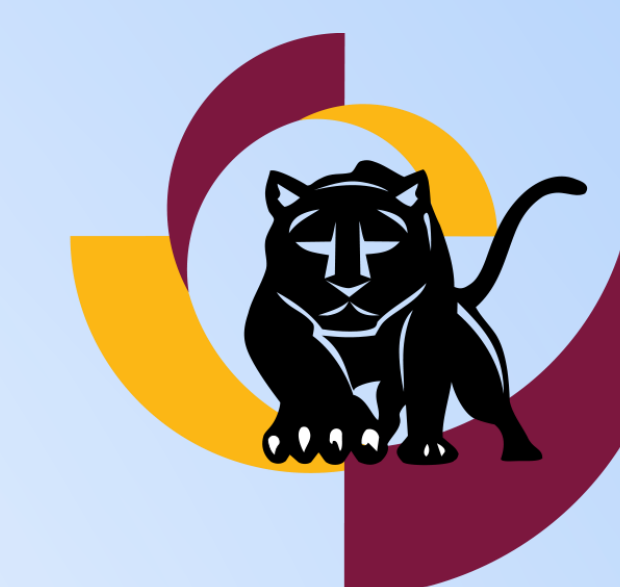




Improving Simulation Data Processing Pipeline using Conduit via HDF5



HARTNELL COLLEGE



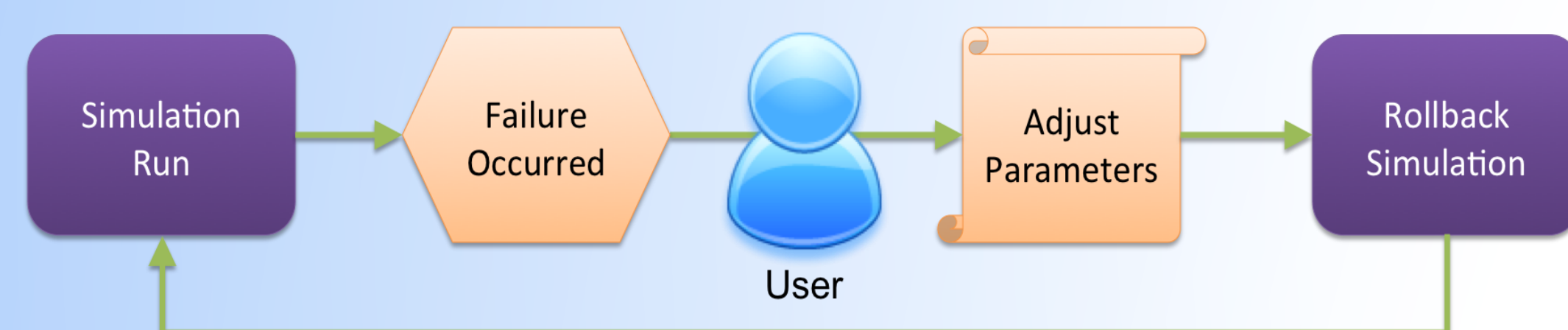
UNIVERSITY OF CALIFORNIA SANTA CRUZ

Gavin Sonne¹, Ming Jiang², Brian Gallagher², Daniel Laney², Cyrus Harrison²
¹Hartnell College, UC Santa Cruz; ²Lawrence Livermore National Laboratory

Background:

Arbitrary Lagrangian-Eulerian (ALE) Codes

- Simulation workflows are complex.
- Involves tuning which requires significant user time and effort.
- Process that relies heavily on the knowledge of the user.



Learning Algorithm-Generated Empirical Relaxer (LAGER)

- Semi-automate this process by using a data analytics approach to reduce reliance on the user.
- Exploit machine learning and develop novel data visualization techniques.
- Develop an *in situ* infrastructure which runs predictive analytics alongside the simulation so that failures can be avoided.

Project Goals:

- Preliminary work to establish data processing pipeline between simulation and analytics.
- Convert simulation data into the Hierarchical Data Format 5 (HDF5).
- Implement HDF5 with NumPy array output.
- Implement Conduit with HDF5 Output.
- Evaluate Conduit's potential as a data exchange method at simulation run time.

