





PCW2022: LINCC Data to Software to Science

Andy Connolly (UW) Rachel Mandelbaum (CMU) Jeremy Kubica (CMU)

White paper: <u>https://arxiv.org/abs/2208.02781</u>

LINCC Data to Software to Science Agenda

- Outcomes from the DSS meeting (Rachel Mandelbaum and Andy Connolly)
- **Current and planned RSP functionality for time series data** (Leanne Guy and Gregory Dubois Feldsman)
- How do we analyze time series data with the RSP? (Eric Bellm, Neven Caplar)
- How do we make time series data accessible to researchers? (Mario Juric, Colin Slater)
- LINCC summary (Jeno Sokoloski)



The LINCC Frameworks Project

LSST Interdisciplinary Network For Collaboration And Computing

A collaboration between UW, CMU, LSSTC, U Pitt, and NOIRLab to build software systems for key LSST science

Pls: Andy Connolly (UW), Rachel Mandelbaum (CMU) Director of Engineering: Jeremy Kubica (CMU)

Science software infrastructure: combining user algorithms & code, astro packages, and industry tools to build scalable science analysis packages

Additional LINCC faculty here at the PCW: Mario Juric (UW), Michael Wood Vasey (Pitt)

LSST Science Pipelines



Science Platform Research



Inference





Algorithms

New LINCC Frameworks Team Members

Software Engineering Team

- Jeremy Kubica (at PCW)
- Carl Christofferson (TL: UW)
- Max West
- Doug Branton
- Drew Oldag
- Emmanuel Sarpong
- 4 more to come

Project Scientists

- Colin Chandler (at PCW)
- Neven Caplar (at PCW)
- Sam Wyatt
- Alex Malz (at PCW remotely)
- 1 more at CMU, to be hired
- 2 more to come from the University of Pittsburgh

Workshop: From Data to Software to Science with the Rubin Observatory LSST



Workshop goals:

- 1. Enabling *interactive development* of exciting scientific use cases for early LSST data, and identifying the common computational/technical challenges and enabling technologies associated with them.
- 2. Promoting the development of a broad and inclusive community of researchers engaged with LINCC Frameworks.

Program design, plenary talk content, and communication channels for the meeting were developed with both goals in mind.

https://indico.flatironinstitute.org/event/2777/

Science use cases

Divided the science into 7 research areas (not a 1:1 mapping to the LSST Science Collaborations)

- Solar System Science: 6 cases (active asteroids, TNOs)
- Local Universe Static Science: 5 cases (IMF, accreted stellar pops, dwarf gals)
- Local Universe Variable and Transient Science: 9 cases (YSO, microlensing)
- Extragalactic Static Science: 7 cases (morphologies, extinction, LSB dwarfs)
- Extragalactic Variable Science: 8 cases (AGN, lensing)
- Extragalactic Transient Science: 7 cases (SNe, TDEs, classification)
- Cosmology: 6 cases (weak lensing, SNe classification, spectroscopic followup)

~50 use cases for science in the first 2 years of Rubin

	Cross- matching	Photo-z	Selection functions	Time series	Image reprocessing	Image analysis
Cosmology	$\checkmark\checkmark$	$\checkmark\checkmark$	$\checkmark\checkmark$	$\checkmark\checkmark$	\checkmark	\checkmark
Extragalactic static	$\checkmark\checkmark$	$\checkmark\checkmark$	$\checkmark\checkmark$		$\checkmark\checkmark$	\checkmark
Extragalactic transient	$\checkmark\checkmark$	$\checkmark\checkmark$	\checkmark	$\checkmark\checkmark$	\checkmark	\checkmark
Extragalactic variable	$\checkmark\checkmark$	\checkmark	\checkmark	$\checkmark\checkmark$	\checkmark	\checkmark
Local Universe transient & variable	$\checkmark\checkmark$		\checkmark	\checkmark		
Local Universe static	11		$\checkmark\checkmark$		\checkmark	\checkmark
Solar system	\checkmark		$\checkmark\checkmark$	$\sqrt{}$	\checkmark	\checkmark

Table 1. Table highlighting the connection between scientific and technical areas discussed at the workshop. Rows are science areas while columns are for infrastructure capabilities. A double checkmark $(\checkmark \checkmark)$ signifies that some infrastructure capability is essential to enable a particular scientific area, while a single checkmark (\checkmark) signifies that the infrastructure capability would enhance or expand scientific discovery within that area but is not necessary to enable all of it.

Common technical areas identified at the meeting

1. Scalable Cross-matching: real-time (low-latency) positional matching of ~10k sources to ~10 catalogs of ~1Bn sources; offline/batch match and join of ~1Bn sources to catalogs of ~1Bn sources.

2. Photometric redshifts: run and update photo-z's tailored to specific science cases; outputting PDFs for error estimates (~10TB for LSST data); run in parallel

3. Selection function determination: build on DM selection function capabilities; extend to broad science cases (scalar and vector selection functions)

4. Scalable job execution system: run time series, image analysis, classification, model fitting at an LSST scale ~1Bn sources in parallel

Common technical areas identified at the meeting

5. Sky image access and reprocessing at scale: reprocessing of subsets of images (cutouts and full-focal plane data); requires scalable data access services, processing infrastructure, and processing software (built from DM software)

6. Object image access and analysis at scale: processing individual (object-level) images (e.g. deblending, classification); requires scalable image cutout service of arbitrary size; ability to link results to archival data; run in parallel

7. Time series analysis support infrastructure: extract features and classify the captured time-series; enable parametric and model fitting; enable anomaly detection; run in parallel; store, link, and update outputs

We want your feedback on the white paper! (https://arxiv.org/abs/2208.02781)

Does the whitepaper miss any high priority technical cases?

What gaps do you see or functionality that we should focus on?

Are you already working on any of these technical cases and infrastructure?

Are you looking to collaborate on any of the use cases?

Are there starter projects (1-3 months) that would enable science for you today?

The rest of the session will be devoted to more in-depth discussion of time series analysis

Current and planned RSP functionality for time series Leanne Guy and Gregory Dubois Felsmann



A set of integrated web applications & services deployed at Data Access Centers through which the scientific community will access, visualize, subset and perform next-to-the-data analysis of Rubin Data products.





