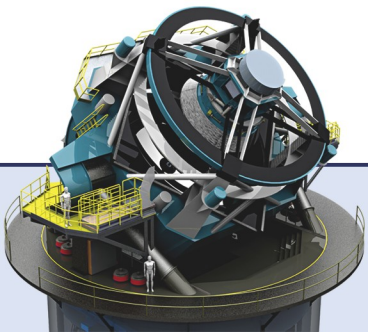# Data Management Status

## Jim Bosch
## Princeton University

# What Data Management has been up to (that's relevant for DESC).

# What Data Management and DESC can do for each other.

# Planning

1. Update design documents

   – Data Products Definition Document (LSE-163): outputs

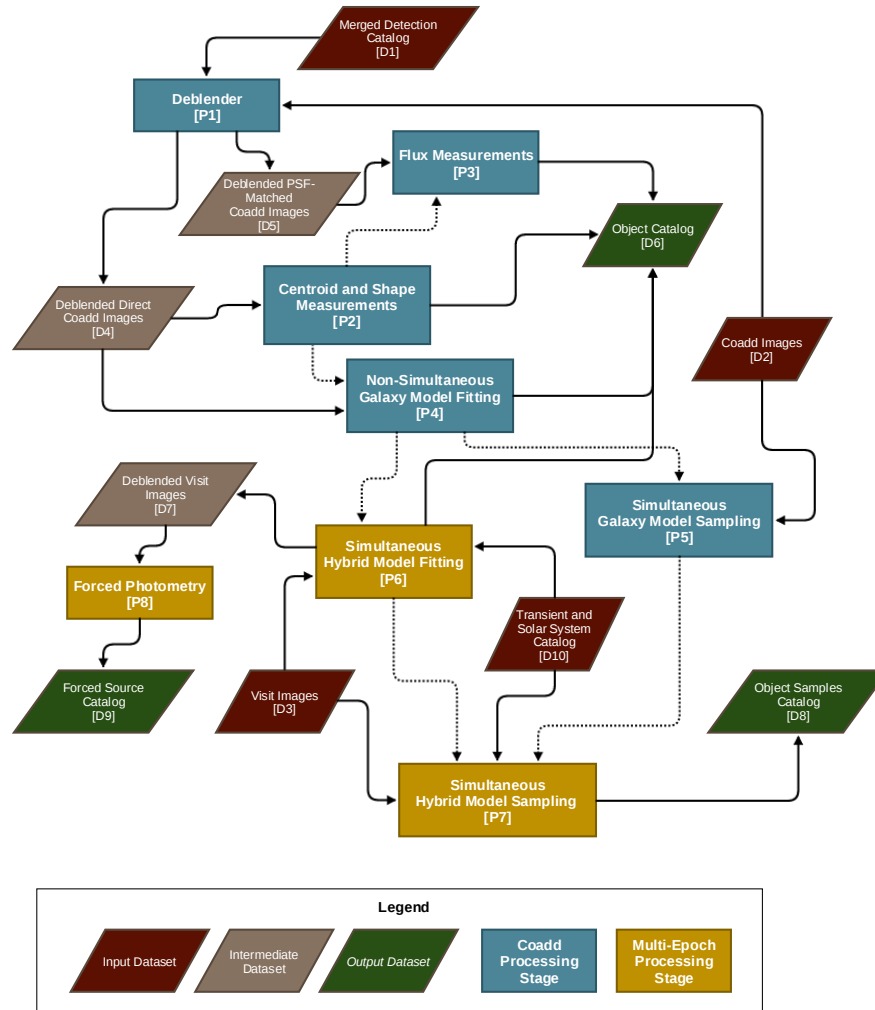   – DM Applications Design (LDM-151): algorithms

2. Identify some intermediate goals ("Define Milestones")

   – Some of these allow different teams within DM to know when they can expect a dependency to be ready.

   – Some of these allow external users (like DESC) to know when to expect a feature.

