

Yoda Technical Details and Current State

J. Taylor Childers (Argonne)

https://github.com/PanDAWMS/panda-yoda



Yoda



Shared File System

Oct 2018 ATLAS S&C Week J. Taylor Childers

Queue message



MPI Communication





Yoda Work Manager

-Processing Requests -Caching Requests -Trigger Requests to Harvester -Send Jobs, Event Ranges to Droids

Queue message

MPIService

MPI Communication



Yoda Module Design Details

- Yoda-Droid is a server/client model, meaning Droids make requests to Yoda and Yoda fulfills those request.
 - The reason for this choice is to avoid needing a database for tracking the current state of each Droid.
 - Since our processing rate is on the scale of minutes per event, we don't expect Yoda to be deluged by incoming requests.
- Yoda launches the WorkManager, which processes all the incoming messages from Droids and triggers the requests from jobs or event ranges from Harvester. The Harvester interaction is treated as a plugin, currently using shared file system for communication, but could move to communication model.
 - Shared File System communication in this case is small and since only Rank 0 (Yoda) is writing/reading it is not a bottleneck.
 - However, it does create race conditions. •
 - Should move to communication where possible, just need to implement the plugins for Harvester & Yoda.
- Yoda can be given an wall-clock time limit. It will run for the specified number of minutes, then kill all AthenaMP processes and exit cleanly.
 - Needed for supporting running AthenaMP in temporary storage locations. •
 - Athena logs must be saved for posterity.
- Yoda launches a FileManager process that receives details from Droid ranks on files produced, which are then communicated to Harvester for stage out.

Oct 2018 ATLAS S&C Week J. Taylor Childers



Note about MPIService Module

- All MPI traffic is handled in the MPIService
- This was driven because I have seen unpredictable behavior when MPI was spread out among Python Threads and Processes.
- Therefore all subprocesses use multiprocessing.Queue communicate.
- Messages are always in the form of a pickled python dictionary.



objects to





Oct 2018 ATLAS S&C Week J. Taylor Childers

Queue message



MPI Communication



Droid Module Design Details

- Droid sends a request to Yoda for a Panda Job definition, after which it launches a TransformManager sub-process which runs AthenaMP on the local node, in or out of a container as needed.
 - Containers are important for scaling the Atlas Software.
 - Provide a 7x speed-up when performing many meta-data operations.



Spencer Williams ANL Summer Student

Oct 2018 ATLAS S&C Week J. Taylor Childers



Droid Module Design Details

- Droid sends a request to Yoda for a Panda Job definition, after which it launches a TransformManager sub-process which runs AthenaMP on the local node, in or out of a container as needed.
 - Containers are important for scaling the Atlas Software.
 - Provide a 7x speed-up when performing many meta-data operations.
- TransformManger simply monitors the AthenaMP process to ensure it is alive.
 - There is a lot to configure here, see slide on TransformManger
- Once a job has begun, Droid launches the JobComm sub-process which handles the YAMPL communication with the AthenaMP workers, thus providing work.
- JobComm manages the event ranges received from Yoda
 - sending requests to Yoda when more are needed
 - responding to AthenaMP workers with event ranges when requested
 - forwarding output file information to Yoda from AthenaMP workers



TransformManager Module

- Example for Singularity submission is here.
- Entries from PandaJob dictionary formatted into template to fill things
- Using this template (following Harvester's exam enables local sites to configure their runt environment.
 - Notice the DB settings.
- Current templates can be a bit Sim focused and need to be made more flexible as we include n workflows, Gen/Reco.
- Or provide a way to change the template based transforms.
- This assumes Event Service would be ported to or transforms.

