absolute magnitude, per band.
G	float[6]	mag	Fitted slope parameter, per band <sup>34</sup>
hErr	float[6]	mag	Uncertainty in estimate of H
gErr	float[6]	mag	Uncertainty in estimate of G
flags	bit[64]	bit	Flags

The LSST database will provide functions to compute the phase (Sun-Asteroid-Earth) angle  $\alpha$  for every observation, as well as the reduced ( $H(\alpha)$ ) and absolute ( $H$ ) asteroid magnitudes.

<sup>33</sup><http://www2.lowell.edu/users/elgb/moid.html>

<sup>34</sup>The slope parameter for the large majority of asteroids will not be well constrained until later in the survey. We may decide not to fit for it at all over the first few DRs, and add it later in Operations. Alternatively, we may fit it with a strong prior.

### 2.3.4 Difference Images

The users will be able to download any difference image (or a select subset thereof, to the individual CCD level) within 300 seconds of the end of observation.

The images will remain accessible on low-latency storage (disk) for at least 30 days, after which they will be accessible from the mass storage (tape). The primary difference is in latency of access: while the time from request to download for an image on disk will be on order of seconds, the time to retrieve an image from tape may be minutes to hours.

### 2.3.5 Image Differencing Templates

Templates for difference image analysis will be created by co-adding 6-months to a year long groups of visits. The co-addition process will take care to remove any transients or fast moving objects (eg., asteroids) from the templates.

The input images may be further grouped by airmass and/or seeing<sup>35</sup>. Therefore, at DR11, we will be creating 11 groups templates: two for the first year of the survey (DR1 and DR2), and then one using imaging from each subsequent year.

Difference image analysis will use the appropriate template given the time of observation, airmass, and seeing.

### 2.3.6 Statistical Considerations

Unless noted otherwise, maximum likelihood values will be quoted for all fitted parameters (measurements). Together with covariances, these will allow the end-user to apply whatever prior they deem appropriate when computing posteriors<sup>36</sup>.

We employ a convention where estimates of standard errors have the suffix **Err**, while the estimates of inherent widths of distribution (or functions in general) have the suffix **Sigma**<sup>37</sup>. The latter are defined as the square roots of the second moment of the quantity at hand.

---

<sup>35</sup>The number and optimal parameters for airmass/seeing bins will be determined in Commissioning.

<sup>36</sup>There's a tacit assumption that a Gaussian is a reasonably good description of the likelihood surface around the ML peak.

<sup>37</sup>Given  $N$  measurements, standard errors scale as  $N^{-1/2}$ , while widths remain constant.



For fluxes, we recognize that a substantial fraction of astronomers will just want the posteriors marginalized over all other parameters, trusting the LSST experts to select an appropriate prior<sup>38</sup>. For example, this is nearly always the case when constructing color-color or color-magnitude diagrams. We will support these use cases by providing additional pre-computed columns, taking care to name them accordingly so as to minimize incorrect accidental usage. For example, a column named `gFlux` may be the expectation value of the g-band flux, while `gFluxML` may be the maximum likelihood value.

### 2.3.7 Fluxes and Magnitudes

Because flux measurements on difference images are performed against a template, the measured flux of a source on the difference image can be negative. The flux can also go negative for faint sources in the presence of noise. Negative fluxes cannot be stored as (Pogson) magnitudes (log of a negative number is undefined). We therefore store fluxes rather than magnitudes, in database tables.

We quote fluxes in units of “maggie”. A maggie, as introduced by SDSS, is a linear measure of flux. An object with flux of one maggie (integrated over the bandpass) has an AB magnitude of 0:

$$m_{AB} = -2.5 \log_{10}(f/\text{maggie}) \quad (1)$$

We chose to use maggies (as opposed to Jansky) to allow the user to differentiate between two different sources of calibration error: error in relative calibration of the survey, and error in absolute calibration (the knowledge of absolute flux of photometric standards).

We realize that the large majority of users will want to work with magnitudes. For convenience, we plan to provide columns with (Pogson) magnitudes<sup>39</sup>, where values with negative flux will evaluate to `NULL`. Similarly, we will provide columns with flux expressed in Jy (and its error estimates).

### 2.3.8 Precovery

When a new `DIASource` is detected, it’s useful to perform (PSF) forced photometry at the location of the new source on images taken prior to discovery,

---

<sup>38</sup>It’s likely that most cases will require just the expectation value alone.

<sup>39</sup>These will most likely be implemented as “virtual” or “computed” columns

colloquially know as “*precovery*”<sup>40</sup>. Doing precovery in real time over all previously taken visits is too I/O intensive to be feasible. We therefore plan the following:

1. For all newly discovered objects, perform precovery PSF forced photometry on visits taken over the previous 30 days<sup>41</sup>.
2. Make available a “precovery service” to request precovery for a limited number of **DIASources** across all previous visits, and make it available within 24 hours of the request. Web interface and machine-accessible APIs will be provided.

