

DM sizing model and purchase plan for the remainder of construction.

Michelle Butler, Kian Tat Lim, William O'Mullane 2020-03-25

1 Introduction

The sizing and cost model for DM has not been revised since 2014. This document presents a simplified sizing model in Section 6.1 based on detailed sizing presented in Section 7. Section 2 presents a very high level budget summary for DM hardware.

As we begin to consider operations one possibility would be to consider the DR1 hardware as an operations cost and only consider the commissioning hardware a construction cost. This is laid out in Section 3.

2 Proposed budget

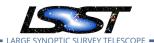
A high level bottom line is given in Table 1. Since Xeon is more expensive that is the number used for the budget calculation, should we get Rome and it really performs we may save a little. The remainder of the document is all the details that went into that.

Table 1: This table pulls together all the information in a high level summary - in this table Xeon pricing is used since that is the more expensive but better known option. Price factors, defined in Table 18 are applied post 2020.

Year	2020	2021	2022	2023
Compute (2019 pricing)	\$690,000	\$0	\$1,510,000	\$2,820,000
Storage (2019 pricing)	\$190,702	\$126,563	\$1,216,867	\$7,864,125
Qserv (2019 pricing)			\$560,000	\$3,240,000
Total (2019 pricing)	\$880,702	\$126,563	\$3,286,867	\$13,924,125
Compute (applying price factor)	\$621,000	\$0	\$1,057,000	\$1,692,000
IN2P3 (50% of compute in ops)				-\$846,000
Storage (applying price factor)	\$181,167	\$113,907	\$1,034,337	\$6,291,300
Qserv (applying price factor)			\$434,000	\$2,268,000
Hosting cost NCSA	\$110,802	\$62,802	\$238,012	\$536,801
Total budget (using price factors)	\$912,969	\$176,709	\$2,763,350	\$9,942,100

In Table 1 we should note that IN2P3 do 50% of processing so we reduce the processing cost by half. This does not reduce the storage cost. We have applied a modest cost reduction assuming that processors and disks get a little cheaper - that percentage is given in Table 18 along with many other parameters. Table 18 also contains the number of nodes we assume

DRAFT 1 DRAFT



to need for Qserv.

Specific costs for storage are detailed in Table 19 and for compute in Table 20 the following budgets can be considered. The detailed annual purchasing based on those prices is given for storage in Table 8 and for compute in Table 6.

3 Potential scope option

In the 2019 JSR we discussed the possibility of purchasing DR1 hardware as part of operations rather than DM construction. Table 2 defines what this would be worth using the cost/sizing model in this document. The table also calls out explicitly the storage and compute costs in Chile for the Data Access centre, these are not deducted but shown to give an idea as to the potential cost there.

Table 2: Considering a scope option of delaying the purchase of LOY1 processing hardware and only purchasing what is needed for commissioning we would only purchase up to and including 2022 hardware of Table 1. If we consider that amount and the current remaining construction budget for hardware the potential worth of such a scope option is given in this table.

Commissioning Budget (to 2022)	\$3,853,028
DM construciton budet remaining	\$13,892,220
Total potential to delay to ops	\$10,039,192
Includes Chilean Storage (2023)	\$1,236,604
Include Chilean Compute (2023)	\$169,200

Should we do this some contingency for extra hardware must be kept in DM construction as well as some manpower budget to aid with the transition to operations.

4 Operations budget estimate

Based on the needs in Table 17 and the costs in Table 13 and Table 20 we get the estimate presented in Table 3

Table 3: This table pulls together all the information in a high level summary for operations - in this table Xeon pricing(see Table 10) is used since that is the more expensive but better known option. Price factors, defined in Table 18 are applied in all cases - other input values come from Table 17, Table 13.

Year (all prices Million\$)	2023	2024	2025	2026	2027	2028	2029	2030	2031	2032
Compute (2019 pricing)	\$2.82	\$2.98	\$4.92	\$6.18	\$6.34	\$7.16	\$6.72	\$6.72	\$7.16	\$7
Qserv (2019 pricing)	\$3.24	\$4.84	\$3.72	\$4.80	\$5.48	\$7.20	\$4.20	\$4.36	\$5.56	\$6.24
Storage (2019 pricing)	\$7.86	\$9.24	\$10.82	\$10.80	\$12.38	\$18.82	\$20.26	\$21.85	\$21.83	\$22.48
Total (2019 pricing)	\$13.92	\$17.06	\$19.46	\$21.78	\$24.20	\$33.18	\$31.18	\$32.93	\$34.55	\$35
Applying price factor (CPU)	\$1.85	\$1.76	\$2.61	\$2.96	\$2.73	\$2.77	\$2.34	\$2.11	\$2.02	\$1.71
IN2P3 (50% of US compute)	-\$0.92	-\$0.87	-\$1.29	-\$1.46	-\$1.35	-\$1.38	-\$1.16	-\$1.04	-\$1.00	-\$0.85

DRAFT 2 DRAFT

Osery (applying factor)

_	\$2.37	\$3.28	\$2.33	\$2.78	\$2.94	\$3.57	\$1.93	\$1.85	\$2.18	\$2.26

quer (applying factor)	42.57	45.20	42.55	42.70	42.5	+5.57	4 1 1 3 5	4 1 100	420	72.20
Applying price factor (Storage)	\$6.41	\$7.15	\$7.96	\$7.54	\$8.21	\$11.86	\$12.13	\$12.43	\$11.79	\$11.54
Hosting Overhead NCSA	\$0.54	\$0.79	\$1.01	\$1.21	\$1.38	\$1.61	\$1.71	\$1.85	\$2.01	\$2.23
Total budget (using price factors)	\$10.25	\$12.11	\$12.61	\$13.02	\$13.91	\$18.44	\$16.95	\$17.20	\$17.01	\$16.89
Total Operations hardware to 2032	\$148.39	million								

Again in Table 3 we assume IN2P3 do 50% of processing. We have applied a compounded modest cost reduction assuming that processors and disks get a little cheaper - that percentage is given in Table 18.

It must be noted that the price of disk and tape have a profound effect over 10 years. We have been fairly conservative on the base prices in Table 19. An even bigger effect is in the compounding of the presumed fall in storage cost. Here we have used an extremely conservative 5% per year (Table 18) - changing this to 15% halves the cumulative ops estimate, setting it to 10% brings the total down by about 30%.

More details on the inputs are in Section 5.1

4.1 US and Chile

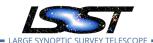
While Table 3 present the total ops cost for Rubin Observatory a fraction of this is in Chile and would potentially remain an NSF cost in operations. Table ?? presents just the US Data Facility budget and Table 5 presents the Chile budget.

Table 4: This table pulls together all the information in a high level summary for USDF operations - in this table Xeon pricing(see Table 11) is used since that is the more expensive but better known option. Price factors, defined in Table 18 are applied in all cases - other input values come from Table 17, Table 14.

Year (all prices Million\$)	2023	2024	2025	2026	2027	2028	2029	2030	2031	2032
Compute (2019 pricing)	\$2.79	\$2.95	\$4.87	\$6.12	\$6.28	\$7.10	\$6.66	\$6.66	\$7.10	\$6.66
Qserv (2019 pricing)	\$1.62	\$2.42	\$1.86	\$2.40	\$2.74	\$3.60	\$2.10	\$2.18	\$2.78	\$3.12
Storage (2019 pricing)	\$6.41	\$7.46	\$8.79	\$9.12	\$10.72	\$15.71	\$16.83	\$18.17	\$18.50	\$19.16
Total (2019 pricing)	\$10.82	\$12.83	\$15.52	\$17.64	\$19.74	\$26.41	\$25.59	\$27.01	\$28.38	\$28.94
Applying price factor (CPU)	\$1.83	\$1.74	\$2.59	\$2.93	\$2.70	\$2.75	\$2.32	\$2.09	\$2.01	\$1.69
IN2P3 (50% of compute)	-\$0.92	-\$0.87	-\$1.29	-\$1.46	-\$1.35	-\$1.38	-\$1.16	-\$1.04	-\$1.00	-\$0.85
Qserv (applying factor)	\$1.19	\$1.64	\$1.17	\$1.39	\$1.47	\$1.78	\$0.96	\$0.92	\$1.09	\$1.13
Applying price factor (Storage)	\$5.22	\$5.77	\$6.46	\$6.37	\$7.11	\$9.90	\$10.08	\$10.33	\$9.99	\$9.84
Hosting Overhead NCSA	\$0.54	\$0.79	\$1.01	\$1.21	\$1.38	\$1.61	\$1.71	\$1.85	\$2.01	\$2.23
Total budget (using price factors)	\$7.86	\$9.08	\$9.93	\$10.43	\$11.31	\$14.67	\$13.91	\$14.16	\$14.10	\$14.04
Total Operations hardware to 2032	\$119.48	million								

Table 5: This table pulls together all the information in a high level summary for Chile operations - in this table Xeon pricing(see Table 12) is used since that is the more expensive but better known option. Price factors, defined in Table 18 are applied in all cases - other input values come from Table 17, Table 15.