	1	#!/bin/bash
	2	<pre>echo [\$SECONDS] Start inside Singularity</pre>
	3	<pre>echo [\$SECONDS] DATE=\$(date)</pre>
	4	WORKDIR=\$1
1.	5	<pre>echo [\$SECONDS] WORKDIR=\$WORKDIR</pre>
this	6	
	7	cd \$WORKDIR
	8	
nple)	9	<pre>export ATHENA_PROC_NUMBER=128 # AthenaMP workers per node</pre>
'P'C'	10	
time	36	echo [\$SECONDS] Setting up Atlas Software
	37	RELEASE={release}
	38	PACKAGE={package}
	39	CMTCONFIG={cmtConfig}
TAT:11	40	echo [\$SECONDS] RELEASE=\$RELEASE
VV 111	41	echo [\$SECONDS] PACKAGE=\$PACKAGE
nore	42	<pre>echo [\$SECONDS] CMTCONFIG=\$CMTCONFIG</pre>
	49	DBBASEPATH=\$ATLAS_DB_AREA/DBRelease/current
	50	<pre>export CORAL_DBLOOKUP_PATH=\$DBBASEPATH/XMLConfig</pre>
1 on	51	<pre>export CORAL_AUTH_PATH=\$DBBASEPATH/XMLConfig</pre>
	52	<pre>export DATAPATH=\$DBBASEPATH:\$DATAPATH</pre>
	53	mkdir poolcond
.1	54	<pre>export DBREL_LOCATION=\$ATLAS_DB_AREA/DBRelease</pre>
other	55	<pre>cp \$DBREL_LOCATION/current/poolcond/*.xml poolcond</pre>
	56	export DATAPATH=\$PWD:\$DATAPATH
	57	unset FRONTIER_SERVER
	70	<pre>echo [\$SECONDS] Starting transformation</pre>
	71	{transformation} {jobPars}
	72	echo [\$SECONDS] Fransform exited with return code: \$?
	73	<pre>echo [\$SECONDS] Exiting</pre>

https://github.com/PanDAWMS/panda-yoda/blob/master/templates/ThetaSubmitSingularity.sh

TransformManager Module

- When running on temp file systems, like RAM-disk or node-local SSDs, Yoda must clean up and stage out log files.
- To achieve this the user specifies the run path for the AthenaMP application.
- For instance, RAM-disk on Theta is '/tmp'
- The user also specifies BASH Globbing strings to specify which logs to target.
- For instance,

"log.EVNTtoHITS, athenaMP-workers-*/*/AthenaMP.log"





https://github.com/PanDAWMS/panda-yoda/blob/master/pandayoda/droid/TransformManager.py

TransformManager Module

- Sites sometimes need local modifications to tl PandaJob dictionary in order to operate properly
- For instance, the Titan operators needed to remove portion of the Athena command line options to get the database to work properly.
- Job mod is a plugin which takes the PandaJ dictionary, can modify it, and return the modified version for Yoda to run. 220

Oct 2018 ATLAS S&C Week J. Taylor Childers

	49	[TransformManager]	
he	50	loglevel	= INFO
	51	loop_timeout	= 60
	52	template	<pre>= /lus/theta-fs0/projects/AtlasADSP/atlas/harvester/templates/ThetaSubmitSing</pre>
	53	run_script	= job_submit.sh
	54	<pre>subprocess_stdout</pre>	<pre>= transform_stdout_r{rank}_pid{PandaID}.log</pre>
	55	<pre>subprocess_stderr</pre>	<pre>= transform_stderr_r{rank}_pid{PandaID}.log</pre>
	56	use_clean_env	= true
a	57	use_container	= true
1	58	<pre>container_prefix</pre>	= /usr/bin/singularity exec -B /lus/theta-fs0/projects/AtlasADSP:/lus/theta-f
he	59	run_directory	= /tmp
	60	run_elsewhere	= false
	61	logs to stage	<pre>= log.EVNTtoHITS.athenaMP-workers-*/*/AthenaMP.log</pre>
	62	jobmods	<pre>= no_file_validation,shorten_input_evnt_file_list,remove_dbrelease</pre>
h	60		

https://github.com/PanDAWMS/panda-yoda/blob/master/pandayoda/droid/TransformManager.py

```
def apply_jobmods(self):
      apply job mods specified by the configuration settings '''
   package_name = 'pandayoda.jobmod.'
   for mod in self.jobmods:
      module_name = package_name + mod
      logger.info('importing jobmod: %s', module_name)
      try:
         local_mod = importlib.import_module(module_name)
      except ImportError:
         logger.exception('Failed to import jobmod: %s\n curret python path: %s',module_name,sys.path)
         raise
      try:
         self.job_def = local_mod.apply_mod(self.job_def)
      except Exception:
         logger.exception('Failed to execute jobmod: %s',module_name)
         raise
```