- Enable peta-scale analysis of LSST data
- Exploratory analysis via browsing & visualisation
- Enable discovery 'bring the analysis to the data'
- Supports User-Generated product creation
- Integration with extant archives via IVOA protocols
- Collaborative working environment
- Provision of backend computation & analysis resources



Rubin Science Platform – Three Aspect Design

Portal Aspect

Exploratory analysis and visualization of the LSST archive

Notebook Aspect

In-depth 'next-to-data' analysis and creation of added-value data products

API Aspect

Remote access to the LSST archive via Virtual Observatory interfaces





DPO is the first of three planned data previews between now and Operations.

Rubin's DPO Goals

- enable the community to prepare for early LSST science with the RSP
- test integration of the LSST science pipelines and the RSP
- use feedback on data products and RSP functionality to inform future development

DPO Data Set

- simulated LSST-like images and catalogs from the DESC's Data Challenge 2 (DC2)
- future DP data sets will be based on LSST commissioning data from Rubin Observatory

DPO Timeline

- DP0.1, June 2021: DC2 as processed by the DESC available in the RSP
- DP0.2, June 2022: DC2 as reprocessed by Rubin Data Production available in RSP



DPO: Single time series analysis



ADQL TAP query joining CcdVisit & ForcedSource tables and selecting a single Object

SELECT src.ccdVisitId, src.band, visinfo.expMidptMJD scisql_nanojanskyToAbMag(psfFlux) as psfMag, FROM dp02_dc2_catalogs.ForcedSource as src JOIN dp02_dc2_catalogs.CcdVisit as visinfo ON visinfo.ccdVisitId = src.ccdVisitId WHERE src.objectId = 1651589610221899038

	ccdVisitId	band	psfMag	expMidptMJD
29	2334102	u	20.119284	59583.120963
23	5882102	У	18.420364	59588.091815
313	7999130	z	18.357822	59591.081810
81	8030161	z	18.376561	59591.097111
323	12467085	У	18.321208	59597.089947













ADQL TAP query joining DiaObject and DiaSource tables and applying selection criteria

1. g-band measurements only

2. sigma flux/flux > 0.25 -- the scatter in measured fluxes is larger than 25% relative to the mean

3. sigma_flux/flux < 1.25 -- the scatter in measured fluxes is no larger than 125% relative to the mean

4. 18 < gmag < 23 -- mean g magnitude between 18-23

5. gPSFluxNdata > 30 -- at least 30 observations in g band

6. gPSFluxStetsonJ > 20 -- StetsonJ index greater than 20

7. within 5 degrees of our chosen RA, Dec position



"From Data to Software to Science" session | Rubin Observatory PCW | 08-12 August 2021



Time Series in the RSP

Gregory Dubois-Felsmann, Caltech/IPAC **RSP** Product Owner

Rubin Observatory PCW, 11 August 2022













At the most basic level, there isn't a specific "time series" entity in the DPDD data model.

Instead, time series are generally represented by a combination of information from:

- A *single* row in a summary-over-time table: Object, DIAObject, SSObject
- A corresponding *set* of rows in a row-per-epoch table: ForcedSource, Source, DIASource, SSSource

Matching elements are linkable by the use of object IDs as foreign keys in JOIN operations.



DIAObject to DIASources link on DP0.2

dp02_c	dc2_catalog	s.DiaObject ×					K	◀ 1 of 3 ▶ ▶	(1 - 10					
	diaO /	bjectId long	ra (deg) <i>double</i>	decl (deg) double	nDiaSources	radecTai double	gPSFluxMean double	gPSFluxMeanE double	rr					
A					•									
	1736269	9597746659454	60.43118	-35.0911744	232	61394.1793962	249.1326373	89.52500	095					
	1736269	9597746659447	60.39539	-35.1199186	178	61394.1793962	2149.6076468	76.52492	222					
	1736269	9597746659442	60.37182	-35.0977601	152	61394.1615222	-2095.7588751	110.8328	194					
	1736269	9597746659594	60.38809	-35.1184075	74	61394.1606082	1757.4983003	65.77340	034					
	1736269	9597746659751	60.37220	-35.1076471	29	61394.1606082	-504.5312189	112.1243	113					
	dp02	2_dc2_catalogs.Dia	Object - data.	× dp02_dc2_cat	alogs.DiaSource	×								
							 	(1 - 74 of 74)						•
		diaObject	(d	diaSourceId	filterName char	midPointTai	psFlux (nJy) double	psFluxErr (nJy) double	apFlux (nJy) <i>double</i>	apFluxErr (nJy) double	snr double	ccdVisitId long	ra (deg) double	decl (deg) <i>double</i>
	- 9		-		_									
		173626959774	46659594	21650450713411730) r	59634.0870892	-162.4165258	493.3907271	-474.4172627	566.1506331	-0.83797	40327107	60.3880926	-35.1184083
		173626959774	46659594	103695629347192947	'r	59839.3345202	782.4670571	491.0154481	-463.4886483	682.508526	-0.6790958	193148161	60.3880885	-35.1184321
		173626959774	46659594	104089161631269026	i r	59840.2634672	129.6234774	481.6366517	-103.3563425	570.0203827	-0.1813204	193881172	60.3880992	-35.1184027
		173626959774	46659594	113880060884156582	! r	59867.2389572	-606.504622	498.3523518	-711.1592372	816.5674336	-0.8709131	212118143	60.3880568	-35.1184211
		173626959774	46659594	114629457515380850) r	59869.1437982	1377.7409913	467.870113	169.5410355	544.6786202	0.311268	213514003	60.3880994	-35.1183974
		173626959774	46659594	115139499377295429) i	59870.2191692	2268.7542599	451.2647375	1791.4775525	771.5089272	2.3220444	214464030	60.3882083	-35.1183562
		173626959774	46659594	115190038794338442	! i	59870.2628022	2547.3883046	499.9141078	114.1604098	817.3599102	0.1396697	214558167	60.3882059	-35.1183591
		173626959774	46659594	117672458487595165	j g	59876.1526762	2310.4245603	358.1040394	836.1425164	392.5975848	2.1297701	219182034	60.3880227	-35.1184146
		173626959774	46659594	118074634930225305	i r	59877.1308142	-463.2969718	465.400646	-1325.8760761	546.9054646	-2.4243245	219931146	60.3881175	-35.118393
		173626959774	46659594	130193935812264085	i r	59914.1503912	-2787.6273709	458.2782084	-1526.8971649	668.821203	-2.2829681	242505103	60.3881838	-35.1184036
		173626959774	46659594	145986031806578912	g	59958.1119582	1715.5003027	379.7710178	460.1891598	387.3741091	1.1879709	271920174	60.3880918	-35.1184118
		173626959774	46659594	146004282733232378	g	59958.1288562	585.2172472	375.220935	19.4731102	386.4123448	0.0503946	271954169	60.3881031	-35.1184134
		17362695977	46659594	220646458704003180	i	60162 3682742	1272 4296824	678 9425943	1674 0205101	1142 1047422	1 465733	410986056	60 3881125	-35 1183818