The former should satisfy the most common use cases (e.g., SNe), while the latter will provide an opportunity for more extensive immediate precovery of targets of special interest.

### 2.3.9 Annual Reprocessings

In what we’ve described so far, the Level 1 database is continually being added to as new images are taken and **DIASources** identified. Every time a new **DIASource** is associated to an existing **DIAObject**, the **DIAObject** record is updated to incorporate new information brought in by the **DIASource**. Once discovered and measured, the **DIASources** would never be re-discovered and re-measured at the pixel level.

This is not optimal. Newer versions of LSST pipelines are likely to improve detection and measurements on older data. Also, PSF forced photometry should be performed on the position of the **DIAObject** on all pre-discovery images.

We therefore plan to reprocess all image differencing-derived data (the Level 1 database), at the same time as we perform the annual Level 2 data release productions. This will include all images taken since the start of observation, to the time when the DR production begins. The reprocessed images will be processed with a single version of the image differencing and measurement software, resulting in a consistent data set.

---

<sup>40</sup>When Solar System objects are concerned, precovery has a slightly different meaning: predicting the position of a newly discovered **SSObject** on previous images, and associating with it **DIASources** consistent with its predicted position.

<sup>41</sup>We will be maintaining a cache of 30 days of processed images to support this feature.

As reprocessing is expected to take on order of  $\sim 9$  months, more imaging will be acquired in the meantime. These data will be reprocessed as well, and added to the new Level 1 database generated by the data release processing. The reprocessed database will thus “catch up” with the Level 1 database currently in use, possibly in a few steps. Once it does, the existing Level 1 database will be replaced with the new one, and all future alerts will refer to the reprocessed Level 1 database. Alerts for new sources “discovered” during data release processing and/or the catch-up process will *not* be issued.

Note that Level 1 database reprocessing and switch will have *significant* side-effects on downstream users. For example, all **DIASource** and **DIAObject** IDs will change in general. Some **DIASources** and **DIAObjects** will disappear (e.g., if they’re image subtraction artifacts that the improved software was now able to recognize as such). New ones may appear. The **DIASource/DIAObject/Objects** associations will change as well.

While the annual database switches will undoubtedly cause technical inconvenience (eg., a **DIASource** detected at some position and associated to one **DIAObject** ID on day  $T - 1$ , will now be associated to a different **DIAObject** ID on day  $T + 0$ ), the resulting database will be a more accurate description of the astrophysics that the survey is seeing (eg., the association on day  $T + 0$  is the correct one; the associations on  $T - 1$  and previous days were actually made to an artifact that skewed the **DIAObject** summary of measurements).

To ease the transition, third parties (VO Event brokers) may choose to provide positional-crossmatching to older versions of the Level 1 database. A set of best practices will be developed to minimize the disruptions caused by the switches (e.g., when writing event-broker queries, filter on position, not on **DIAObject** ID, if possible, etc.). A Level 1 database distribution service, allowing for bulk downloads of the reprocessed Level 1 database, will need to be established to support the brokers who will use it locally to perform more advanced brokering<sup>42</sup>.

Older versions of the Level 1 database will be archived following the same rules as for the Level 2 databases. DR1, the most recent DR, and the penultimate data release will be kept on disk and loaded into the database. Others will be archived to tape and available as bulk downloads.

---

<sup>42</sup>A “bulk-download” database distribution service will be provided for the Level 2 databases as well, to enable end-users to establish and run local mirrors (partial or full).

### 2.3.10 Repeatability of Queries

We require that queries executed at a known point in time against some version of the Level 1 database be repeatable at a later date. The exact implementation of this requirement is under consideration by the DM database team.

One possibility may be to make the key tables (nearly) append-only, with each row having two timestamps – `createdTai` and `deletedTai`, so that queries may be limited through a `WHERE` clause:

```
SELECT * FROM DIASource WHERE 'YYYY-MM-DD-HH-mm-SS' BETWEEN
createdTAI and deletedTAI
```

or, more generally:

```
SELECT * FROM DIASource WHERE ''data is valid as of YYYY-MM-DD''
```

A perhaps less error-prone alternative, if technically feasible, may be to provide multiple virtual databases that the user would access as:

```
CONNECT lsst-dr5-yyyy-mm-dd
SELECT * FROM DIASource
```

The latter method would probably be limited to nightly granularity, unless there's a mechanism to create virtual databases/views on-demand.

### 2.3.11 Uniqueness of IDs across database versions

To reduce the likelihood for confusion, all `Source`, `Object`, `DIASource`, and `DIAObject` IDs shall be unique across database versions. For example, DR4 and DR5 reprocessings will share no identical IDs.

Note, however, that exposure and visit IDs will remain the same across releases.

## 2.4 Alerts to DIASources

### 2.4.1 Information Contained in Each Alert

For each detected `DIASource`, LSST will emit an “Event Alert” within 60 seconds of the end of exposure. These alerts will be issued in `VOEvent` format<sup>43</sup>, and should be readable by `VOEvent`-compliant clients.

---

<sup>43</sup>Or some other format that is broadly accepted and used by the community at the start of LSST commissioning.