TransformManag	er Module	https://github.com/Pa	anDAWMS/panda-yoda/bl	lob/master/p	andayoda/yoda_template.cf
 Sites sometimes need loca PandaJob dictionary in order f For instance, the Titan operation of the Athena commandatabase to work properly. 	<pre>49 [TransformManager] 50 loglevel 51 loop_timeout 52 template 53 run_script 54 subprocess_stdout 55 subprocess_stderr 56 use_clean_env 57 use_container 58 container_prefix 59 run_directory 60 run_elsewhere 61 loos to stage 62 jobmods</pre>	<pre>= INF0 = 60 = /lus/theta-fs0/projects/AtlasADSP/atlas/harvester/templates/ThetaSubmitSin = job_submit.sh = transform_stdout_r{rank}_pid{PandaID}.log = transform_stderr_r{rank}_pid{PandaID}.log = true = true = true = /usr/bin/singularity exec -B /lus/theta-fs0/projects/AtlasADSP:/lus/theta- = /tmp = false = log.EVNTtoHITS.athenaMP-workers=*/*/AthenaMP.log = no file validation.shorten input evnt file list.remove dbrelease</pre>			
Jol 📮 PanDAWMS / panda-yoda		O Unwatch ▼	6 🖈 Star 0 😵	Fork 2	
Code I Issues 8 1 Pull request	s 0 III Projects 0 III Wiki III	nsights 🔅 Settings			dayoda/droid/TransformMana
Branch: master - panda-yoda / pandayod	Branch: master panda-yoda / pandayoda / jobmod / final jtchilders adding plugin capability for job definition modification			History on Aug 15	
■initpy	adding plugin capability for job definition me	odification	a m	onth ago	
no_file_validation.py	ding plugin capability for job definition modification		a month ago		
remove_dbrelease.py	adding plugin capability for job definition me	odification	a m	onth ago	<pre>python path: %s',module_nam</pre>
shorten_input_evnt_file_list.py	adding plugin capability for job definition me	odification	a m	onth ago	
unlimited_athenamp_messages.py	adding plugin capability for job definition me	odification	a m	onth ago	
Oct 2018 ATLAS S&C Week J Taylor Childers	237 238	<pre>logger.exception(' raise</pre>	Failed to execute jobmod	: %s',module	_name)





Custom Python Multiprocessing Module

- After investigating some 'bug' reports from Wen, I that while the Python version inside the A environment claims to be 2.7 it contains some me from 2.6 which are very different in implementatio
- To cope with this, I wrote the yoda_multiproce module to obfuscate this difference.
- Must set environment variable RUN_YODA_IN_A
 to get this custom module.
- Otherwise it is just a pass through for all standard
- Ideally somewhere down the road we can delete th only use standard Python releases

https://github.com/PanDAWMS/panda-yoda/blob/master/pandayoda/common/yoda

Oct 2018 ATLAS S&C Week J. Taylor Childers

found	1	import os
	2	<pre>from multiprocessing import *</pre>
AILAS	3	
odulas	4	<pre># running inside the Athena release, at least for 21.0.15,</pre>
ouures	5	# the python version can be strange. While it reports Python :
n.	6	# some behaviors are like 2.6.
.	7	<pre># Especially in multiprocessing module</pre>
ssing	8	# This file tries to masks these differences by overridding t
	9	<pre># module objects if the environment variable 'RUN_YODA_IN_ATH " '</pre>
	10	# 1S SET.
THENA	11	
	12	I DUN VODA TH ATUENAL in an anuinen.
	13	<pre>IT 'RUN_YUDA_IN_AIHENA' in os.environ: print(!Dunning with custom multipressessing module!)</pre>
	14	from multiprocessing import synchronize
	15	# for the Event object in 2.6 the Event wait(timeout-None
nis and	17	# always returns None
	18	# diways returns None # for Python 2 7+ it returns True if the internal flag is
	10	# it returns False otherwise. However the Python version 2
	20	# that comes with Athena has the 2.6 behavior
	20	# chat comes with Athena has the 210 behavior
	22	<pre>class MyEvent(synchronize.Event):</pre>
	23	# overload the wait to behave like the 2.7+ version
	24	<pre>def wait(self,timeout=None):</pre>
	25	<pre>super(MyEvent, self).wait(timeout)</pre>
	26	<pre>return self.is set()</pre>
	27	
	28	# override multiprocessing Event class with my version
	29	Event = MyEvent
_multiprocessing.p	y 30	else:
	31	<pre>print('Running with standard multiprocessing module')</pre>
	2.0	



.7

- Instructions are on the github, people should send me their notes if there are hints for specific cases where the install needs to change.
- Typically I recommend installing inside a Python virtualenv.



https://github.com/PanDAWMS/panda-yoda

E README.md

panda-yoda

Requirements

- python 2.7+ (haven't actually tested with Python 3)
- mpi4py
- yampl
- python-yampl

Installation

If you are using a non-standard MPI library, which is common for most supercomputers, you will need to compile the mpi4py module against your local MPI libraries. In order to compile the mpi4py module you will need to compile the Cython module which mpi4py uses to wrap the C++ MPI library functions in Python.