Vera C. Rubin Observatory | Project & Community Workshop | 8-12 August 2022



Searchable by position as well as object ID

- Searches based on *Object IDs rely on the pipeline software having performed associations on the single-epoch detections (Source, DIASource)
 - ForcedSource, of course, has an inherently "perfect" association the measurements were taken explicitly driven by the parent object
- There are judgement calls inherent in those associations
- You can get an independent take on that by searching the single-epoch tables by position instead of ID
- The tables are/will be indexed to facilitate access either way



Data are not segregated by filter

- The "Object-like" tables are not only summaries over time but also over filter
- The associated "Source-like" tables will generally contain data from multiple filters (a single filter per row)
- Plotting or otherwise analyzing time-series data requires dividing the single-epoch data by filter



Variability metrics and features

- The "Object-like" tables are planned to contain summary data on variability
 - E.g., in the DPDD for Object:

		•
lcPeriodic	float[6×32]	Periodic features extracted from dif-
		ference image-based light-curves us- 🔤
		ing generalized Lomb-Scargle peri- 📄
		odogram [Table 4, 17].
lcNonPeriodic	float[6 \times 20]	Non-periodic features extracted from
		difference image-based light-curves
		[Table 5 17]

• Eric Bellm is speaking about the development of these metrics in this session



Links to images

dp02_dc2_catalogs.DiaObject - data... × dp02 dc2 catalogs.DiaSource ... × 1 of 1) (1 - 74 of 74) diaObjectId diaSourceId filterName midPointTai 🔺 psFlux psFluxErr apFlux apFluxErr snr ccdVisitId decl ra (nJy) double (deg) (deg) (nJy) double (nJy) double (nJy) double lona long char double double long 8 21650450713411730 59634.0870892 -162,4165258 493.3907271 -474.4172627 566.1506331 -0.83797 40327107 60.3880926 -35.1184083 1736269597746659594 15148161 .6790958 60.3880885 -35.1184321 Data Product: ivoa.ObsCore - data... Coverage XQ | 9 💵 🗗 🗆 📽 🖬 🔨 193881172 60.3880992 -35.1184027 1 of 1 (1 - 2 of 2) .1813204 212118143 60.3880568 -35.1184211 .870913 HDU (#1): IMAGE <> 1/3 HDU (#1): IMAGE <> 1/3 5117 -35.1183974 213514003 60.3880994 Primary product (#this) FOV: 23 Primary product (#this) FOV: 23' .3220444 214464030 60.3882083 -35.1183562 $\oplus \bigcirc \oslash \bigcirc$ Links to the template images for DIASources should be available by **DP1**. [2] EO-J2000: 4h02m17.07s, -35d13m12.9s Value: -6.609688 DN Lock by click dp02 dc2 catalogs.DiaObject - data... × dp02 dc2 catalogs.DiaSource - dat... × dp02 dc2 catalogs.CcdVisit - data.l... × ivoa.ObsCore - data isst.cloud/. 1 of 1 (1 - 2 of 2) dataproduct_type lsst filter dataproduct subtype calib level isst band min em max lsst visit isst detector t_exptime t_min (m) (m) (s) double (d) double chai int cha lona char lona 8 lsst.calexp 2 r 5.52e-7 6.91e-7 40327 r sim 1.4 107 30 59634.08691561111 596 image lsst.goodSeeingDiff_differenceExp 3 r 5.52e-7 107 596 image 6.91e-7 40327 r_sim_1.4 30 59634.08691561111