Year (all prices Million\$)	2023	2024	2025	2026	2027	2028	2029	2030	2031	2032
Compute (2019 pricing)	\$0.03	\$0.04	\$0.05	\$0.07	\$0.07	\$0.07	\$0.07	\$0.07	\$0.07	\$0.07



Qserv (2019 pricing)	\$1.62	\$2.42	\$1.86	\$2.40	\$2.74	\$3.60	\$2.10	\$2.18	\$2.78	\$3.12
Storage (2019 pricing)	\$1.45	\$1.78	\$2.03	\$1.68	\$1.66	\$3.11	\$3.43	\$3.69	\$3.33	\$3.32
Total (2019 pricing)	\$3.10	\$4.24	\$3.94	\$4.15	\$4.47	\$6.78	\$5.60	\$5.94	\$6.18	\$6.51
Applying price factor (CPU)	\$0.02	\$0.02	\$0.03	\$0.03	\$0.03	\$0.03	\$0.02	\$0.02	\$0.02	\$0.02
Qserv (applying factor)	\$1.19	\$1.64	\$1.17	\$1.39	\$1.47	\$1.78	\$0.96	\$0.92	\$1.09	\$1.13
Applying price factor (Storage)	\$1.18	\$1.37	\$1.49	\$1.17	\$1.10	\$1.96	\$2.06	\$2.10	\$1.80	\$1.70
Total budget (using price factors)	\$2.39	\$3.04	\$2.69	\$2.60	\$2.60	\$3.77	\$3.04	\$3.05	\$2.91	\$2.85
Total Operations hardware to 2032	\$28.94	million								

5 Cost details

The summary table (Table 1) uses Xeon pricing for compute as shown in Table 6.

Table 6: Implementation with Intel Xeon

Year	2020	2021	2022	2023
Number of Xeon	69	0	151	282
Approximate cost	\$690,000.00	\$0.00	\$1,510,000.00	\$2,820,000.00

An alternative architecture would be Rome - we have not tested this but if it performs well it may provide savings. Table 7 gives the price of compute based on Rome -small and large only for comparison purposes - not used in the calculations.

Table 7: Implementation with AMD Rome (we have no good proce for these reallly)

Year	2020	2021	2022	2023
number of small rome	49	0	75	200
Approximate cost of small rome	\$637,000.00	\$0.00	\$975,000.00	\$2,600,000.00
number of large rome	16	0	25	66
Approximate cost of large rome	\$379,200.00	\$0.00	\$592,500.00	\$1,564,200.00

Table 8 gives the price of storage using all types that we need. This would be needed regardless of the compute chosen.

Table 8: Total storage cost estimate

Year	2020	2021	2022	2023
Fast Storage	\$11,842.11	\$11,842.11	\$26,070.00	\$156,420.00
Fast Storage Chile	\$0.00	\$0.00	\$0.00	\$156,420.00
Normal Storage	\$91,764.54	\$9,199.77	\$741,634.14	\$4,015,129.31
Latent Storage	\$14,333.41	\$25,083.47	\$184,042.71	\$1,074,948.27
Latent Storage Chile	\$0.00	\$0.00	\$0.00	\$1,298,407.87
High Latency Storage	\$72,762.21	\$80,438.13	\$265,120.43	\$1,162,799.08
Total	\$190,702.27	\$126,563.48	\$1,216,867.29	\$7,864,124.52

Table 9 gives the annual cost of hosting compute. This includes purchasing racks to house new nodes etc.

Table 9: Overheads(NCSA) per year based on number of cores in Table 18 and costs in Table 21 assuming Xeon density from Table 20.

Year	2020	2021	2022	2023	1
------	------	------	------	------	---

DRAFT 4 DRAFT



Total Incremental cores (USA)	1,836	0	4,026	7,521
Total owned cores (USA)	3,528	3,528	7,554	15,075
Total owned nodes	111	111	251	567
Cost for hosting nodes	\$62,802	\$62,802	\$142,012	\$320,801
Total new nodes	58	0	140	317
Total new racks	2	0	4	9
Rack install cost	\$48,000.00	\$0.00	\$96,000.00	\$216,000.00
Total Overhead (NCSA)	\$110,802.27	\$62,802.27	\$238,012.35	\$536,800.80

5.1 Ops Cost details

Table 10 gives the price of compute based on Xeons. This is broken down further for US in Table 11 and Chile in Table 12

Table 13 gives the price of storage using all types that we need. This is broken down further for US in Table 14 and Chile in Table 15 This would be needed regardless of the compute chosen.

Table 10: Implementation with Intel Xeon for full Rubin Observatory

Year	2023	2024	2025	2026	2027	2028	2029	2030	2031	2032
Number of Xeon	282	298	492	618	634	716	672	672	716	672
Approximate cost (2019 Mdollars)	\$2.82	\$2.98	\$4.92	\$6.18	\$6.34	\$7.16	\$6.72	\$6.72	\$7.16	\$6.72

Table 11: Implementation with Intel Xeon for USDF

Year	2023	2024	2025	2026	2027	2028	2029	2030	2031	2032
Number of Xeon USDF	279	295	487	612	628	710	666	666	710	666
Approximate cost (2019 Mdollars)	\$2.79	\$2.95	\$4.87	\$6.12	\$6.28	\$7.10	\$6.66	\$6.66	\$7.10	\$6.66

Table 12: Implementation with Intel Xeon for Chile Compute

Year	2023	2024	2025	2026	2027	2028	2029	2030	2031	2032
Number of Xeon Chile	3	4	5	7	7	7	7	7	7	7
Approximate cost (2019 Mdollars)	\$0.03	\$0.04	\$0.05	\$0.07	\$0.07	\$0.07	\$0.07	\$0.07	\$0.07	\$0.07

Table 13: Total storage cost estimate for operations of Rubin Observatory USDF and CHile

Year (all in M\$)	2023	2024	2025	2026	2027	2028	2029	2030	2031	2032
Fast Storage	\$0.31	\$0.35	\$0.43	\$0.16	\$0.16	\$0.44	\$0.48	\$0.55	\$0.28	\$0.25
Normal Storage	\$4.02	\$3.91	\$4.26	\$4.25	\$4.97	\$8.25	\$8.17	\$8.54	\$8.54	\$8.54
Latent Storage	\$2.37	\$3.17	\$3.64	\$3.19	\$3.38	\$5.57	\$6.36	\$6.83	\$6.39	\$6.39
High Latency Storage	\$1.16	\$1.80	\$2.50	\$3.19	\$3.87	\$4.56	\$5.25	\$5.93	\$6.62	\$7.31
Total (M\$)	\$7.86	\$9.24	\$10.82	\$10.80	\$12.38	\$18.82	\$20.26	\$21.85	\$21.83	\$22.48

Table 14: Total storage cost estimate for operations at USDF

Year (all in M\$)	2023	2024	2025	2026	2027	2028	2029	2030	2031	2032
Fast Storage USDF	\$0.16	\$0.16	\$0.22	\$0.08	\$0.09	\$0.22	\$0.22	\$0.27	\$0.14	\$0.12
Normal Storage USDF	\$4.02	\$3.91	\$4.26	\$4.25	\$4.97	\$8.25	\$8.17	\$8.54	\$8.54	\$8.54
Latent Storage USDF	\$1.07	\$1.59	\$1.82	\$1.60	\$1.78	\$2.67	\$3.18	\$3.41	\$3.19	\$3.19
High Latency Storage USDF	\$1.16	\$1.80	\$2.50	\$3.19	\$3.87	\$4.56	\$5.25	\$5.93	\$6.62	\$7.31
Total (M\$)	\$6.41	\$7.46	\$8.79	\$9.12	\$10.72	\$15.71	\$16.83	\$18.17	\$18.50	\$19.16

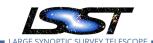


Table 15: Total storage cost estimate for operations in Chile

Year (all in M\$)	2023	2024	2025	2026	2027	2028	2029	2030	2031	2032
Fast Storage Chile	\$0.16	\$0.19	\$0.22	\$0.08	\$0.07	\$0.22	\$0.25	\$0.27	\$0.14	\$0.12
Latent Storage Chile	\$1.30	\$1.59	\$1.82	\$1.60	\$1.60	\$2.89	\$3.18	\$3.41	\$3.19	\$3.19
Total (M\$)	\$1.45	\$1.78	\$2.03	\$1.68	\$1.66	\$3.11	\$3.43	\$3.69	\$3.33	\$3.32

Table 9 gives the annual cost of hosting compute. This includes purchasing racks to house new nodes etc.