In all cases you will need to compile the yampl program and the python-yampl wrapper for Yoda to use during communications with AthenaMP.

Here are rough instructions for this process, there are specific instructions for some sites further below.

Generic Instructions

```
# setup environment
# typically I setup the virtualenv of Harvester at this point
# if your HPC system uses modules to setup environments you may
# need to source some modules for an older GNU or move from
```



- Instructions are on the github, people should send me their notes if there are hints for specific cases where the install needs to change.
- Typically I recommend installing inside a Python virtualenv.
- Build Cython, mpi4py, yampl, and python-yampl inside environment
- pip install yoda

https://github.com/PanDAWMS/panda-yoda

<pre># Custom build Cython w git clone git@github.co cd cython git checkout tags/0.25. # here cc and CC are pl # local system requires CC=cc CXX=CC python set export PYTHONPATH=\$VIRT</pre>		
<pre># this installs the py # as harvester, etc. CC=cc CXX=CC python se</pre>	<pre># Custom build mpi4py with local compilers for worker nodes git clone git@github.com:mpi4py/mpi4py.git cd mpi4py git checkout tags/2.0.0 # later tags may work (haven't tested) ## edit mpi.cfg ## under [mpi]</pre>	
<pre># next you need to ins git clone git@github.c cd yampl # on some systems like # if you find it causi ./configureprefix=\$ make -j 10 install # in my case this fail # so go into the zerom cd zeromq</pre>	in this example	
<pre>./configureprefix=\$ make -j 10 install cd # continue with yampl make -j 10 install</pre>	<pre># now you need to install the python bindings for yampl git clone git@github.com:vitillo/python-yampl.git cd python-yampl # then run python setup.py build_ext # the last command may fail, if so, # edit setup.py such that inside the 'include_dirs' inside the # both the path to the '\$VIRTUAL_ENV/install/include/yampl' # and '\$VIRTUAL_ENV/install/include' and the 'extra_link_args' # '-L/path/to/yampl/install/lib'. Then run it again # Run this to install the module python setup.py install</pre>	e 'Extension' have has

install yoda into the \$VIRTUAL_ENV/python2.7/site-packages pip install git+git://github.com/PanDAWMS/panda-yoda



- Configuration includes a configuration file with many settings, all documented in the README.md
- configuration is divided by python module configparser
- many things can be tuned, how many events Yoda should request from Harvester, or Droid from Yoda.
- log filenames, template names, log levels... and so on.

https://github.com/PanDAWMS/panda-yoda

1	[yoda_droid]		
2	loop_timeout	= 60	
3	wallclock_expiri	ng_leadtime = 300	
4			
5	[Yoda]		
6	loglevel	= INFO	
7	loop_timeout	= 60	
8	messenger_plugin	_module = pandayoda.yo	oda.shared_file_messenger
9			
0	[WorkManager]	[shared file messenger]	
1	loglevel	loglevel	= TNFO
2	loop_timeout	harvester config file	<pre>= /lus/theta_fs0/projects/AtlasADSP/atlas/harvester/etc/panda/panda_har</pre>
3	send_n_eventrang	hurvester_conrig_rite	
4	request_n_eventr	[FileManager]	
.5		loglevel	= TNFO
.6	[RequestHarveste	loop timeout	= 60
.7	loglevel	output file type	
.8	loop_timeout	barvester output timeout	= 10
9		harvester_output_timeout	- 10
0	[RequestHarveste	[Droid]	
1	loglevel	loglevel	= TNFO
2	loop_timeout	loon timeout	= 60
3	eventrange timeo		
4	[TransformMana	iger]	
	loglevel	= INFO	
	loop_timeout	= 60	
	template	= /lus/theta-fs0/pro	jects/AtlasADSP/atlas/harvester/templates/ThetaSubmitSingularity.sh
	run_script	= job_submit.sh	
	subprocess_sto	<pre>lout = transform_stdout_r</pre>	{rank}_pid{PandaID}.log
	subprocess_sto	err = transform_stderr_r	{rank}_pid{PandaID}.log
	use_clean_env	= true	
	use_container	= true	
	<pre>container_prefix = /usr/bin/singularity</pre>		ty exec -B /lus/theta-fs0/projects/AtlasADSP:/lus/theta-fs0/projects/AtlasA
	run_directory	= /tmp	
	run_elsewhere	= false	
	logs_to_stage	= log.EVNTtoHITS,ath	enaMP-workers-*/*/AthenaMP.log
	jobmods	= no_file_validation	,shorten_input_evnt_file_list,remove_dbrelease
	[NDTConvice]		
	[MPIService]		
	debug message	char length = 200	
	default_message_	$r_{\rm cond} = 200$	
	loop timeout	- 20 Je_Duilei_STS6 = 10000000	
	coop_cimeouc	= 50	