Vera C. Rubin Observatory | Project & Community Workshop | 8-12 August 2022



... and links to image metadata

diaObjectId long diaSourceId long filterName char midPointTai + double psFlux (n)y double psFlux (n)y double psFlux (n)y double psFlux (n)y double psFlux (n)y double psFlux (n)y double snr ccdVisitId double filterName (n)y filterName (n)y snr ccdVisitId double snr ccdVisitId double <th></th> <th></th> <th></th> <th></th> <th>【 ◀ 1 of 1</th> <th>(1 - 74 of 74)</th> <th></th> <th></th> <th></th> <th></th> <th></th> <th></th> <th></th>					【 ◀ 1 of 1	(1 - 74 of 74)							
Image: Control of the second seco	diaObjectId long	diaSourceId long	filterName char	midPointTai	psFlux (nJy) double	psFluxErr (nJy) double	apFlux (nJy) double	apFluxErr (nJy) double	snr double	ccdVisitId long	ra (deg) <i>double</i>	dec (deg doub	l I) Die
1736269597746659594 21650450713411730 r 59634.0870892 -162.4165258 493.3907271 -474.4172627 566 160331 -0.83797 40327107 60 1736269597746659594 103695629347192947 r 59839.3345202 782.4670571 491.0154481 -463.4886483 682.508526 -0.6790958 193148161 60 1736269597746659594 104089161631269026 r 59840.634672 129.6234774 481.6366517 103.3563425 570.0203827 -0.1813204 193881172 60 1736269597746659594 113880060884156582 r 59867.2389572 -606.504622 498.3532618 -711.1592372 816.5674336 -0.8709131 212118143 60 1736269597746659594 114629457515380850 r 59870.2191692 22667542599 451.2647375 1791.4775525 771.508272 2.3220444 214464030 60 1736269597746659594 115139499377295429 i 59870.2191692 22667542599 451.2647375 1791.4775525 771.508272 2.3220444 214464030 60 1736269597746659594 1913492_catalogs.DiaSurce - dat × dp02_ac2_catalogs.CcdVisit - d × <td< td=""><td>▼</td><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td></td><td></td></td<>	▼												
1736269597746659594 103695629347192947 r 59839.3345202 782.4670571 491.0154481 -463.4886482 682.508526 0.6790958 19314816 66 1736269597746659594 104089161631269026 r 59840.2634672 129.6234774 481.6366517 140.3563425 570.0203827 0.1813204 193881172 66 1736269597746659594 113880060884156582 r 59867.2389572 -606.504622 498.3532518 -711.1592372 816.5674336 -0.8709131 212118143 66 1736269597746659594 114629457515380850 r 59869.1437982 1377.740913 467.870113 169.5410355 544.6786202 0.311268 213514003 66 1736269597746659594 115139499377295429 i 59870.2191692 2267.542599 451.2647375 1791.4775525 771.5089272 2.3220444 214464030 66 1072_cc_catalogs.Dia/Object - data × dp02_dc2_catalogs.CotVisit - d × </td <td>1736269597746659594</td> <td>21650450713411730</td> <td>r</td> <td>59634.0870892</td> <td>-162.4165258</td> <td>493.3907271</td> <td>-474.4172627</td> <td>566 1500331</td> <td>-0.83797</td> <td>40327107</td> <td>60.388092</td> <td>6 -35.1</td> <td>18408</td>	1736269597746659594	21650450713411730	r	59634.0870892	-162.4165258	493.3907271	-474.4172627	566 1500331	-0.83797	40327107	60.388092	6 -35.1	18408
17362695977466559594 104089161631269026 r 59840.2634672 129.6234774 481.6366517 1503.3563425 570.0203827 -0.1813204 193881172 66 17362695977466559594 113880060884156582 r 59867.2389572 -606.504622 498.353518 -711.1592372 816.5674336 -0.870913 212118143 66 17362695977466559594 114629457515380850 r 59869.1437982 1377.7409913 467.870113 169.5410355 544.6786202 0.311268 213514003 66 1736269597746659594 115139499377295429 i 59870.2191692 2266.7542599 451.2647375 1791.4775525 771.5089272 2.3220444 214464030 66 1736269597746659594 1902-dc2_catalogs.Cut/sit - d × 59870.2191692 2266.7542599 451.2647375 1791.477552 771.5089272 2.3220444 214464030 66 1 0402-dc2_catalogs.Dua/subject - data ×	1736269597746659594	103695629347192947	r	59839.3345202	782.4670571	491.0154481	-463.4886483	682.508526	0.6790958	193148161	60.388088	-35.1	1843
17362695977466559594 113880060884156582 r 59867.2389572 -606.504622 498.353618 -711.1592372 816.5674336 -0.8709131 212118143 66 17362695977466559594 114629457515380850 r 59869.1437982 1377.7409913 467.870113 169.5410355 544.6786202 0.311268 213514003 66 1736269597746659594 115139499377295429 i 59870.2191692 2266.7542599 451.2647375 1791.4775525 771.5089272 2.3220444 214464030 66 1 01736269597746659594 10902-dc2_catalogs.CcdVisit - d × 467.87013 169.5410355 544.6786202 0.311268 213514003 66 1 01736269597746659594 115139499377295429 i 59870.2191692 2266.7542599 451.2647375 1791.4775525 771.5089272 2.3220444 214464030 66 1 0173626950746659594 dp02-dc2_catalogs.CcdVisit - d × × </td <td>1736269597746659594</td> <td>104089161631269026</td> <td>r</td> <td>59840.2634672</td> <td>129.6234774</td> <td>481.6366517</td> <td>103.3563425</td> <td>570.0203827</td> <td>0.1813204</td> <td>193881172</td> <td>60.388099</td> <td>2 -35.1</td> <td>1840</td>	1736269597746659594	104089161631269026	r	59840.2634672	129.6234774	481.6366517	103.3563425	570.0203827	0.1813204	193881172	60.388099	2 -35.1	1840
17362695977466559594 114629457515380850 r 59869.1437982 1377.7409913 467.870113 169.5410355 544.6786202 0.311268 213514003 66 1736269597746659594 115139499377295429 i 59870.2191692 2266.7542599 451.2647375 1791.4775525 771.5089272 2.3220444 214464030 66 cc_catalogs.DiaObject - dat × dp02_dc2_catalogs.CcdVisit - d × <	1736269597746659594	113880060884156582	r	59867.2389572	-606.504622	498.3522518	-711.1592372	816.5674336	0.8709131	212118143	60.388056	8 -35.1	1842
1736269597746659594 115139499377295429 i 59870.2191692 2266.7542599 451.2647375 1791.4775525 771.5089272 2.3220444 214464030 66 c2_catalogs.DiaObject - data × dp02_dc2_catalogs.CcdVisit - d × dp02_dc2_catalogs.CcdVisit - d × <td< td=""><td>1736269597746659594</td><td>114629457515380850</td><td>r</td><td>59869.1437982</td><td>1377.7409913</td><td>467.870113</td><td>169.5410355</td><td>544.6786202</td><td>0.311268</td><td>213514003</td><td>60.388099</td><td>4 -35.1</td><td>1839</td></td<>	1736269597746659594	114629457515380850	r	59869.1437982	1377.7409913	467.870113	169.5410355	544.6786202	0.311268	213514003	60.388099	4 -35.1	1839
c2_catalogs.DiaObject - data × dp02_dc2_catalogs.CcdVisit - d × I of 1) (1 - 1 of 1) nd ccdVisitId darkTime detector expMidpt expMidpt (5), (5), (5), (5), (5), (5), (5), (7), (7), (7), (7), (7), (7), (7), (7	1736269597746659594	115139499377295429	i	59870.2191692	2268.7542599	451.2647375	1791.4775525	771.5089272	2.3220444	214464030	60.388208	3 -35.1	1835
2_catalogs.DiaObject - data × dp02_dc2_catalogs.DiaSource - dat × dp02_dc2_catalogs.CcdVisit - d × I													
nd ccdVisitId darkTime detector expMidpt expMidpt (1 of 1) (1 - 1 of 1) (3) (3) (5) (5) (5) (6) (6) (6) (7) (1 - 1 of 1)													
nd ccdVisitId darkTime detector expMidpt expMidptMJD expTime (s) cbsStartMJD physical_filter psfSigma seeing (d) (s) cbsStartMJD physical_filter (pixel) (arcsec) (cbsStartMJD cbsStartMJD cbs	2_catalogs.DiaObject - data ×	t dp02_acz_catalogs.DiaS	ource - dat ×	dp02_dc2_catalogs.Cc	dVisit - d ×								
ar long double long char double double char double char float double f	2_catalogs.DiaObject - data >	د dp02_acz_catalogs.DiaS	'ource - dat ×	dp02_dc2_catalogs.Cc	dVisit - d ×	of 1)					9		i
	c2_catalogs.DiaObject - data > nd ccdVisitId darkTim ar long double	dp02_dc2_catalogs.DiaS	ource - dat × expMidpt char	dp02_dc2_catalogs.Cc id d expMidptMJD (d) double	dVisit - d × 1 of 1) (1 - 1 of 1) (1 -	of 1) obsStart <i>char</i>	obsStartMJI double	D physical_filter	psfSigma (pixel) float	seeing (arcsec) double	SkyBg (DN) float	skyNoise (DN) float	i) skr