3. Figure out how long it will take ("Resource Load")
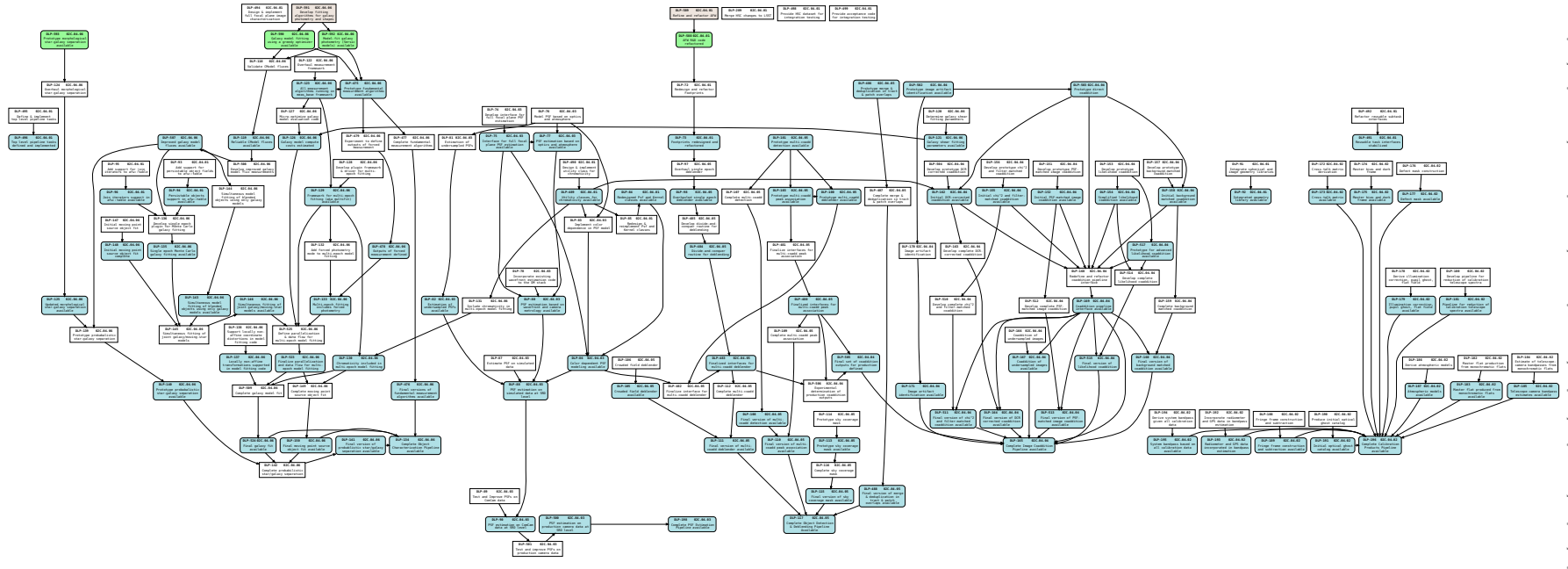
   – Do we have enough people?
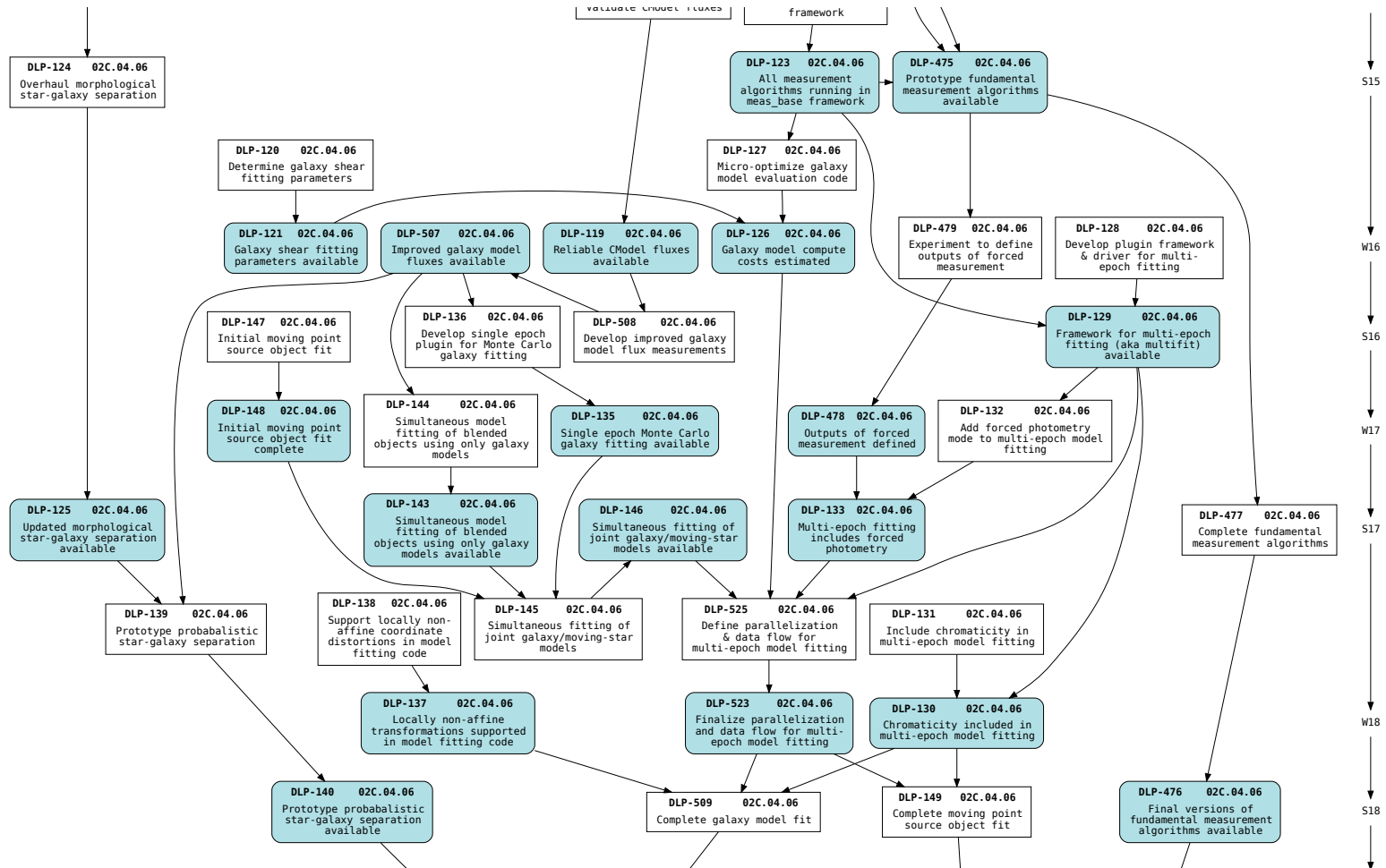
# Flowcharts for the Algorithms