Table 16: Overheads(NCSA) per year based on number of cores in Table 17 and costs in Table 21 assuming Xeon density from Table 20.

				1	T		ī	1	
Year	2023	2024	2025	2026	2027	2028	2029	2030	
Total Incremental cores (USA)	7,521	7,958	8,979	8,979	8,979	8,979	8,979	8,979	
Total owned cores (USA)	15,075	23,033	32,012	40,990	49,969	58,947	67,926	76,904	
Total owned nodes	567	936	1,310	1,629	1,926	2,294	2,559	2,812	
Cost for hosting nodes	\$320,801	\$529,576	\$741,180	\$921,666	\$1,089,704	\$1,297,914	\$1,447,847	\$1,590,991	\$1
Total new nodes	317	370	374	401	418	461	386	390	
Total new racks	9	11	11	12	12	13	11	11	
Rack install cost	\$216,000.00	\$264,000.00	\$264,000.00	\$288,000.00	\$288,000.00	\$312,000.00	\$264,000.00	\$264,000.00	\$28
Total Ops Overhead (NCSA)	\$536,800.80	\$793,575.92	\$1,005,179.97	\$1,209,665.78	\$1,377,704.29	\$1,609,913.63	\$1,711,846.98	\$1,854,990.90	\$2,01

Various other inputs to ops costing are given in Table 17.

Table 17: Various inputs for deriving costs in operations - these drive the costs in Table 3. This is based on Table 10, Table 13

Year	2023	2024	2025	2026	2027	2028	2029	2030	2031	2032
Core-hours Needed Total (DRP)	4.5E+07	8.2E+07	1.2E+08	1.6E+08	2.0E+08	2.5E+08	2.9E+08	3.3E+08	3.7E+08	4.1E+08
Core-hours Annual Increase	3.40E+07	3.6E+07	4.1E+07							
Time to Process days	200	200	200	200	200	200	200	200	200	200
Time to Process hours	4,800	4,800	4,800	4,800	4,800	4,800	4,800	4,800	4,800	4,800
Cores (DRP) Annual increase	7,093	7,594	8,512	8,512	8,512	8,512	8,512	8,512	8,512	8,512
Cores (DRP) Annual refresh			2,837	7,093	7,594	8,512	8,512	8,512	8,512	8,512
Cores (DRP) Annual purchase	7,093	7,594	11,349	15,605	16,106	17,024	17,024	17,024	17,024	17,024
Cores (Alerts)	1,188	1,188	1,188	1,188	1,188	1,188	1,188	1,188	1,188	1,188
Cores (Alerts) Annual refresh			1,188			1,188			1,188	
Cores (US DAC/ Staff)	568	933	1,399	1,866	2,332	2,798	3,265	3,731	4,198	4,664
Cores (US DAC/ Staff) Annual increase	428	364	466	466	466	466	466	466	466	466
Cores (US DAC/ Staff) Annual refresh			141	428	364	466	466	466	466	466
Cores (US DAC/ Staff) Annual purchase	428	364	607	894	831	933	933	933	933	933
Cores (Chilean DAC)	103	187	280	373	466	560	653	746	840	933
Cores (Chilean DAC) Annual increase	78	83	93	93	93	93	93	93	93	93
Cores (Chilean DAC) Annual refresh	1.5		26	78	83	93	93	93	93	93
Cores (Chilean DAC) Annual purchase	78	83	119	171	177	187	187	187	187	187
Oserv nodes (US DAC/ Staff)	95	216	309	348	364	451	436	408	367	418
Qserv nodes (US DAC/ Staff) Annual Increase	81	121	93	120	137	180	105	109	139	156
Qserv nodes (Chilean DAC)	95	216	309	348	364	451	436	408	367	418
Qserv nodes (Chilean DAC) Annual Increase	81	121	93	120	137	180	105	109	139	156
Total Cores Annual Increase	7,599	8,042	13,264	16,670	17,113	19,332	18,144	18,144	19,332	18,144
Fast Storage (TB)	206	371	586	667	735	798	859	918	974	1029
Annual Increase (Fast)	156	164	215	81	68	63	60	59	57	55
Annual Refresh (Fast)					26	156	164	215	81	68
Annual Purchase (Fast)	156	164	215	81	94	220	225	275	138	123
Normal Storage (TB)	38,983	67976	99538	131025	162321	193727	225288	256991	288788	320731
Annual Increase (Normal)	29,742	28993	31563	31487	31296	31406	31560	31703	31797	31943
Annual Refresh (Normal)					5,494	29,742	28,993	31,563	31,487	31,296
Annual Purchase (Normal)	29,742	28,993	31,563	31,487	36,790	61,148	60,553	63,266	63,284	63,239



Latent Storage (TB)	28,854	64,086	104,491	139,969	175,447	210,925	246,403	281,881	317,359	352,837
Annual Increase (Latent)	23,888	35,232	40,405	35,478	35,478	35,478	35,478	35,478	35,478	35,478
Annual Refresh (Latent)					4,090	23,888	35,232	40,405	35,478	35,478
Annual Purchase (Latent)	23,888	35,232	40,405	35,478	39,568	59,366	70,710	75,884	70,956	70,956
High Latency (TB)	63,245	135,129	234,931	362,490	517,473	699,899	909,793	1,147,179	1,412,074	1,704,500
Annual Increase (High Latency)	46,512	71,884	99,803	127,559	154,983	182,426	209,894	237,386	264,894	292,426
Chilean DAC Fast Storage (TB)	156	347	562	643	711	774	835	894	951	1,006
Annual Increase (Fast Chilean DAC)	156	190	215	81	68	63	60	59	57	55
Annual Refresh (Fast Chilean DAC)						156	190	215	81	68
Annual Purchase (Fast Chilean DAC)	156	190	215	81	68	220	251	275	138	123
Chilean DAC Latent Storage (TB)	28,854	64,086	104,491	139,969	175,447	210,925	246,403	281,881	317,359	352,837
Annual Increase (Latent Chilean DAC)	28,854	35,232	40,405	35,478	35,478	35,478	35,478	35,478	35,478	35,478
Annual Refresh (Latent Chilean DAC)						28,854	35,232	40,405	35,478	35,478
Annual Purchase (Latent Chilean DAC)	28,854	35,232	40,405	35,478	35,478	64,332	70,710	75,884	70,956	70,956

6 Models

6.1 Sizing model

An exhaustive and detailed mode is provided in [LDM-138; LDM-144] - here we concentrate on the needs for the final years of construction. We explore the compute and storage needed to get us through commissioning and suggest a 2023 purchase for DR1,2 processing which could be pushed to operations.

Table 18 gives the annual requirements for the next few years.

Table 18: Various inputs for deriving costs - 2019 represents currentl holdings.