- Yoda should not require any coding unless you need custom 'jobmods'
- Otherwise, the template run script for AthenaMP is the major piece.



https://github.com/PanDAWMS/panda-yoda

1	[yoda_droid]		
2	loop_timeout	= 60	
3	wallclock_expiri	ng_leadtime = 300	
4			
5	[Yoda]		
6	loglevel	= INFO	
7	loop_timeout	= 60	
8	messenger_plugin	_module = pandayoda.yo	oda.shared_file_messenger
9			
0	[WorkManager]	[shared file messenger]	
1	loglevel	loglevel	= TNFO
2	loop_timeout	harvester config file	<pre>= /lus/theta_fs0/projects/AtlasADSP/atlas/barvester/etc/papda/papda_bar</pre>
3	send_n_eventrang	harvester_config_fite	
.4	request_n_eventr	[FileManager]	
.5		loglevel	= TNEO
.6	[RequestHarveste	loop timeout	= 60
.7	loglevel	output file type	
.8	loop_timeout	harvester output timeout	= 10
9		harvester_output_timeout	- 10
0	[RequestHarveste	[Droid]	
1	loglevel	loglevel	= TNFO
2	loop_timeout	loop timeout	= 60
3	eventrange timeo	coop_clineou c	
4	[TransformMana	ager]	
	loglevel	= INFO	
	loop_timeout	= 60	
	template	= /lus/theta-fs0/pro	jects/AtlasADSP/atlas/harvester/templates/ThetaSubmitSingularity.sh
	run_script	= job_submit.sh	
	subprocess_std	<pre>iout = transform_stdout_r</pre>	{rank}_pid{PandaID}.log
	subprocess_std	<pre>lerr = transform_stderr_r</pre>	{rank}_pid{PandaID}.log
	use_clean_env	= true	
	use_container	= true	
	container_pref	ix = /usr/bin/singulari	ty exec -B /lus/theta-fs0/projects/AtlasADSP:/lus/theta-fs0/projects/AtlasA
	run_directory	= /tmp	
	run_elsewhere	= false	
	logs_to_stage	= log.EVNTtoHITS,ath	enaMP-workers-*/*/AthenaMP.log
	jobmods	<pre>= no_file_validation</pre>	<pre>,shorten_input_evnt_file_list,remove_dbrelease</pre>
	[MPIService]		
	loglevel	= INFO	
	debug_message_	_char_length = 200	
	default_messag	e_buffer_size = 10000000	
	loop_timeout	= 30	





Performance on Theta



- Yoda has been running continuously since about May on Theta with Harvester.
- I have run scaling test of Yoda (outside of Harvester) up to 1024 nodes of Theta with very good performance inside a container.
- Grafana Link

9/16





18



Performance on Theta/NERSC

- One additional Study: with/without running AthenaMP work directory in RAM-disk
- No difference in runtime. Implies we are not limited by IO at this time.
- So are we limited by Memory? AthenaMP uses ~ 1.5GB + 0.25GB * (# AthenaMP workers)
- KNL has 16GB high bandwidth memory (MCDRAM) which has 4x the bandwidth as the 192GB DDR.
- Can run using 16GB in 'L3-cache' mode. If we are cache limited, reducing number of workers to fit AthenaMP inside 16GB (so < 32 workers) we should see improved performance. However, we do not see a big improvement at that step. So
- that is not a problem.
- Most likely remaining issue is L1/2 cache thrashing. Each worker requires different codes at any given time, therefore we never keep the same code in L1 cache for very long.
- Breaking up our event processing into smaller pieces may help us gain some performance by ensuring we limit the amount of code loaded on a CPU at any given time.

Oct 2018 ATLAS S&C Week J. Taylor Childers



To Dos:

- Move from shared_file_messenger to communication layer for Harvester interaction.
 - more reliable, less race conditions
- More testing of 'run elsewhere' functionality to ensure it is generalized.
- Check Python 3 compatibility