Each alert (a `VOEvent` packet) will at least include the following:

- Level 1 database id (example: DR5-Level1)
- alertTimestamp (A timestamp that can be used to execute a query against the Level 1 database as it existed when this alert was issued)
- Science Data:
  - The `DIASource` record that triggered the alert
  - The entire `DIAObject` (or `SSObject`) record
  - All previous `DIASource` records
- $30 \times 30$  pixel cut-out of the difference image (10 bytes/pixel, FITS MEF)
- $30 \times 30$  pixel cut-out of the template image (10 bytes/pixel, FITS MEF)

#### 2.4.2 Receiving and Filtering the Alerts

Alerts will be transmitted in `VOEvent` format, using standard IVOA protocols (e.g., `VOEvent Transport Protocol`; VTP). As a very high rate of alerts is expected, approaching  $\sim 2$  million per night, we plan for public `VOEvent Event Brokers`<sup>44</sup> to be the primary end-points of LSST’s VTP streams. End-users will use these brokers to classify and filter events on the stream for those fitting their science goals. End-users will *not* be able to subscribe to full, unfiltered, alert streams coming directly from LSST<sup>45</sup>.

For the end-users, LSST will provide a basic, limited capacity, alert filtering service. This service will run at the LSST archive center (at NCSA). It will let astronomers create simple filters that limit what alerts are ultimately forwarded to them<sup>46</sup>. These *user defined filters* will be possible to specify

---

<sup>44</sup>These brokers are envisioned to be operated as a public service by third parties who will have signed MOUs with LSST. An example may be the VAO or its successors.

<sup>45</sup>This is due to finite network bandwidth available: for example, a 100 end-users subscribing to a  $\sim 100$  Mbps stream (the peak full stream data rate at end of the first year of operations) would require 10Gbps WAN connection from the archive center, just to serve the alerts.

<sup>46</sup>More specifically, to their VTP clients. Typically, a user will use the Science User Interface (the web portal to LSST archive center) to set up the filters, and use their VTP client to receive the filtered `VOEvent` stream

using an SQL-like declarative language, or short snippets of (likely Python) code. For example, here’s what a filter may look like:

```
# Keep only never-before-seen events within two
# effective radii of a galaxy. This is for illustration
# only; the exact methods/members/APIs may change.

def filter(alert):
    if len(alert.sources) > 1:
        return False
    nn = alert.diaobject.nearest_neighbors[0]
    if not nn.flags.GALAXY:
        return False
    return nn.dist < 2. * nn.Re
```

We emphasize that this LSST-provided capability will be limited, and is *not* intended to satisfy the wide variety of use cases that a full-fledged public Event Broker could. For example, we do not plan to provide any classification (eg., “is the light curve consistent with an RR Lyra?”, or “a Type Ia SN?”). No information beyond what is contained in the `VOEvent` packet will be available to user-defined filters (eg., cross-matches with other catalogs). The complexity and run time of user defined filters will be limited by available resources. Execution latency will not be guaranteed. The number of `VOEvents` transmitted to each user per user will be limited as well (eg., at least up to  $\sim 20$  per visit per user, dynamically throttled depending on load). Finally, the total number of simultaneous subscribers is likely to be limited – in case of overwhelming interest, a TAC-like proposal process may be instituted.

## 2.5 Open Issues to be Closed before Baselining this Document

What follows is a (non-exhaustive) list of issues, technical and scientific, that are still being discussed and where changes are likely. Input on any of these will be appreciated.

- *What light-curve metric should we compute and provide with alerts?* We strive to compute general purpose metrics which will facilitate classification. We have not baselined any yet.

- *Should we measure on individual snaps (or their difference)?* Is there a demonstrable science case requiring immediate followup that would be triggered by the flux change over a  $\sim 15$  second period? Is it technically feasible?
- *Should we choose `nearbyObjs` differently?* One proposal is to find the brightest `Object` within  $XX$  arcsec (with  $XX \sim 10$ arcsec), and the total number of `Objects` within  $XX$  arcsec.
- *When should we (if ever) stop performing forced photometry on positions of `DIAObjects`?* Depending on the rate of false positives, unidentified artifacts, or unrecognized Solar System objects, the number of forced measurements may dramatically grow over time.
- *Can we, should we, and how will we measure proper motions on difference images?* This is a non-trivial task (need to distinguish between dipoles that are artifacts, and those due to proper motions), without a clear science driver (since high proper motion stars will be discoverable using Level 2 catalogs).
- *Is Level 1 database required to be relational?* A no-SQL solution may be more appropriate given the followup-driven use cases. Even if it is relational, the Level 1 database will *not* be sized or architected to perform well on large or complex queries (eg. complex joins, full table scans, etc.).
- *Can users query the Level 1 database for all `DIASources` next to an `Object`?* Is this technically feasible?
- *Do we have to, and can we, use 128 bit integers for IDs?* If 64 bit integers are provably sufficient, they will take up less space and be better supported (technologically).