Year	2019	2020	2021	2022	2023
Core-hours Needed Total (DRP)		4.41E+06	4.41E+06	1.12E+07	4.53E+07
Annual Increase		4.41E+06	0.00E+00	6.81E+06	3.40E+07
Time to Process days		100.0	100.0	100.0	200
Time to Process hours		2,400	2,400	2,400	4,800
Instantaneous cores (DRP) Annual in-	1152	1,836	0	2,837	7,093
crease					
Instantaneous cores (Alerts)		0	0	1188	1188
Cores (Alerts) Annual increase		0	0	1188	0
Instantaneous cores (US DAC/ Staff)	540	540	540	141	568
Cores (US DAC/ Staff) Annual increase		0	0	0	428
Instantaneous cores (Chilean DAC)		0	0	26	103
Cores (Chilean DAC) Annual increase		0	0	26	78
Qserv nodes (US DAC/ Staff)				14	95
Qserv nodes (US DAC/ Staff) Annual				14	81
Increase					
Qserv nodes (Chilean DAC)				14	95
Qserv nodes (Chilean DAC) Annual In-				14	81
crease					
Total Cores Annual Increase		1,836	0	4,051	7,599
Fast Storage (TB)		12	24	50	206
Annual Increase (Fast)		12	12	26	156
Normal Storage (TB)	3000	3680	3748	9241	38983
Annual Increase (Normal)		680	68	5494	29742
Latent Storage (TB)		319	876	4966	28854



Annual Increase (Latent)	319	557	4090	23888
High Latency (TB)	2910	6128	16733	63245
Annual Increase (High Latency)	2910	3218	10605	46512
Chilean DAC Fast Storage (TB)				156
Annual Increase (Fast Chilean DAC)				156
Chilean DAC Latent Storage (TB)				28854
Annual Increase (Latent Chilean DAC)				28854
Annual price decrease CPU	10%			
Annual price decrease Storage	5%			
Annual price decrease Qserv	8%			

6.2 Compute and storage

We which to base our budget on reasonable well know machines for which we have well know prices. Table 20 gives an outline of a few standard machines we use and a price. This table also gives a FLOP estimate for those machines. Table 19 gives costs for different types of storage we will require various latency for different tasks and those have varying costs. These tables are used as look ups for the cost models in Section 2

Table 19: Storage types and costs used as inputs used for calculations

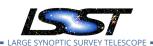
Storage type	cost
fast – NVME (50GB/ s each) / TB	\$1,000.00
normal - SATA GPFS file systems/ TB	\$135.00
latency – slower but on disk	\$45.00
high latency – very slow – on tape	\$25.00

In Table 19 we should consider for NVME for each TB with file system servers two DDN NVME box with GPFS servers. The price is based on the TOP performer with best price. The Normal price is for each TB with file system disks and servers locally attached to production resources.

In the latency and high latency prices are only at NCSA: for each TB with file systems and all people/services. The complete service not usually attached. S3 bucket type. Can be mounted if needed but not for production worthy speeds. The complete service with data flowing to tape using policies.

Table 20: Machine types and costs used as inputs used for calculations

Type of machine	Cores	Memory(GB)	Eff cores/ node	Cost	purpose/ use
xeon	32	192	27	\$10,000.00	current K8 node
qserv	12	128	12	\$20,000.00	current qserv node
small rome	64	256	38	\$13,000.00	https://www.microway.com/product/navion-1u-amd-epyc-gpu-server/
large rome	128	512	116	\$23,700.00	
current compute node	24	128	24	\$9,000.00	current compute node



There is also an associated running cost for machines included in the total cost of ownership. These overheads are listed in Table 21.

Table 21: Overhead costs per rack

Item	Number/ Cost
Compute nodes in a rack	36
Rack initial cost has power, networking switches, networking cables, ready for machine installation-switches last 5 years. Will need to refresh, but rack should last entire project.	\$24,000.00
** need to add annually: floor space for rack for 1 years. need to renew after new nodes are racked/ stacked	\$300
** Need to add annually: power for 1 node for 1 yr - kw * rate * hours/ year *	\$348
** need to add annually: cooling for 1 node for 5 years kw* chillded wa- ter per MTBU* hours/ year * 1KW in (MTBU)	\$210
** Need to add annually: mainte- nance for node s – can't purchase more than what the contract has in time left. could be included in the price of the machine, and might not be added in here.	\$1,500
Cost for each machine for 1 year in a rack.	\$566
**** need to add in at an annual basis. software maintenance (ora- cle and other software not associated with specific node annually) Oracle li- cense, VM licensing.	\$35,000
Power per Rack (for Chile) Watts	12096
Approx PB per Storage Rack	8
Compute node lifetime (years)	3
Storage lifetime (years)	5
Chile Power CLP / KW-hr	105.65

7 Sizing inputs

The following simplified sizing was used to give the input sizes for the cost model in Section 2. The storage sizes are given in Table 24 and Table 25 while the compute is given in Table 27 and Table 28.

7.1 Processing Plan

This model assumes the following processing:

• Precursor data (HSC RC2 and a similarly-sized DESC DC2 subset) is reprocessed each

DRAFT 9 DRAFT



month during the Construction period using the Data Release Production (DRP).

- A large precursor reprocessing of HSC PDR2 (or equivalent) is completed twice a year.
 Products from one of these reprocessings will be released as Data Preview 0. One or
 more of these processings during Commissioning could be devoted instead to ComCam
 science data for Data Preview 1 or LSSTCam science data in preparation for Data Preview
 2.
- Alert Production (AP) processing happens continuously as LSSTCam science images are obtained. AP hardware is purchased in FY22 to support this, presumably on a limited basis during that year and then to whatever extent possible during the first year of the survey.
- Commissioning processing of LSSTCam science images for Data Preview 2 is assumed to happen towards the end of FY22 as a single execution of the DRP. The hardware for this can be purchased early in that fiscal year.
- Annual DRP execution starts at the beginning of LSST Operations Year 2 with the processing for DR2. The hardware for each year's processing must be purchased and ready for use at the beginning of the year, so it is allocated in the tables to the prior fiscal year, when the images for that processing were taken.
- DR1 processing begins after the first 6 months of the survey; the hardware for this can be part of the DR2 purchase during FY23.

Some storage for raw data needs to be in place at the beginning of the fiscal year, but it can be ramped up over the course of the year. As a simplification it is allocated to the fiscal year in which it will be used.

7.2 Storage Model

Table 22: Inputs used to calculate storage needs

Parameters	unit	FY2020	FY2021	FY2022	FY2023/ LOY1	Notes
Objects	number			4.58E+09	2.75E+10	from LSE-81, scaled to 2 months for 2022, ComCam ignored
Sources	number			1.50E+11	9.01E+11	from LSE-81, scaled to 2 months for 2022, ComCam ignored
ForcedSources	number			4.85E+11	2.91E+12	from LSE-81, scaled to 2 months for 2022, ComCam ignored
Science users	users	50	100	5000	5000	"Stack Club" to 2021, DP users thereafter
Storage per science user	TB	0.1	0.2	0.2	0.4	ramp to LSE-81 number; includes oversubscription
LSSTCam image size	TB	0.0152				uncompressed, 32 bit, with overscan and corner rafts
Raw image compression	factor	0.42				lossless-compressed divided by uncompressed for raws
Lossy image compression	factor	0.250				lossy-compressed divided by lossless-compressed for PVIs
Observing nights per year	nights	300				maximum
Visits per night	visits	1000				maximum



Images per visit	images	2					
Calibration images per day	images	500					
LSSTCam Science images	images			100000	600000	test images until 2 months of science in 2022	
LSSTCam Test images	images	25000	50000	50000		ramp to science images	
LSSTCam Engineering images	images	12500	12500	15000	6000	decreasing ramp	
LSSTCam Calibration images	images	12500	25000	37500	150000	estimates based on science and test images; actual for 2023	
Number of coadd data products	number	2				deep and good-seeing	
Object table row size	bytes			1896	1896	from LDM-141	
Object_Extra tables row size	bytes			21005	21005	from LDM-141	
Source table row size	bytes			467	467	from LDM-141	
ForcedSource table row size	bytes			41	41	from LDM-141	
Qserv replication factor	factor	3.0	3.0	3.0	3.0		

7.2.1 Overview

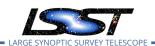
This simplified storage model eliminates many details in the previous storage model [LDM-141] that end up being insignificant. There are relatively few data products that require significant amounts of fast SSD or slower disk or tape storage; the others complicate the model without giving much insight. In addition, it is assumed that bandwidth is not a significant constraint, other than the distinction between SSD and spinning disk. With the advent of highly-parallel shared and object storage, having large numbers of spindles solely to achieve high bandwidth for certain operations is not thought to be necessary.

Values are computed for the amount of storage expected to be "on the floor" at the beginning of each fiscal year from FY2020 through FY2023 (which is LSST Operations Year 1). Not included is any storage already present at the end of FY2019 holding past data.

Key scientific and algorithmic assumptions made include:

- All significant intermediates and data products generated by Data Release Production processing need to be kept on filesystem disk until the DRP is complete. Some scratch space is provided to hold small, temporary intermediates. If some intermediates could be removed during DRP when it is known they will no longer be needed, some space savings could be realized.
- HSC RC2 processing is representative of the outputs that DRP will generate. In particular, the presence or absence of "heavy footprints" is assumed to be correct. The coadd storage is doubled to account for an additional "good-seeing" coadd along with the existing "deep" coadd.
- Processed visit images (PVIs) and catalogs in Parquet format start on "normal" filesystem

DRAFT 11 DRAFT



disk but then move to object storage at the completion of the DRP, with lossy compression of the PVIs at that time. This is in accordance with RFC-325, although the relevant LCR has not yet been approved. Object storage is expected to be cheaper and more scalable for read-only data products; filesystem storage is used for data that is being generated or modified.