• Seeing, PSF metrics, background levels, etc.



Fundamentally, tabular data in the RSP will be available in two ways:

- Via ADQL database queries, mediated by an IVOA TAP service
 - This is part of the "API Aspect" of the Science Platform
 - It is accessible from both the Portal Aspect and the Notebook Aspect (as illustrated by DP0.2 tutorials)

SELECT diaObjectId,diaSourceId,filterName,midPointTai,psFlux,psFluxErr,apFlux,apFluxErr,snr,ccdVisitId,ra,decl FROM dp02_dc2_catalogs.DiaSource WHERE (diaObjectId = 1736269597746659594)

- As spatially-sharded Parquet files, accessible via the "Butler" middleware
 - Internally, this is aimed at the Notebook Aspect and at usability in the user batch system
 - Externally, we are working with the IVOA on means to support this via appropriate metadata and API Aspect services mainly for access from beyond the RSP



Example from DP0.2

From the DP0.2 Portal Lightcurve tutorial



Vera C. Rubin Observatory | Project & Community Workshop | 8-12 August 2022

Acronyms & Glossary



Data Access philosophy

- You can think of these, TAP/ADQL access and Parquet-shard access, as "query" and "bulk" analysis interfaces, respectively.
 - However, remember that the Qserv database behind the TAP system is optimized for performing full-table scans efficiently, not just "targeted", indexed searches, so some types of bulk analysis will work just fine as TAP queries.
 - The value of this TAP access is very strongly tied to the usefulness of the object-level metrics and feature extraction that we are able to do, because...
 - It is still not thought to be realistic to do full time-series analysis (e.g., template-based searches, periodograms, etc.) in the database, so...
- A robust system for analyzing the Parquet data in bulk is essential; and
- Your input on appropriate feature analyses to perform would be most valuable.



API Aspect direct support for time series

- In addition to the basic ability to query the tables via TAP, the RSP will include value-added services in the API Aspect to retrieve the Source-like time-series data based on an object-like ID without the user having to construct ADQL
 - In other words, interfaces with URL query parameters like "?OBJECTID=1736269597746659594"
- We will also have services to return image data and metadata based on a source
 These will generally just be trivial convenience wrappers for TAP queries
- Queries on the Object-like and Source-like tables will contain IVOA DataLink annotations that enable client software (e.g., PyVO) to follow these links easily
- The Portal Aspect will provide "one-click" UI support for this
- Prototypes will be deployed in DP0.2 later this year stay tuned for tutorials!



- At present there have been many proposals, but there is no fully worked-out community standard for a "time series" query result that combines summary/feature information with the actual time-dependent photometry
- We are engaged with the IVOA in standards development efforts and will actively
 promote progress in this area
- We have a service architecture in the RSP that will facilitate rapid prototyping, and following a standard if/when it develops
- In the mean time, as we move to the DR1 era we will progressively annotate queries on the time series tables with more metadata, following existing standards



- The service architecture will easily allow for additional, value-added services like periodogram calculation to be linked to light curve data in the API Aspect and then easily discovered from the Portal Aspect.
 - We expect to enhance the environment along these lines during operations.
 - LINCC contributions could be made available to users of all Aspects by wrapping them in this (Python) framework

• Project requirements also call for us to provide a forced-photometry-on-demand service. This will also be readily available through the API Aspect and Portal Aspect, in addition to being straightforward to code using the Science Pipelines stack in the Notebook Aspect.



- We have been communicating with the IVOA on formalizing support for Parquet in a way that preserves interoperability with existing and future IVOA metadata and data models...
- And that facilitates access to very large, wide-area spatially-sharded catalog datasets through existing services like ObsTAP and the Registry.

• Mario Juric will also be talking about this



Computation on time series in the RSP

- We are obligated by project requirements to provide a bulk-data-analysis environment for Rubin/LSST data users
- This was deliberately left vaguely defined in early project design documents because we recognized how rapidly the software and big-data-architecture landscape was changing
- We are only explicitly committed to provide a "conventional batch" service
- However, we recognize the importance of providing a structured environment in the RSP for bulk catalog analysis (e.g., Dask, Spark), and we have successfully prototyped Dask access from the RSP. This is still under investigation and we will have more to say about this at future meetings.



- What gaps do you see in tools that will be available? Which additional functionality would you like to see? What things should go directly into the RSP?
- Are there other systems besides the RSP that you are planning to use for time series analysis?
- * This slide is about systems used for analysis; after subsequent presentations we'll talk about more specific analysis algorithms

How do we analyze time series data with the RSP? Eric Bellm, Neven Caplar



Rubin will pre-compute time series features in both Alert and Data Release Production.

Alert Production: Difference Image (DIAObject) lightcurve features

- computed on 12 months of DIASources during Prompt Processing (< 60 second latency)
- included in alerts & the Prompt Products Database

Data Releases: both difference & direct imaging (Object) features

- computed on all DIAForcedSources and ForcedSources during DRP
- included in Data Release catalogs

The Data Products Definition Document allocates space:

lcPeriodic	float[6×32]	Periodic features extracted from DIA-
		Source light-curves using generalized
		Lomb-Scargle periodogram [Table 4,
		17 <mark>1⁴⁸.</mark>
lcNonPeriodic	float[6×20]	Non-periodic features extracted from
		DIASource light-curves [Table 5, 17].



ls.st/dpdd



In AP, features will be constrained by latency, data, and compute.

