ARK ENERGY	_
JRVEY	

The Dark Energy Survey Data Management System

Robert Gruendl NCSA/University of Illinois

	Outline	
DARK ENERGY SURVEY		

- A Brief History of Pipelines
- Overview of the Main Pipeline Modules
- Some Lessons Learned and Challenges that LSST might still have time to avoid



Data Challenge 6B (2011) 12 months before DECam first light

DARK ENERGY SURVEY

- Reduced 10 simulated nights (relatively ideal)
- SE-pipeline (a monolith) took ~36-48 hrs to process one night (required success for all exposures in a night)
- Calibrations depended on good nightly data (e.g. nightly illumination/fringe correction)
- SNDiff pipeline needed flat with ~5-10x fidelity achievable with a single night's calibrations
- Many "science codes" (e.g. WL and photo-z) were undergoing heavy, continuous development by WGs (poor turn around of tests)



Pipeline Parallelism



Crosstalk Block: X pipeline jobs Modules in pipeline: Crosstalk

CreateCor Block: Y pipeline jobs Modules in pipeline: mkbiascor, mkflatcor Note: mkflatcor is repeated for each band

Detrend Block: Z pipeline jobs (Z >> Y) Modules in pipeline: imcorrect Example for Execution Paths



06/14/2017

DARK ENERGY

LSST School & Workshop



The Intermediate Solution

DARK ENERGY SURVEY

Bring the computation closer to the data preferably on resources with a well provisioned file system

- The original DESDM proposal was that it would use HPC resources available through TeraGrid/XSEDE centers. In reality these were not ideal for most of the processing needs
 - Always required heavy file transport
 - Typically had LUSTRE file systems (not good for heavy i/o)
- Arranged for use of iForge (part of the PSP program at NCSA)
- NCSA also stood up a GPFS storage condominium (to better serve both DES and other facility projects)

Revamp DB infrastructure

- Eliminate live replica versions at secondary sites.
 - Helped removes latency/bottlenecks when ingesting large amounts of data
- Replace web-portal DB access with light-weight tools
 - Open access for the DES Collaboration to all data



The Intermediate Solution

DARK ENERGY SURVEY

Break up the pipeline computation into native chunks with realizable goals

- SE processing was altered to work on individual exposures
 - one problematic exposure no longer halts all processing
- Pipelines responsible for building calibrations decoupled from science pipeline processing
 - Long term boon because DECam is an incredibly stable instrument
 - Removed the one-and-done mentality for calibration
 - Allows knowledge about the changing state of the instrument to filter back and inform calibration process
- FIRSTCUT → rapid feedback/QA
- FINALCUT \rightarrow best reduction possible (after calibration is understood)
- Pipelines with heavy development were temporarily off-loaded to WGs
 - Most have been re-united with DESDM processing as afterburners as of the Y3 production campaign



Current Dataflow and Pipelines





Detrending in the DES SE pipelines





The Long-Term Solution

DARK ENERGY SURVEY

Revamp Middleware

- PERL \rightarrow Python
- Allow for use of Open Science Grid resources (e.g. FermiGrid)
 - + More locations to run
 - + Use of a local filesytem where only local job(s) are competing for i/o.
 - More data transfer necessary (all calibrations to each node)
- Require unique filenames for all products
- Allow for provenance tracking (using the Open Provenance Model) so that the interdependence of products can be traced

Revamp Science Codes so that Collaboration more able to contribute

- Move from C/C++ \rightarrow Python (w/ wrapped C where needed)
- Bring DES Collaboration expertise to bear
 - More eyes, more minds (spot and understand problems/"features")
 - Provide solutions (up to and including actual code)
- EUPS \rightarrow to provide more regimented production environment
 - Also provides for stand-alone development (and even stand-alone production)



The Long-Term Solution (hardware)

DARK ENERGY SURVEY

Replace iForge with dedicated DES production nodes

- Illinois Campus Cluster (with 40 Gb/s, 3-hop to main file storage)\
- Well provisioned nodes 28 cores, 8GB/core, 12 TB local disk
- Use FermiGrid and BlueWaters for peaking

Second Generation DB Hardware

- Replaces a cobbled together system (stone-soup DB cluster)
- Adds SSDs for user and temporary table spaces

More nimble Disaster Recovery

- Spinning Backup (using JBODs composed of shingle media)
- Provides nimble than recovering from tape



Distributed Development





Pipelines: FirstCut vs. FinalCut

DARK ENERGY SURVEY

Minimal processing needed to provide data quality assessment:

• Detrending:

MIN: demonstrates that artifacts are properly understood +++: early science is possible

• Astrometry:

MIN: exposure location verified, distortion is stable +++: moving and variable object discovery

Cataloging

MIN: rough PSF analysis (FWHM, ellipticity, 2nd moments)

MIN: depth/sensitivity are adequate

- MIN: preliminary photometric solution
- +++: preliminary science, moving and variable object discovery



Pipelines: FirstCut vs. FinalCut

DARK ENERGY SURVEY

FirstCut: uses a preliminary set of calibrations and produces a result suitable to for Data Quality Assessment and feedback to the observing team.

FinalCut: uses a best set of calibrations (draws from experience and QA gained running FirstCut) and produce results that are release to DES Collaboration (and eventually the astronomical community).

By Year 1, the pipelines were streamlined to the point that software changes incorporated in FinalCut were routinely made available for FirstCut processing.