- Raw images and coadd images are only temporarily stored on filesystem disk and are then rapidly moved to object storage, where they are retained.
- Intermediates like warped images for coaddition are not survey data products and do not need to be kept beyond the end of the DRP and subsequent QA.

All data is backed up to tape permanently, including annual snapshots of filesystems. Any incremental backups are assumed to be reusable or otherwise purged and hence not significant.

7.2.2 Parameters

The key parameters in Table 22 are described below.

The numbers of Objects, Sources, and ForcedSources are taken from LSE-81, with the FY2022 numbers reduced by a factor of 2/12 to account for the anticipated 2 months of on-sky science validation time for LSSTCam before the survey begins. These numbers are ultimately based on models for stars in the galaxy and galaxies in the universe that are dependent on the limiting magnitude achieved in each year.

The numbers of science users are estimates, using "Stack Club" users and Commissioning users for FY2020 and 2021, followed by US science users in FY2022 and FY2023 for Data Preview data. The bulk of US science users are not expected to arrive until after Data Release 1 at the beginning of FY2024.

Storage per science user is estimated based on today's usage at NCSA, scaled up as users become more active, and approaching the number given in LSE-81 as Operations begins. Note that it is expected that there will be a wide distribution of usage by user, with some using almost none and some using much more than their proportional share.

The LSSTCam image size is uncompressed and includes overscan, 4 bytes of raw data per

DRAFT 12 DRAFT



pixel, and both science and corner rafts.

The raw image compression factor was measured on simulated LSST images. The lossy image compression factor for processed visit images is the ratio between the lossy-compressed file size (estimated at 1/6 of uncompressed) and the lossless-compressed file size (estimated at 66% of uncompressed). Note that PVIs do not compress losslessly as well as raw images due to their floating point planes.

The number of observing nights per year and the number of visits per night are maximal estimates. 2 images per visit is still the baseline and a possibility that must be accounted for. The number of calibration images per day was derived from the calibration plan.

Two complete all-sky coadds are assumed, one for "good seeing" and one deep.

As stated above, the number of LSSTCam science images is scaled by 2/12 for FY2022 given the length of science validation time. The number of test images, taken on test stands, is estimated as a ramp up to the full science cadence. The numbers of engineering (unprocessed) and calibration images are estimated as ramping-down fractions of the number of science and test images, with calibration images ending at the number per day given previously.

Sizes of rows in various data product tables are taken from LDM-141, which was in turn derived from the DPDD.

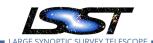
Qserv replicates its data for fault tolerance; a typical replication factor is selected here.

7.2.3 Data Product Sizing

Images and the results of processing them are the dominant factor controlling the storage sizing which is outlined in Table 23. Precursor survey and LSSTCam images are the largest; ComCam, at less than 5% of the size of LSSTCam and with little on-sky science time is negligible, as is LATISS, which is less than 1% of the size of LSSTCam, though it has considerable on-sky time.

The sizing of the Alert Production Database (APDB) is based on experiments in Salnikov (DMTN-113) which found that 57,000 visits took 4.5 TB including indexes. A simple linear scaling to a full year's visits was performed, with half that purchased in 2020 for large (but not full) scale

DRAFT 13 DRAFT



testing.

HyperSuprime-Cam (HSC) RC2 is a relatively small dataset used for monthly processing tests, but it is highly representative of the currently-known DRP work and so is used as the basis for scaling. The size of the input images was taken from Wood-Vasey et al. (DMTN-091); the size of the outputs (image and Parquet/other non-image files) was measured from the latest execution. A similar size dataset based on DESC DC2 is assumed to be being used for an additional monthly processing test. Note that this is a very small subset of the full DESC DC2, which is expected to cover 300 square degrees to 10-year LSST depth (approximately 1000 epochs per point on the sky). The full DESC DC2 is not currently scheduled to be reprocessed by the construction team. Instead, twice-a-year processings of the full HSC SSP PDR2 dataset (including PDR1) are assumed to occur. The size of this dataset was measured on disk; it is 2,564,358 CCD images, each at 18.2 MB (approximately three times the size of PDR1 alone).

Output sizes are assumed to scale linearly with input size, and by the same factor for each instrument, except for coadds which scale by the sky area processed. While the Object catalog ought to be proportional to sky area as well, its size is expected to be dominated by Source and ForcedSource, so we conservatively make them all proportional to input size (visits) for the precursor data where we do not have object count estimates. For LSSTCam, we use the catalog row estimates to derive Qserv table sizes, but the Parquet file sizes are scaled based on HSC, as they may differ from the Qserv schema.

Scratch space is set at 10% of the output image storage for LSSTCam processing; it is assumed to be already present for precursor processing.

Qserv Czar fast (SSD) storage is assumed to be used for the primary Object table; additional space for the so-called "secondary index" mapping object identifiers to spatial chunks is negligible in comparison.

The main Qserv database storage is based on the Parquet file sizing for precursor data and on the estimated numbers of Objects, Sources, and ForcedSources for LSSTCam data.

Note that no space is explicitly reserved for Qserv query result storage.

An additional 20% disk and tape storage is added to account for all other needs.

DRAFT 14 DRAFT

Table 23: Inputs on dataset sizes used to calculate storage needs

Dataset Sizing	unit	FY2020	FY2021	FY2022	FY2023/ LOY1	Notes
HSC RC2 Area	deg2□	3.0	3.0	3.0	3.0	
HSC SSP PDR2 Area	deg2□	300	300	300	300	
DESC DC2 Area	deg2□	300	300	300	300	
LSSTCam Area	deg2□			2000	17000	
APDB	TB	12	24	24	24	4.5/ 57K TB per visit; 1 year retention; 6 months in 2020
HSC RC2 Input Images	TB	0.8	0.8	0.8	0.8	428 visits * 104 CCDs * 18.2 MB uncompressed
HSC RC2 Output Images	TB	2.4	2.4	2.4	2.4	lossless-compressed, not including warps
HSC RC2 Output Coadd Images	TB	0.7	0.7	0.7	0.7	lossless-compressed
HSC RC2 Output Catalogs	TB	1.4	1.4	1.4	1.4	
HSC SSP PDR2 Input Images	TB	93.3	93.3	93.3	93.3	2564358 CCDs * 18.2 MB uncompressed (3 * PDR1)
DESC DC2 Input Images	TB	455	455	455	455	300 sq deg, 10 year depth
Object store datasets:						
LSSTCam Raw Images	TB	319	557	1290	4816	compressed, immediate object store
LSSTCam Output Coadd Images	TB			909	7727	lossless-compressed, immediate object store
Normal disk datasets:						
Precursor Input Images	TB	549	549	549	549	HSC RC2, HSC PDR2, DC2
Precursor Output Images	TB	916	916	916	916	monthly RC2 and DC2 subset plus biannual PDR
Precursor Output Parquet	TB	361	361	361	361	
LSSTCam Output Images	TB			2248	13485	lossless-compressed, moves to object store
LSSTCam Output Parquet	TB			1329	7973	moves to object store
Scratch	TB			225	1349	10% of output images
Qserv Czar/ Object	TB			26	156	based on row sizes and counts
Qserv Database	TB	1088	1088	585	3510	based on Parquet for preliminary; based on row sizes and counts
Science User Home	TB	5	20	1000	2000	
Other/ Misc	TB	761	814	2003	8684	20% of total

7.2.4 Storage Sizing

Finally, storage is allocated to specific types as shown in Table 24. Fast storage (SSD) is used for the APDB and Qserv Czar, which accumulates data from year to year until Data Releases are retired. Normal storage is used for the datasets labeled as such, including output images (initially), output catalogs, and scratch. Local Qserv storage is used for Qserv catalogs. It is assumed that precursor data will be removed from Qserv once LSST data is available, but the LSST data accumulates from year to year. Object storage is used for raw images, lossy-compressed output images, lossless-compressed coadd images, and output catalogs. It also accumulates from year to year. Tape is used for long-term archiving of all the raw images, all the data products (both filesystem and object store), and filesystem backups. Again, this accumulates from year to year.

Note that no replication is assumed in the object store.