High-level 60 second latency requirement ⇒

- 12 months of history = ~80 epochs total (6 filters)
- Limited CPU & memory
- Only a few seconds to compute features!

Still enough data for features useful for query, alert filtering, user classification





Andy Tzanidakis



Rubin is working to develop a feasible and useful feature set for AP.

Provide general-purpose feature set:

- Generic summary statistics
- Basic period estimation
- Transient parameterization
- Charaterize aperiodic variability

DMTN-118 discusses technical considerations & open questions

<u>Is.st/fkr</u> is a work-in-progress draft feature set; see discussion on <u>community.lsst.org</u>



DMTN-118



DRP feature computation is less constrained:

- Computed on forced photometry from entire survey
- More flexibility with computation environment and latency

DRP features will likely evolve from Data Release to Data Release based on community usage and feedback.

Expect we'll start with the AP feature set in DP2-DR1 era.

Good opportunity for discussion/exchange with LINCC.



Specialized science cases

E.g., searches for very short period binaries; changing periods

Computationally intensive tasks Large-scale fitting of template lightcurves MCMC

GPU implementations

Rubin is not trying to develop a general-purpose timeseries library but we are interested in discussions with others who might be



- What are the requirements for a good candidate?
 - Code applicable for LSST science and scale
 - Active developers/groups using the code (I.e., good and readily deployable ideas)
 - Possibility for a tight interaction between scientists and programmers

We have identified 31 existing timeseries codes in these broad areas:

• Explosive transients

- Transient Classifiers
- Lightcurve fitting for SN standardization
- Lensing
 - Microlensing
 - Strong lensing
- AGN
- Periodograms



Examples of ideas by current maintainers

• Explosive transients

- SuperNNova
 - Implement unit and integration tests
 - Major rewrite to make compatible with pytorch updates
 - Optimize running of the code on the alert type data products
- Lensing
 - Lensastronomy
 - JAX and adaptive mesh supported micro-lensing code; JAX for established macro-lensing codes
- AGN
 - EzTao
 - Optimize current JAX implementation, stress testing
 - Mutliband analysis
- Periodograms
 - Astropy implementations
 - Add possibility for mutliband analysis



- 1. What features do you need to measure on time series data?
- 2. What algorithms do you want to run on time series data (clustering, modeling, etc.)? What software packages are you currently using that apply these algorithms to time series data?
- 3. How will you sample lightcurves for analysis (uniform random, based on features, no-sampling, etc.)?

How do we make time series data accessible to researchers for large-scale analyses?

Mario Juric & Colin Slater

С	Tech	nical areas in detail	176
	C .1	Introduction	. 176
	C.2	Cross Matching	. 176
	C.3	Selection Functions	. 181
	C.4	Time Series	. 183
	C.5	Image Reprocessing	. 195
	C.6	Image Analysis	. 200
	C.7	Photometric Redshifts	. 212

Example

Left: a search for "dipper" objects (stars that exhibit dimming episodes, inconsistent with binarity), in 2Bn light curves of ZTF by running a custom filter:

Stetzler et al. (2022)



Analysis at (whole-dataset) scale

• Relational databases are not ideal for this type of work; it's better to operate on files.



- Rubin plans to provide data in the Parquet (<u>https://parquet.apache.org/</u>) file format both to support bulk download and large-scale analyses.
 - Why Parquet: it the de-facto standard format for tabular big-data analyses. Supported by all major analytics frameworks from Spark, Dask, to Pandas and astropy.
 - It is well suited to keeping data in the cloud (i.e., in "object stores").

import pandas as pd
pd.read_parquet('example_pa.parquet', engine='pyarrow')

- The dataset will be partitioned delivered in multiple files, enabling tools to operate in parallel.
 - Think 1,000 cores, each analyzing ~40 files worth of light curves.
 - Or a massively parallel download, for bulk distribution.

How to partition? Historically, we haven't generally given this much thought...