DARK ENERGY SURVEY

Raw Exposure from the telescope















Crosstalk & Overscan

DARK ENERGY SURVEY

Remove overscan:

- Currently uses line-by-line average with outlier rejection
- Functional fitting and splines would require knowledge of bias jumps for backplanes containing focus chips

Crosstalk:

 Crosstalk removal (mostly inter-ccd) but has a non-linear behavior and super-saturated sources cannot be corrected.















Detrend

DARK ENERGY SURVEY

- Bias (either nightly or super-bias)
- Linearity Correction
- Gain Correction (added for Y3)
- Brighter Fatter (Y3A1 and later)
- Flat (super-flat)
 - Y3A1 switched to normalization across focal plane to enable full focal plane sky subtraction (not shown in figure)
- Pupil/Illumcor derived from starflats
 Replaced in Y3
- Fringe (zY-bands only)
- Y3: replaced by full focal plane sky subtraction using a PCA fit including fringe
- Y3: Illumumination Correction → Star Flat 06/14/2017
 LSST School & Workshop













Intentional Dome Misalignment: ~2 mmag effect

DARK ENERGY SURVEY



Experiment by G. Bernstein and D. James

Flat Field Monitor



LSST School & Workshop



Astrometry



<u>Up through Y3A1:</u> Used SCAMP and UCAC-4 along with a single predetermined distortion correction \rightarrow 250 mas RMS

<u>Y4:</u>

Switched to GAIA-DR2 for a reference

 \rightarrow 70 mas RMS

Add per epoch distortion estimates (from star-flats) \rightarrow 25-50 mas RMS













Astrometric Distortion

DARK ENERGY SURVEY

Y1:

 update addressed band dependence issues but also identified temporal changes in astrometric distortion pattern

Y4:

- Temporal variation is now accomodated within the pipelines.





Astrometry (longer term)

DARK ENERGY SURVEY

Current (Y5) effort will attempt to refine relative astrometry still further/ Tests show systematic residuals of 50-80 mas can be modeled and removed.









Bleed trail identification by searching for extended structures stemming form saturated island.

Y1 additions:

- Mask dilation in the cross-trail direction to better remove strong bleeds.
- Search for edge-bleed conditions for trails that intercept the read registers.











06/14/2017





Bleed trail identification by searching for extended structures stemming form saturated island.

Y1 additions:

- Mask dilation in the cross-trail direction to better remove strong bleeds.
- Search for edge-bleed conditions for trails that intercept the read registers.











06/14/2017



CR-reject & streak finder

DARK ENERGY SURVEY

- Early CR-rejection used neural net identification (only partially effective).
- SV: Single-Image CR-rejection was via gradient (better)
- Y1: Implemented LSST-stack CR-rejection algorithm within DESDM pipelines.
- Streak finder deployed in Y1 uses identification via Hough transform
- Y4: adding a truth-table/testbed to investigate further improvements for Y5











06/14/2017



Exposure Based Assessment

DARK ENERGY SURVEY

Evaluate each reduced exposure based on SE products. Primary goal is to determine whether each observation meets basic survey requirements/ standards.

Primary decision based on the effective exposure time:

 $T_{eff} = (0.9 \text{ k / FWHM})^2 (Bkgd_{dark} / Bkgd) (10^{-2} \text{ cloud / 2.5})$ $= F_{eff} B_{eff} C_{eff}$

Current cutoffs are $T_{eff} > 0.3$ (riz-bands), > 0.2 (gY-band)

Further cuts can be placed based on individual components or other QA (e.g. astrometry, PSF) to form input TAG for COADD or other analysis.

Coverage Check Prior to COADD



06/14/2017

LSST School & Workshop



Projected Y1A1 Footprint



06/14/2017



Y3A1 COADD Footprint (Y1+Y2+Y3)

DARK ENERGY SURVEY



06/14/2017



Y3A1 COADD Footprint (DEEP)

DARK ENERGY SURVEY





COADD

DARK ENERGY SURVEY

COADDITION of single-epoch images requires a global calibration based on single epoch photometry.

- Y1: GCM \rightarrow 2-3 percent
- Y3: FGCM → sub-percent accuracy

Y1 added astrometric refinement reduce the relative astrometric residuals to ~50 mas.

Y3 added:

- pipeline generation of MEDs products (ties COADD detection to single-epoch pixels)
- Afterburner pipelines use MEDs for:
 - MOF: multi-epoch, -object,-band photometry
 - Ngmix: shear estimates for WL









Some Lessons Learned

DARK ENERGY SURVEY

DES has been extremely successful but there was a vast amount of improvement that was needed.

- 1. Do not assume that you can have the final word in calibration at the time an observation occurs
- 2. Simulations are good but real data are better
- 3. Build QA metrics that can identify observations taken in poor conditions so that pipelines have a means to exit gracefully
 - Record the occurrence so that you can easily re-identify problematic data
 - FAILURE should be an option
- 4. Build pipelines so that lessons learned in DRP can influence subsequent nightly processing
- 5. Do not completely eliminate the ability to use human judgment to determine what data will go through DRP
 - DES had to add both TAGs and BLACKLISTs to identify/remove problematic observations above and beyond data flagged by QA metrics



DARK ENERGY SURVEY

- 6. Feedback from users may eventually help but it cannot be relied on in the short-term.
 - DES's next DRP campaign generally begins just as the collaboration is beginning to dig into the data.
- 7. Do not underestimate the value of carrying provenance with the production
 - Allows a means to trace the the downstream impact of a problematic (exposure/ cal/code)
- 8. The perfect can be the enemy of the good
 - Expect that automation and QC assessment will mature over time
 - Afterburners can be your friend
- 9. "Data management" should not be an oxymoron



Unspoken Thoughts

DARK ENERGY______SURVEY