Table 24: On floor LDF storage estimates based on Table 23 and Table 22 $\,$

Storage Sizing (on the floor)	unit	FY2020	FY2021	FY2022	FY2023/ LOY1	Notes
APDB	TB	12	24	24	24	
Qserv Czar/ Object	TB			26	182	accumulates with time
Total Fast	ТВ	12	24	50	206	SSD, sum of previous two rows
Total Normal	ТВ	3680	3748	9241	38983	enterprise-grade SATA
Total Qserv Storage	ТВ	1088	1088	585	4094	local consumer-grade SATA, accumulates with time



LSSTCam Raw Images	TB	319	876	2166	6982	accumulates with time	
LSSTCam Output Images	TB			562	3933	lossy-compressed, accumulates with time	
LSSTCam Output Coadd Images	TB			909	8636	accumulates with time	
LSSTCam Output Parquet	TB			1329	9302	accumulates with time	
Total Object Store	ТВ	319	876	4966	28854	consumer-grade SATA, sum of previous four rows	
LSSTCam Raw Images	ТВ	319	876	2166	6982	accumulates with time	
All Data Products/ Backup	TB	2502	====	10.550	17000	normal storage minus Qserv/ scratch, accumulates with t	
All Data Products/ Backup	16	2592	5252	13658	47626	normal storage minus Qserv/ scratch, accumulates with time	
All Object Store-Only Products	TB	2592	5252	13658 909	47626 8636	normal storage minus Qserv/ scratch, accumulates with time accumulates with time	

An additional table (Table 25) gives the storage needs in the Chilean Data Access Center (DAC). This comprises Qserv fast and local storage plus the data products in object storage. Since no DRP computation occurs in Chile, no "normal" filesystem disk is required. Chilean user home directories are assumed to be negligible at this level.

Table 25: On floor Chile storage estimates for Base Data Center

Chile Storage (on the floor)	unit	FY2020	FY2021	FY2022	FY2023/ LOY1	Notes
Fast	TB				156	SSD
Normal	TB				0	Enterprise-grade SATA
Qserv Storage	TB				4094	Local consumer-grade SATA
Object Store	TB				28854	
Tape	TB				0	

7.3 Compute Model

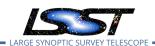
Table 26: Inputs used to calculate compute needs

Parameters	units			Notes
HSC PDR1 Input Images	ТВ	13.7		7238 visits of 104 CCDs
HSC PDR1 small-memory compute	core-hours	64392		measured on E5-2680 v3 @ 2.50GHz
HSC PDR1 high-memory compute	core-hours	78523		measured on E5-2680 v3 @ 2.50GHz
Small-memory DRP algorithm ratio	factor	1.5		image differencing, etc.
High-memory DRP algorithm ratio	factor	2.5		stackfit, etc.
DRP compute per TB	core-hours/ TB	2.1E+04		
Percent DRP on high-memory	factor	67%		
ap_pipe single-core sec/ CCD	core-sec/ CCD	166		measured 83 on DECam 2Kx4K CCD
Additional AP steps	factor	1.25		DCR, real_bogus, etc.
AP compute per visit	core-hours/ visit	1.1E+01		
Qserv data/ node	TB	43.2		1 GB/ sec for 12 hours

7.3.1 Overview

This simplified computing model (Table 26) divides computation into three classes: Data Release Production (DRP), Alert Production, and LSST Science Platform (for the US DAC, Chilean DAC, and LSST staff internal use). Calibration Products Production is assumed to be negligible.

The pipelines have advanced considerably in terms of fidelity and science performance since



the previous computing model [LDM-138] was developed. Scaling compute needs based on an execution of the nascent DRP pipeline on HSC PDR1 data and nightly executions of the nascent ap_pipe pipeline on HiTS2015 data is thus appropriate, but the fact that several steps are still missing from these pipelines must be taken into account.

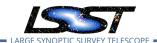
Elapsed times are measured on existing hardware and converted into core-hours on a nominal CPU (Intel Xeon E5-2680v3 at 2.50 GHz). For example, if a pipeline running on precursor data took an average of one hour on a 32-core nominal CPU, 32 core-hours would be used as its compute requirement. This estimation methodology incorporates all I/O, memory bandwidth, cache miss, and other overheads into the core-hour measurement, simplifying calculations. Note that the nominal CPU does not evolve with time; if future CPUs do more work per core, the actual core-hours may be less than estimated here.

Scaling to other CPUs of the same architecture is based on the ratios of nominal GHz clock rates and core counts. For different architectures (e.g. Rome), the scaling is based on the ratio of industry-reported achievable FLOPS for the two architectures.

Key scientific and algorithmic assumptions are:

- DRP compute time is proportional to the input data size (or, equivalently, the number of visits). While certain tasks are undoubtedly proportional to sky area or number of Objects, overall the pipeline elapsed times are a better fit to the number of visits. Some of this may be because the Object density increases as the number of visits to the same sky patch increases.
- HSC PDR1 processing is generally representative of the final DRP, with an allocation for future additional steps as described below.
- Qserv node counts should remain proportional to the size of data loaded into the database
 in order to maintain sufficient disk bandwidth and query processing capability, but the
 proportionality constant changes with time as new generations of system bus with greater
 bandwidth become available.
- The US DAC LSP is sized at 10% of the DRP compute budget in core-hours, readjusted to be spread over an entire year. The Chilean DAC LSP is sized at 20% of the US DAC (as in LDM-138). The LSST staff LSP is sized at 10% of the US DAC.

DRAFT 17 DRAFT



7.3.2 Parameters

The key parameters in Table 26 are described below.

HSC PDR1 was executed on the NCSA verification cluster, which uses the nominal CPU. The Alert Production executes on Kubernetes nodes, which are a bit slower; to be conservative, this is neglected.

The most recent run of DRP on HSC PDR1 data is described at https://confluence.lsstcorp.org/x/WpBiB. The input data size is measured; note that the input data files are lossless-compressed. Most jobs (but not most of the time) could run on relatively small-memory machines with 24 cores and 5 GB RAM per core. The largest and longest-running jobs, however, required up to 4 times as much memory, using half or a quarter of the cores. To be conservative, we assume that half the cores were used for the large-memory jobs. The percentage of DRP core-hours that will need to execute on large-memory nodes is estimated.

Since the HSC PDR1 processing did not include several steps from the Science Pipelines Design document [LDM-151] such as image differencing and full multi-epoch characterization, the core-hours used are scaled up to the expected pipeline consumption. Note that these algorithmic adjustments are multiplicative.

The SQuaSH system reports the execution time of ap_pipe in seconds per CCD. A mean was taken over all processed CCDs, and it was assumed that each CCD is processed on a single core. These CCDs are from DECam, which is half the size of an LSST CCD, so the total time is doubled. A factor is added to account for additional steps like differential chromatic refraction compensation and false positive detection that are not well-represented in the current pipeline. Multiplying by the number of LSSTCam science CCDs gives the total number of corehours per visit.

The amount of Qserv data that can be handled by one node is estimated based on the amount of disk that can be scanned in 12 hours at an aggregate rate of 1 GB per second. (Since the Qserv data replicas are not all anticipated to be accessed at the same rate, this is a conservative estimate.)

DRAFT 18 DRAFT



7.3.3 Data Release Production

The number of nominal core-hours per TB of input data is multiplied by the precursor (HSC RC2 and DESC DC2 subset for 12 months and HSC PDR2 twice a year) and LSSTCam input data sizes (with lossless compression) to determine the total number of core-hours needed in each year. This is shown in Table 27. Approximately one-third of these core-hours need to be provided by small-memory (4-5 GB/core) machines; the other two-thirds need to come from large-memory (8-20 GB/core) machines.

Data Release Production	units	FY2020	FY2021	FY2022	FY2023/ LOY1	Notes
Precursor input size	TB	206	206	206	206	
LSSTCam visit input size	TB			319	1911	raw images / images/ visit, lossless-compressed
Precursor compute	core-hours	4.4E+06	4.4E+06	4.4E+06	4.4E+06	
LSSTCam compute	core-hours			6.8E+06	4.1E+07	
Total DRP compute	core-hours	4.4E+06	4.4E+06	1.1E+07	4.5E+07	
Alert Production	units	FY2020	FY2021	FY2022	FY2023/ LOY1	Notes
AP cores	cores			1,188	1,188	minimum necessary to keep up

Table 27: Compute needs for DRP and AP

7.3.4 Alert Production

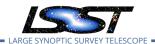
The core-hours per visit are divided by the minimum visit length (30 sec plus 1 sec shutter motion plus 2 sec readout) to give the minimum number of cores needed to keep up with image taking. This is shown in Table 27. These cores are expected to be provided over multiple "strings" of nodes. Note that the current AP design is not readily able to take advantage of more than one core per CCD.