				Oi Alert Archive × +			~
Index of /Gaia/gdr3/gaia_sour	• × +		*	← → C		ů 🌣 🕫 🤩 🖣	🗟 🐮 🗯 🔲 🍘 🗍 Update 🕴
\leftarrow \rightarrow C (\blacktriangle Not Secure cdn.gea.e	sac.esa.int/Gaia/gdr3/gaia_sou	rce/ 🖞 🕸 🥬 🎈 🗟 🐮 🕏 🗍 🍘 Update		•			
Index of Caje/adm2/a	aia courac/			TF ALERT ARCHIVE			DIRAC
muex of /Gala/gur5/g	ala_source/	Index of /sas/dr8/sdss/segue2/ × +		•			
			A A				
/		← → C ■ data.sdss.org//sas/dr8/sdss/segue2/targetAil/301/		What is included?	Known ca	veats	
GaiaSource_000000-003111.csv.gz	05-May-2022			Below you will find compressed tar archives of ZTF event alert	s (observations • The d	lata provided on this site is generate	ad automatically. The files provided
GaiaSource_003112-005263.csv.gz	05-May-2022	Index of /eac/dr8/edee/eague2/	targetAll/301/	detected in image differences). Each tar file contains alerts coll night (LTC-based) with each alert stored in a senarate file in t	ected in the given contai he AVBO format. To conco	in a full, unfiltered, 5-sigma alert str	eam. Depending on your science
GalaSource_005264-006601.csv.gz	05-May-2022	index of /sas/uro/suss/seguez/	algerAll/301/	get you started, we offer a repository with few basic utilities for	reading AVRO- on the	included attributes such as the sig	nal-to-noise ratio or the real-bogus
GaiaSource_007953-010234.csv.gz	05-May-2022			serialized data, as well as an example Jupyter notebook. The s	chema fields are score.		
GaiaSource_010235-012597.csv.gz	05-May-2022	File Name ↓	File Size ↓	described nove.	• Users ZTF D	anterested in un-subtracted archiva Data Releases, accessible at IRSA.	ii photometry should consider the
GaiaSource_012598-014045.csv.gz	05-May-2022	Parent directory/	-	Why this service?	A sub-	set of events obtained through Call	ech time are made public here in
GaiaSource 015370-016240.csv.gz	05-May-2022			We are providing this archive as simple alternative to public ev	ent brokers. Full- "progr	ramid3" tarballs; as of this writing th	ese are additional observations of
GaiaSource_016241-017018.csv.gz	05-May-2022	1000/	-	featured event brokers that provide real-time access to these a	lerts include MARS,	1200 30000.	
GaiaSource_017019-017658.csv.gz	05-May-2022	1006/		Lasair, ANTARES, and ALeRCE.			
GaiaSource_017659-018028.csv.gz	05-May-2022	1009/	•				
GalaSource_018029-018472.csv.gz	05-May-2022	1010/					
GaiaSource 019162-019657.csv.gz	05-May-2022	1011/		Name Search		Last modified	Size
GaiaSource_019658-020091.csv.gz	05-May-2022	1013/					
GaiaSource_020092-020493.csv.gz	05-May-2022	1022/		ztf_public_20220810.tar.gz		10 hours ago	8.3G
GalaSource_020494-020747.csv.gz	05-May-2022	1024/					
GaiaSource 020985-021233.csv.gz	05-May-2022	1024/		ztf_public_20220809.tar.gz		1 day ago	8.3G
GaiaSource_021234-021441.csv.gz	05-May-2022	1033/		-			
GaiaSource_021442-021665.csv.gz	05-May-2022	1035/		1 ztf_public_20220808.tar.gz		2 days ago	1.7G
GalaSource_021966-021919.CSV.gz	05-May-2022	1037/		D atf public 20220207 tor or		2 days ann	100
GaiaSource 022159-022410.csv.gz	05-May-2022	1040/		E zil_public_zozzobov.tal.gz		o uays ago	100
GaiaSource_022411-022698.csv.gz	05-May-2022	1043/		D ztf public 20220806 tar.oz		4 days ago	11G
GaiaSource_022699-022881.csv.gz	05-May-2022	1045/				, ,	
GalaSource_022882=023058.CSV.gz	05-May-2022	1055/		ztf_public_20220805.tar.gz		5 days ago	3.6G
GaiaSource 023265-023450.csv.gz	05-May-2022	1033/					
GaiaSource_023451-023649.csv.gz	05-May-2022	1056/		ztf_public_20220804.tar.gz		5 days ago	6.3G
GaiaSource_023650-023910.csv.gz	05-May-2022	1057/					
GalaSource_023911=024205.csv.gz	05-May-2022	109/		ztf_public_20220803.tar.gz		7 days ago	11G
GaiaSource 024527-025166.csv.gz	05-May-2022	1119/		D // // 00000000		0 4000 0000	170
GaiaSource_025167-025691.csv.gz	05-May-2022	1120/		ztf_public_20220802.tar.gz		8 days ago	176
GaiaSource_025692-026057.csv.gz	05-May-2022	1122/		D ztf public 20220801 tar.oz		9 days ago	3.4G
GalaSource_026058=026390.csv.gz	05-May-2022	1122/		En _public_coccoop nange		o dayo ago	0110
GaiaSource 026649-027106.csv.gz	05-May-2022	1133/		T ztf public 20220731.tar.oz		10 days ago	2.8G
GaiaSource_027107-027517.csv.gz	05-May-2022	1140/	-				
GaiaSource_027518-027832.csv.gz	05-May-2022	1142/		ztf_public_20220730.tar.gz		11 days ago	14G
GalaSource 028077-028338 cev gz	05-May-2022	1231/	-	2012-DEC-19 10:35			
		1233/		2012-Dec-19 10:33			
		1239/		2012-Dec-19 10:33			
		1241/		2012-Dec-19 10:33	but Rubir	n isl (see Gl	PDF's slides -
		125/		2012-Dec-19 10:33		10. (000 01	
		1202/		2012-Dec-10 10:24			those slides)
		1220/	-	2012 Dec 10 10:34			- unese silues
		1329/	-	2012-Dec-19 10:34			
		J 1331/		2012-Dec-19 10·34			

User desiderata

- 1. Keep all data related to the same object together (i.e., the whole time series).
- 2. Partition spatially. Each file should have objects that are close together, enabling spatial queries.
- 3. Ensure similarly-sized partitions (files). When processing in parallel, each file should take roughly same time to process.

A Solution: Hierarchical Healpix Partitioning ("HiPSCat")



Extension of the widely used IVOA HiPS standard.

Some benefits

- 1. Can keep all data related to the same object together.
 - Open implementation question: how? Series of rows? Or arrays?
- 2. Partitioned spatially. Very easy to locate a file an object w. (ra, dec) is stored in. 🖌
- Ensures similarly-sized partitions (files). Partitions are hierarchically split to the next healpix order if they exceed a pre-set size limit.
- 4. Minimal extension of an already well-understood standard (HiPS)
- 5. Can re-uses (significant existing) HiPS infrastructure, including mirroring and caching.
- 6. Can be utilized by (HiPS) image viewers for overplotting of catalog data (including full light curves)
- 7. Enables high-level functionality via HTTP-only interface (no special servers)
- 8. and...







If two or more surveys have catalogs published following this format (ideally on the cloud), **highly parallel, on-the-fly, joining and cross-matching becomes possible**.

Right: the Python that you'd write to cross-match and jointly analyze two large-scale catalogs.



Towards scalable multi-catalog analysis



euclid = spark.table("s3://ipac/euclid/objects.hcat")
lsst = spark.table("s3://noirlab/lsst/objects.hcat")

euclid



Figure 7. A screenshot of the job timeline from the Spark UI when dynamic allocation is enabled. A long-running query is stated, executing with a small number of executors. As the query continues, Spark adds exponentially more executors to the cluster at a user-specified interval until the query completes or the max number of executors is reached. Once the query completes (or its terminated, as shown here), the Spark executors are moved from the cluster.

Above: exponential scale-up of nodes until all data are processed.

Right: weak scaling (nearly perfect).

(Stetzler et al. 2022, using Spark)



Progress

- Learning what others are doing, collecting feedback, possible alternatives, and fellow travelers. This is very much an idea at its initial stages.
 - Been testing these concepts for the past ~3 years (more in Zecevic+ 2019, Stetzler+ 2022)
- To explore it further, looking to:
 - Develop a Python implementation prototype/testbed
 - Creation/query
 - Cross-matching
 - Developing a Dask-based large-scale processing example...
 - ... to be followed by Spark, if successful.
- Goal: A prototyped concept ready for IVOA discussions / Rubin evaluation.