Level of detail in design docs: 1-2 paragraphs for each of these boxes.

This is only about 1/6 of the data release production pipeline (not including middleware).

# Flowcharts for Development

A new version of the full plan, with updated design documents and resource-loaded milestones, ready by mid-November *is not going to happen.*

We may have a rough plan with a good outline in the design documents by then.

# What DM is Up To: Hiring and Onboarding

- We're perhaps a bit more than half of the way to fully staffed.

- We're still missing a few lead positions.

- Still training lots of new developers.

DM Bootcamp a few weeks ago. Materials may be useful
https://community.lsst.org/t/dm-boot-camp-announcement/249

Some new names/roles DESC members should know:

- K-T Lim is now Project Engineer

We're a three-headed dragon now!

- Tim Jenness is Deputy Architect

Hate our APIs? Tim wants to fix that.

- David Nidever is Science/Data Quality Scientist

- Jonathan Sick is documentation tsar

Developing plan for verification datasets and production runs.

Docs are finally *someone's* top priority.

# State of the DM Stack: Current Work

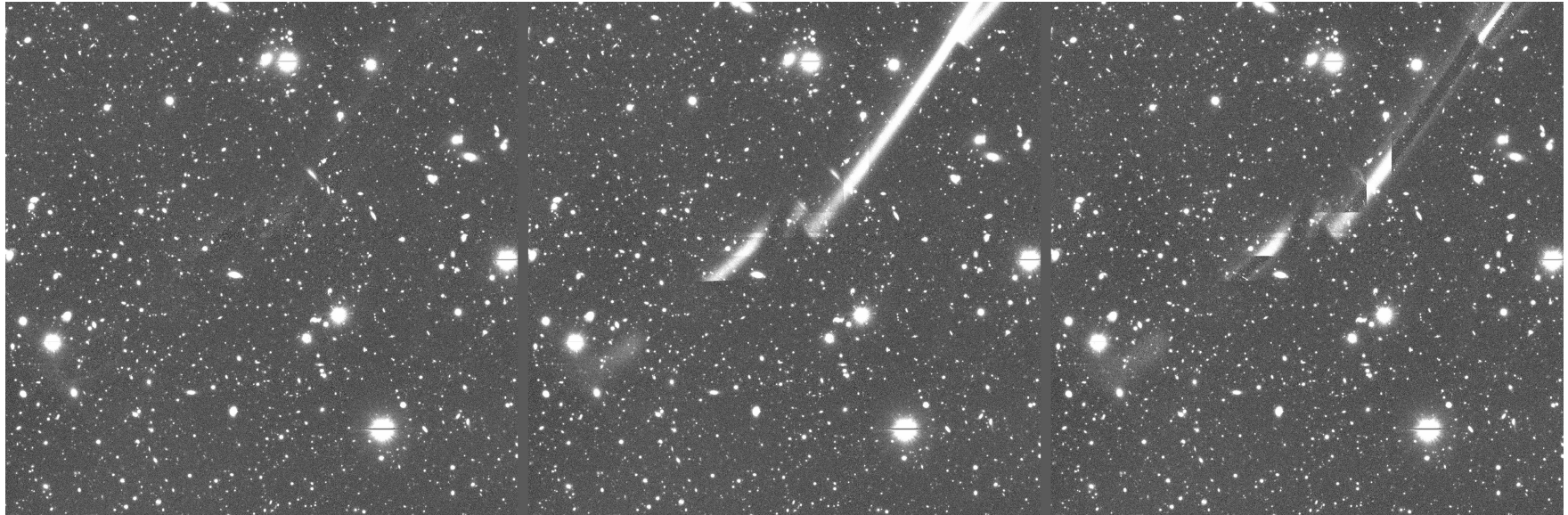About 2 years ago, the HSC pipeline team forked the LSST stack.  Since then:

- The HSC stack has become a state-of-the-art end-to-end pipeline with high-quality (for Stage II) algorithms for HSC data reduction – at the expense of a lot of confusing, poorly-documented code.

- The LSST stack has been refactored and improved considerably in terms of code quality, our testing and workflow have been greatly enhanced, and we've made it work on more cameras.  But algorithms are largely unchanged, or in some cases have bitrotted.

## Now we're putting them back together.

And the code release from September is about halfway there.

# State of the DM Stack: HSC Fork

- We're making the best coadds in the world.

- We detect and deblend consistently in all bands.

- PSF modeling uses a modified *PSFEx*; already meeting LSST minimum design specs for residual ellipticity correlations.

- We have a basic correction for brighter-fatter (~90% corrected).

- Moderately robust galaxy model fluxes and colors.

- *No multi-epoch fitting.*

- *No PSF-matched coadds.*

- *No modern shear estimation (just Regaussianization).*

# "Safe Clipping" in Coadds



safe clipping          direct mean          3σ clipped mean

- build mean and clipped mean coadds
- difference them, find above-threshold regions
- reject whole regions from individual images
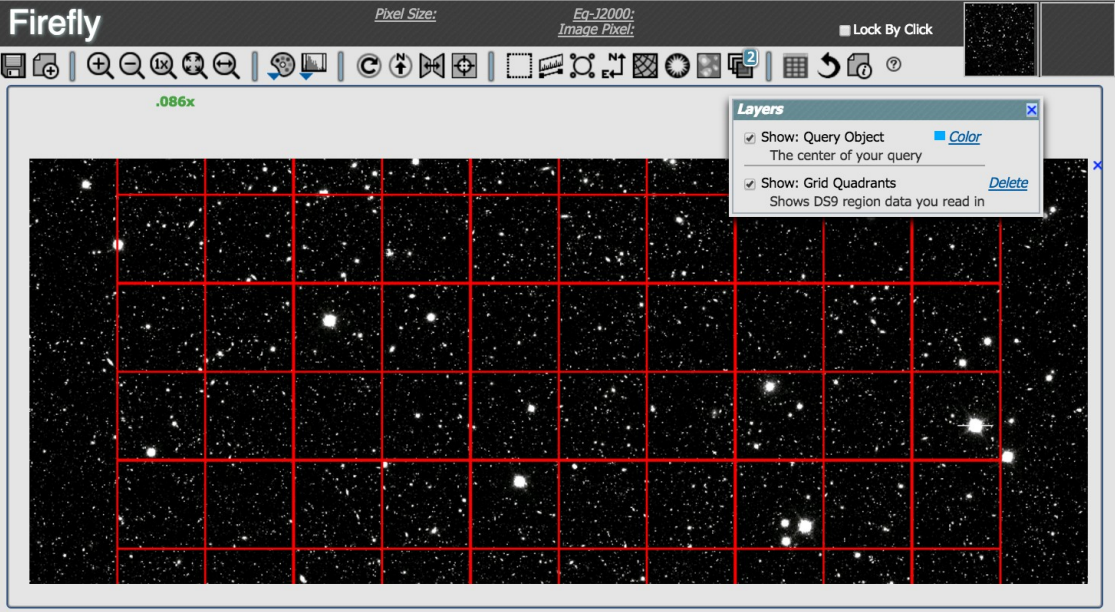
courtesy of Bob Armstrong