7.3.5 LSST Science Platform

LSST Science Platform needs for US DAC science users are derived as 10% of the DRP core-hour requirement and are shown in Table 28. The LSP core-hours are assumed to be spread over a year, giving the total number of nominal cores needed in the DAC. Peak loads are expected to be handled by "borrowing" elastically from the DRP compute pool.

As a reasonableness check, the number of cores per science user is computed, but it must be noted that an oversubscription factor needs to be taken into account since not all users are expected to be simultaneously active.

Similar computations for the Chilean DAC (at 20% of the US DAC) and the LSST staff LSP (at



10% of the US DAC) are also in Table 28.

The number of Qserv nodes needed is computed from the storage devoted to it and the storage per node number. Note that staff use of Qserv is taken into account by loading the Data Release products into an internal-only Qserv instance and then making that instance part of the DAC at Data Release, so the compute sizing is part of the US DAC.

US DAC	units	FY2020	FY2021	FY2022	FY2023/ LOY1	Notes
LSP cores	cores			128	517	10% of DRP, over a year
Qserv nodes	nodes			14	95	
LSP cores/ science user	cores/ user			0.03	0.10	includes oversubscription
Chilean DAC	units	FY2020	FY2021	FY2022	FY2023/ LOY1	Notes
LSP cores	cores			26	103	20% of US DAC
Qserv nodes	nodes			14	95	
Staff LSP	units	FY2020	FY2021	FY2022	FY2023/ LOY1	Notes
LSP cores	cores			13	52	10% of US DAC

Table 28: Compute needs for the Science Platform instances

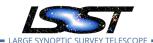
7.3.6 DES Comparison

As another check on the model, core-hour figures for Dark Energy Survey (DES) processing were obtained. These are given in Table 29. The CPUs used for single-frame and coadd processing had slightly slower clock rates but better bandwidths and expected instructions per clock performance, so they were considered equivalent to our nominal core. The CPUs used for Multi-Object Multi-Band Fitting and Single-Object Fitting (MOF/SOF) included a large contribution from the Blue Waters machine at NCSA. Those CPUs (AMD 6276) are somewhat older and were estimated at 0.245 nominal cores.

The single-frame processing measured number of 5.2 core-hours per visit compares well with the 5.4 core-hour per visit parameter used in our sizing model. Similarly, the overall DES compute figure of 21,000 core-hours per terabyte is virtually identical to our estimate (including the factors for additional steps).

DES Comparison	units			Notes
Input data size	TB	50		
Single-frame data size	TB	0.001		
Single-frame processing	core-hours/ visit	5.2		Xeon E5-2680 v4 2.4GHz
Coadd processing	core-hours/ deg20	34.7		Xeon E5-2680 v4 2.4GHz
MOF/ SOF measurement	core-hours/ deg20	108.0		AMD 6276 (313 GFLOPS/ 32 scheduled cores) and Xeon E5-2680 v4
Sky area	deg2□	5707		
DES compute per TB	core-hours/ TB	2.1E+04		

Table 29: Comparison with DES compute



7.4 Operations Sizing

Five tables use some of the parameters from the above model to project LSST storage and compute needs throughout the 10 years of Operations.

7.4.1 Storage in Operations

The Object, Source, and ForcedSource numbers in Table 30 are taken from LSE-81, as before. The number of science users and storage per user is ramped up. Note that the number of images needing storage and processing grows linearly with time. Table row sizes are taken from LDM-141; they include growth over time as columns are added.

The dataset sizes in Table 31 are calculated using the same formulas and proportionality constants as in Table 23.

The on-the-floor storage estimates in Table 32 include fast (SSD) storage for the APDB and Qserv Czar, with the latter being sized for three Data Releases (two being served and one being prepared).

"Normal" filesystem storage holds raw images, data products, scratch space, Qserv data prior to loading, science user workspace, and a 20% allocation for everything else.

Qserv local storage holds catalogs for three Data Releases.

Raw images (lossless-compressed) are written immediately to object storage, as are Parquetformat catalogs. PVIs are lossy-compressed and placed in object storage. The complete set of raw images is available, whereas the catalogs from only the last two Data Releases and the one in preparation are kept, and the PVIs from only the last Data Release and the one in preparation are online.

All data products and new raw images for each Data Release are copied to tape, but scratch space and the Qserv-schema catalogs are not.

Table 33 extracts the Qserv and object store sizing needed to populate the Chilean DAC with a copy of the data products and raw images.

DRAFT 21 DRAFT

Table 30: Inputs used to calculate storage needs during Operations

Parameters	unit	LOY1/ FY23	LOY2/ FY24	LOY3/ FY25	LOY4/ FY26	LOY5/ FY27	LOY6/ FY28	LOY7/ FY29	LOY8/ FY30	LOY9/ FY31	LOY10/ FY32
Objects	number	2.75E+10	3.25E+10	3.57E+10	3.82E+10	4.03E+10	4.22E+10	4.38E+10	4.53E+10	4.64E+10	4.74E+10
Sources	number	9.01E+11	1.80E+12	2.70E+12	3.60E+12	4.51E+12	5.41E+12	6.31E+12	7.21E+12	8.11E+12	9.01E+12
ForcedSources	number	2.91E+12	6.87E+12	1.13E+13	1.61E+13	2.13E+13	2.67E+13	3.24E+13	3.83E+13	4.41E+13	5.01E+13
Science users	users	5000	6000	7000	7500	7500	7500	7500	7500	7500	7500
Storage per science user	TB	0.4	0.5	0.6	0.7	0.8	0.9	1	1.1	1.2	1.3
LSSTCam image size	TB	0.0152	0.0152	0.0152	0.0152	0.0152	0.0152	0.0152	0.0152	0.0152	0.0152
Raw image compression	factor	0.42	0.42	0.42	0.42	0.42	0.42	0.42	0.42	0.42	0.42
Lossy image compression	factor	0.250	0.250	0.250	0.250	0.250	0.250	0.250	0.250	0.250	0.250
Observing nights per year	nights	300	300	300	300	300	300	300	300	300	300
Visits per night	visits	1000	1000	1000	1000	1000	1000	1000	1000	1000	1000
Images per visit	images	2	2	2	2	2	2	2	2	2	2
Calibration images per day	images	500	500	500	500	500	500	500	500	500	500
LSSTCam Science images	images	600000	1200000	1800000	2400000	3000000	3600000	4200000	4800000	5400000	6000000
LSSTCam Engineering images	images	6000	12000	18000	24000	30000	36000	42000	48000	54000	60000
LSSTCam Calibration images	images	150000	300000	450000	600000	750000	900000	1050000	1200000	1350000	1500000
Number of coadd data products	number	2	2	2	2	2	2	2	2	2	2
Object table row size	bytes	1896	1953	2012	2073	2136	2201	2268	2337	2408	2481
Object_Extra tables row size	bytes	21005	21636	22286	22955	23644	24354	25085	25838	26614	27413
Source table row size	bytes	467	482	497	512	528	544	561	578	596	614
ForcedSource table row size	bytes	41	41	41	41	41	41	41	41	41	41
Qserv replication factor	factor	3.0	3.0	3.0	3.0	3.0	3.0	3.0	3.0	3.0	3.0

Table 31: Dataset sizes used to calculate storage needs during Operations

Dataset Sizing	unit	LOY1/ FY23	LOY2/ FY24	LOY3/ FY25	LOY4/ FY26	LOY5/ FY27	LOY6/ FY28	LOY7/ FY29	LOY8/ FY30	LOY9/ FY31	LOY10/ FY32
LSSTCam Area	deg2□	17000	17000	17000	17000	17000	17000	17000	17000	17000	17000
APDB	TB	24	24	24	24	24	24	24	24	24	24
Object store datasets:											
Incremental LSSTCam Raw Images	TB	4816	4816	4816	4816	4816	4816	4816	4816	4816	4816
LSSTCam Output Coadd Images	TB	7727	7727	7727	7727	7727	7727	7727	7727	7727	7727
Normal disk datasets:											
LSSTCam Output Images	TB	13485	26970	40456	53941	67426	80911	94397	107882	121367	134852
LSSTCam Output Parquet	TB	7973	15946	23919	31893	39866	47839	55812	63785	71758	79731
Scratch	TB	1349	2697	4046	5394	6743	8091	9440	10788	12137	13485
Qserv Czar/ Object	TB	156	190	215	238	258	279	298	318	335	353
Qserv Database	TB	3510	5748	8018	10378	12881	15475	18199	21042	23965	27010
Science User Home	TB	2000	3000	4200	5250	6000	6750	7500	8250	9000	9750
Other/ Misc	ТВ	8208	13424	18684	23932	29148	34382	39642	44926	50226	55550