Discussion questions

- 1. How are other groups thinking about the problem of large-scale data/time series storage? What other efforts are out there?
- 2. Tool support and (more) use-cases. What to focus on first?
- 3. Thoughts about this general approach. Is HiPS the right standard to start from? Or an industry-standard format (e.g., Delta Lake, Iceberg, ...)?
- 4. Cross-matching support how to efficiently handle objects at partition boundaries.
- 5. Updates, appends, transactions, etc...

Interested? Join #lincc-dataformats on the LSSTC Slack.

LINCC summary

Jeno Sokoloski



LINCC Frameworks is a key pillar of LINCC An LSSTC initiative

Goal of LINCC: provide the astrophysics community with the tools, training, and collaborative opportunities – beyond and complementary to those provided by the federally funded project – to enable Rubin LSST to fulfill its potential.

Strategy: with input from LSSTC member institutions, the SCs, and the broader astrophysics community, seek private funding to build programs that have broad community impact and could not be carried out by a single university or PI.

Launched: <u>LINCC Frameworks</u> and the <u>LSSTC Catalyst Fellowship</u> Funded by the John Templeton Foundation.

More to come!



- LSE-319: Science Platform Vision Document
- LSE-61: Data Management System Requirements
- LDM-554: Data Management LSST Science Platform Requirements
- LDM-542: Science Platform Design
- DMTN-202: Use cases and science requirements on a user batch facility
- DMTN-086: Next-to-the-Database Processing Use Cases



Rubin Science Platform – Three Aspect Design

Portal Aspect

Exploratory analysis and visualization of the LSST archive

Notebook Aspect

In-depth 'next-to-data' analysis and creation of added-value data products

API Aspect

Remote access to the LSST archive via Virtual Observatory interfaces





The DPO-era RSP provides delegates with access to the data set via the Portal, Notebook, and API Aspects. All three aspects have tools to query, subset, visualize, and analyze the DPO data set, as well as documentation and tutorials for users. The LSST Science Pipelines (and many other common software packages) are pre-installed in the Notebook environment.



'From Data to Software to Science" session | Rubin Observatory PCW | 08-12 August 2021



Data Preview Schedule and Data Products

Rubin Baseline Data Release Scenario	Jun 2021	Jun 2022	Mar 2024 - Jul 2024	Dec 2024 - Mar 2025	Oct 2025 - Jan 2026	Oct 2026 - Jan 2027	Nov 2027 - Jan 2028	Oct 2028 - Jan 2029
	DP0.1	DP0.2	DP1	DP2	DR1	DR2	DR3	DR4
Data Product	DC2 Simulated Sky Survey	Reproces sed DC2 Survey	ComCam On-Sky Data	LSSTCam On-Sky Data	LSST First 6 Months Data	LSST Year 1 Data	LSST Year 2 Data	LSST Year 3 Data
Raw images		\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark	\checkmark
DRP Processed Visit Images and Visit Catalogs		\checkmark			\checkmark	\checkmark	\checkmark	\checkmark
DRP Coadded Images		\checkmark			\checkmark	\checkmark	\checkmark	\checkmark
DRP Object and ForcedSource Catalogs							\checkmark	\checkmark
DRP Difference Images and DIASources						\checkmark	\checkmark	\checkmark
DRP ForcedSource Catalogs including DIA outputs		\checkmark			\checkmark	\checkmark		\checkmark
PP Processed Visit Images					\checkmark	\checkmark		\checkmark
PP Difference Images						\checkmark		\checkmark
PP Catalogs (DIASources, DIAObjects, DIAForcedSources)						\checkmark	\checkmark	\checkmark
PP Alerts (Canned)					\checkmark	\checkmark		\checkmark
PP Alerts (Live, Brokered)					\checkmark	\checkmark	\checkmark	\checkmark
PP SSP Catalogs								\checkmark
DRP SSP Catalogs						\checkmark	\checkmark	



Getting Help - Community Forum



Community forum

https://community.lsst.org/

Support	Rubin Science Platform	all tag	s ▶ all	• La	test	Тор	Bookmarks	Unread (2)	My Posts	+ New T	opic 🐥
: ≣ Topic		Q							Replies	Views	Activity
Question at	Data Preview 0 × 87							۲			May 16
Can I use lsst-	B DP0 RSP Service Issues >	, ₁₂ ;imu	lation dat	:a?				@ @@	5	137	May 11
dm-dev, butler	Rubin Science Platform ×2	3									
RSP Notebool	Camera ×2	is Pa	atch Thurs	sday (202	2-03-03	3 3pm P	T)	462 🥋			Mar 1
	Lasair										
MeasureMulti	lask implementation U							U T	2	154	Jan 21
CPU resources	s for individual RSP users ir	n the ope	erations-e	ra RSP				<u>s</u> 💀	4	127	Jan 18
RSP / data.lsst	cloud updates - 2021-12-0)3						4 fe >>		121	Dec '21
Copying a coll	ection							0 🔊 🔇	9		Dec '21

Installing pywwt into LSP notebook aspect 🖋

Support Rubin Science Platform



Sep '2

I'm interested in seeing if pywwt 3 could be installed into the LSP Jupyter(Lab) framework as an LSST data visualization option (live pywwt demo notebooks; here 2). It's a pretty straightforward Python package, but full integration does require a bit of fiddling with various sorts of Jupyter extensions to get everything working: (install docs: here 2). The LSP docs say that I should make such a request here, so that's what I'm doing. Thanks!

Solved by adam in post #7

🌘 Sep '20 🛛 😡

jsick 🗊 Jonathan Sid

2:

I have been informed by my manager that installation of pywwt (as driven by this thread) did not follow the process for user requests for new functionality. She is understandably wary of supporting a whole new visualization framework. If you are a DPO.1 user, can you please open an issue against 5...

			\heartsuit				<table-cell-rows> Reply</table-cell-rows>
8 replies	306 _{views}		@? (76	9		~
							Sep '20

Hi @pkgw , thanks for suggesting this. We've made an internal ticket for it: https://jira.lsstcorp.org/browse/DM-26634 3 , but we'll update you here too.