# State of the DM Stack: SUI



User interface team at IPAC is experimenting with embedding their `Firefly` image display tool inside IPython notebooks.

# State of the DM Stack

- Pierre Astier and Dominique Boutigny are developing a new relative astrometry and photometry module that's now being integrated into the pipeline; already working on CFHT and HSC data, with DECam coming soon.

- We're putting a lot of work into getting the stack running well on DECam data.

- We're working on the image subtraction code again, focusing on DCR.

- We've greatly improved our continuous integration and automated testing.

Main goals: regular, complete end-to-end tests on real data and improvements that make developers work more efficiently.

What Data Management has been up to (that's relevant for DESC).

**What Data Management and DESC can do for each other.**

# Introducing **community.lsst.org**



- Preferred over HipChat for long-running, in-depth discussions.

- Most non-announce email lists are being phased out.

- Use/search the Q&A tag for StackOverflow-style FAQs.

- Intended for users, not just DM/Project people (self-signup).

# Verification Datasets and Productions

We've started thinking about what real datasets we can use for development and verification before commissioning. No firm plans yet, but we're considering many options:

- DECam – Start with general observer datasets whose owners want to try our code to do their science, *eventually* reprocess significant fraction (all?) of DES.

- HSC – Amount of public data will increase dramatically in the next year.

- CFHTLS – best for early tests of shear estimation, due to publicly available shear catalogs to compare against.

- PanSTARRS – mostly database and user interface testing (probably no reprocessing).

# Should you use the DM Stack?

The DM stack is essentially state-of-the-art from an algorithms standpoint – it's about as good as the PanSTARRS, SDSS, and DES pipelines in most respects.  It's clearly not as good in some areas (but frequently better in others).

Documentation and ease-of-use are in bad shape compared to AstroPy.  They're in pretty good shape compared to any of the above pipelines.

> After a few months you'll mostly just be annoyed, and only occasionally frustrated.

Expect to get a lot out of it, but expect to be very frustrated at first.

# Feedback and Contributing to the Stack

Please tell us what you don't like.  Don't assume you just don't know enough to comment intelligently.  This *absolutely* includes changes to bad interfaces.

If you can, create a GitHub PR to fix it (but ask about it before you do a lot of work).

The DM team has *spirited internal debates* about what we think our users want.  Please, please just tell us.

# The DM/DESC Boundary

If you don't know if DM is responsible for producing something, *please ask*.  We may or may not know, but we also might not have even thought about it.

*Please don't* assume you shouldn't work on a problem just because DM is responsible for it.  We want algorithmic input from outside, both in ideas and in code.

*Please do* talk to us about those algorithms early and often, so we don't duplicate effort or spend time on solved problems.

*Please do* consider writing code directly within the DM stack.  You don't have to be a DM member to contribute to the codebase.

# Design Document Development

What DM is planning in detail for algorithms and data products should become more much more clear in the next few months.

If you're particularly interested, you can also follow our design documents on GitHub and comment on the pull requests:

- `https://github.com/lsst/LSE-163` (data products)

- `https://github.com/lsst/LDM-151` (algorithms)

For now, take the content of these documents with a large bucket of salt.

# Summary

Planning and design document work will give you all a lot to read (soon).

Verification dataset runs will enable a ton of new DESC investigations (soon).

We've turned the corner on hiring – contributions from new people now outweigh the hiring/onboarding effort.

DM stack is steadily improving: still a big learning curve, but it's really a complete pipeline.

Join community.lsst.org!  Ask questions!  Contribute!