Table 32: On floor LDF storage estimates during Operations

LDF Storage (on the floor)	unit	LOY1/ FY23	LOY2/ FY24	LOY3/ FY25	LOY4/ FY26	LOY5/ FY27	LOY6/ FY28	LOY7/ FY29	LOY8/ FY30	LOY9/ FY31	LOY10/ FY32	
APDB	TB	24	24	24	24	24	24	24	24	24	24	
Qserv Czar/ Object	TB	182	347	562	643	711	774	835	894	951	1006	
Total Fast	TB	206	371	586	667	735	798	859	918	974	1029	
Normal	TB	38983	67976	99538	131025	162321	193727	225288	256991	288788	320731	П
Qserv Storage	TB	4094	9257	17275	24144	31277	38734	46555	54716	63206	72017	
LSSTCam Raw Images	TB	6982	11798	16614	21430	26246	31062	35878	40694	45510	50326	
LSSTCam Output Images	TB	3933	10676	16857	23599	30342	37084	43827	50570	57312	64055	
LSSTCam Output Coadd Images	TB	8636	16364	23182	23182	23182	23182	23182	23182	23182	23182	
LSSTCam Output Parquet	TB	9302	25248	47839	71758	95678	119597	143516	167436	191355	215275	
Object Store	TB	28854	64086	104491	139969	175447	210925	246403	281881	317359	352837	
LSSTCam Raw Images	TB	6982	11798	16614	21430	26246	31062	35878	40694	45510	50326	
All Data Products/ Backup	TB	47626	106967	194226	309242	451681	621564	818915	1043758	1296109	1575992	
All Object Store-Only Products	TB	8636	16364	24091	31818	39545	47273	55000	62727	70455	78182	
Tape	TB	63245	135129	234931	362490	517473	699899	909793	1147179	1412074	1704500	

Table 33: On floor Chile storage estimates during Operations

Chile Storage (on the floor)	unit	LOY1/ FY23	LOY2/ FY24	LOY3/ FY25	LOY4/ FY26	LOY5/ FY27	LOY6/ FY28	LOY7/ FY29	LOY8/ FY30	LOY9/ FY31	LOY10/ FY32	
Fast	ТВ	156	347	562	643	711	774	835	894	951	1006	
Qserv Storage	TB	4094	9257	17275	24144	31277	38734	46555	54716	63206	72017	
Object Store	TB	28854	64086	104491	139969	175447	210925	246403	281881	317359	352837	

Table 34: Compute needs during Operations

Data Release Production	units	LOY1/ FY23	LOY2/ FY24	LOY3/ FY25	LOY4/ FY26	LOY5/ FY27	LOY6/ FY28	LOY7/ FY29	LOY8/ FY30	LOY9/ FY31	LOY10/ FY32
LSSTCam visit input size	TB	1911	3822	5733	7644	9556	11467	13378	15289	17200	19111
DRP compute	core-hours	4.5E+07	8.2E+07	1.2E+08	1.6E+08	2.0E+08	2.5E+08	2.9E+08	3.3E+08	3.7E+08	4.1E+08
Alert Production	units	LOY1/ FY23	LOY2/ FY24	LOY3/ FY25	LOY4/ FY26	LOY5/ FY27	LOY6/ FY28	LOY7/ FY29	LOY8/ FY30	LOY9/ FY31	LOY10/ FY32
AP cores	cores	1,188	1,188	1,188	1,188	1,188	1,188	1,188	1,188	1,188	1,188
US DAC	units	LOY1/ FY23	LOY2/ FY24	LOY3/ FY25	LOY4/ FY26	LOY5/ FY27	LOY6/ FY28	LOY7/ FY29	LOY8/ FY30	LOY9/ FY31	LOY10/ FY32
LSP cores	cores	517	933	1,399	1,866	2,332	2,798	3,265	3,731	4,198	4,664
Qserv data per node	TB/ node	43	43	86	86	86	86	173	173	173	173
Qserv nodes	nodes	95	216	309	348	364	451	436	408	367	418
LSP cores/ science user	cores/ user	0.1	0.2	0.2	0.2	0.3	0.4	0.4	0.5	0.6	0.6
Chilean DAC	units	LOY1/ FY23	LOY2/ FY24	LOY3/ FY25	LOY4/ FY26	LOY5/ FY27	LOY6/ FY28	LOY7/ FY29	LOY8/ FY30	LOY9/ FY31	LOY10/ FY32
LSP cores	cores	103	187	280	373	466	560	653	746	840	933
Qserv nodes	nodes	95	216	309	348	364	451	436	408	367	418
Staff LSP	units	LOY1/ FY23	LOY2/ FY24	LOY3/ FY25	LOY4/ FY26	LOY5/ FY27	LOY6/ FY28	LOY7/ FY29	LOY8/ FY30	LOY9/ FY31	LOY10/ FY32
LSP cores	cores	52	93	140	187	233	280	326	373	420	466

DRAFT



7.4.2 Compute in Operations

The DRP compute sizing in Table 34 follows directly from the size of the input data to be processed. The number of cores for Alert Production does not change with time. The DAC and staff LSP instances are sized based on the assumed percentages of DRP compute. The amount of Qserv data that can be handled by a node is assumed to grow with time, doubling every four years (PCI Express has gone from 1.0 GB/sec to 16 GB/sec between 2003 and 2019). The number of Qserv nodes is calculated by dividing each Data Release's storage by the storageper-node figure for its year; older nodes are assumed to be retired.

References

References

- [LDM-141], Becla, J., Lim, K.T., 2013, Data Management Storage Sizing and I/O Model, LDM-141, URL https://ls.st/LDM-141
- [LSE-81], Dubois-Felsmann, G., 2013, LSST Science and Project Sizing Inputs, LSE-81, URL https: //ls.st/LSE-81
- [LDM-144], Freemon, M., Pietrowicz, S., Alt, J., 2016, Site Specific Infrastructure Estimation Model, LDM-144, URL https://ls.st/LDM-144
- [LDM-138], Kantor, J., Axelrod, T., Lim, K.T., 2013, Data Management Compute Sizing Model, LDM-138, URL https://ls.st/LDM-138
- [DMTN-113], Salnikov, A., 2019, Performance of RDBMS-based PPDB implementation, DMTN-113, URL http://dmtn-113.lsst.io
- [LDM-151], Swinbank, J.D., et al., 2017, Data Management Science Pipelines Design, LDM-151, URL https://ls.st/LDM-151
- [DMTN-091], Wood-Vasey, M., Bellm, E., Bosch, J., et al., 2019, Test Datasets for Scientific Performance Monitoring, DMTN-091, URL https://dmtn-007.lsst.io/v/DM-15448/index.html, LSST Data Management Technical Note

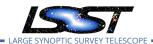
DRAFT DRAFT 24



B Acronyms

AP Alert Production APDB Alert Production DataBase CCD Charge-Coupled Device	
CCD Charge-Coupled Device	
CPU Central Processing Unit	
DAC Data Access Center	
DCR Differential Chromatic Refraction	
DDN Data Delivery Network	
DM Data Management	
DMTN DM Technical Note	
DPDD Data Product Definition Document	
DRP Data Release Production	
FLOP FLoating point Operation	
GB Gigabyte	
GPFS General Parallel File System (now IBM Spectrum Scal	e)
HSC Hyper Suprime-Cam	
IN2P3 Institut National de Physique Nucléaire et de Physiqu	ue des Particules
JSR Joint Status Review	
LATISS LSST Atmospheric Transmission Imager and Slitless S	Spectrograph
LCR LSST Change Request	
LDF LSST Data Facility	
LDM LSST Data Management (Document Handle)	
LSE LSST Systems Engineering (Document Handle)	
LSP LSST Science Platform	
LSST Large Synoptic Survey Telescope	
MB MegaByte	
NCSA National Center for Supercomputing Applications	
PDR Preliminary Design Review	
QA Quality Assurance	
RAM Random Access Memory	
RFC Request For Comment	
SATA Serial Advanced Technology Attachment	
SQuaSH Science Quality Analysis Harness	

DRAFT 25



SSD	Solid-State Disk
SSP	Solar System Processing
ТВ	TeraByte
US	United States
VM	Virtual Machine
deg	degree; unit of angle