HDF5:

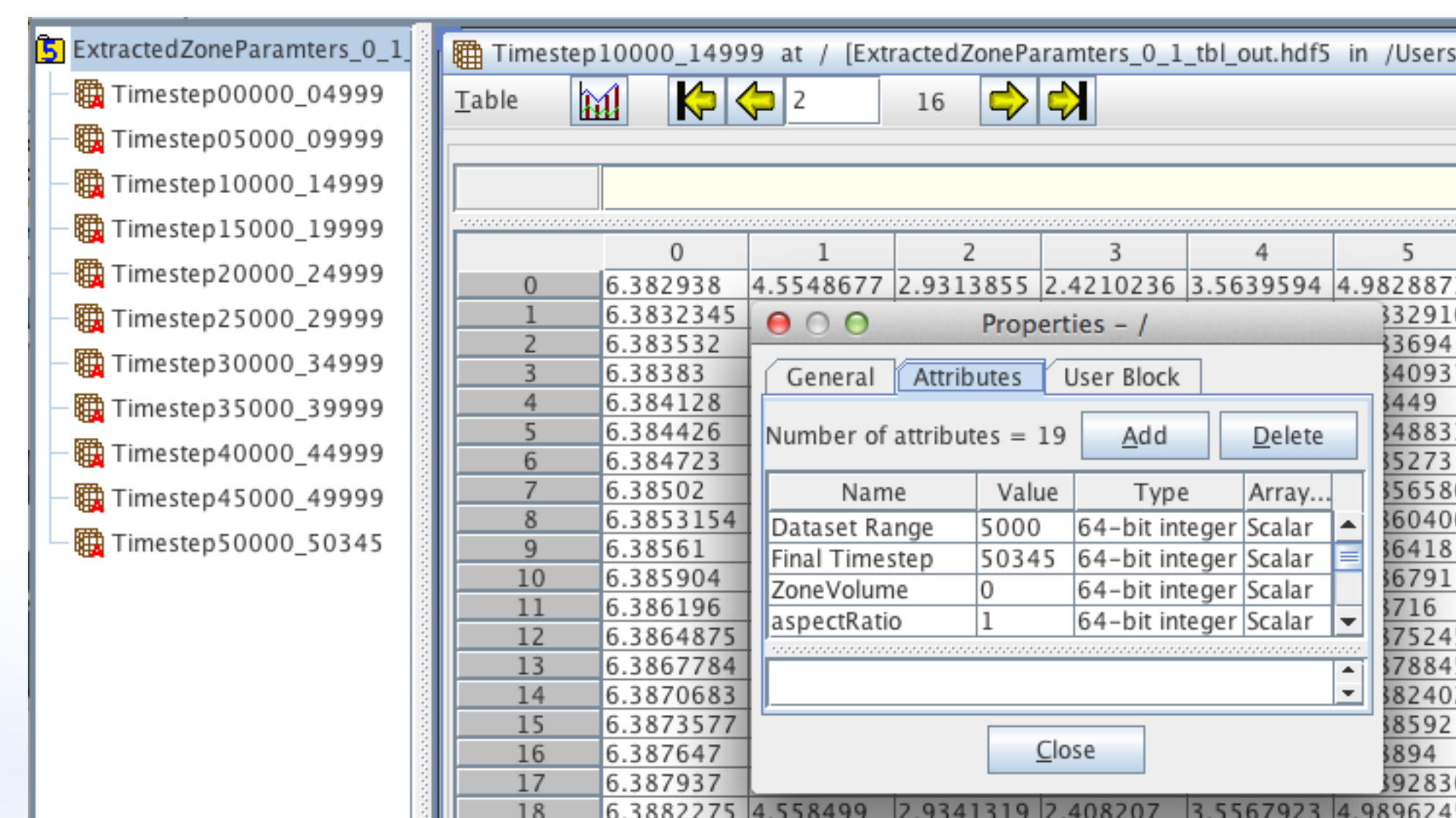
Currently, the data that is extracted from the simulation code is stored in comma separated value (CSV) text files. HDF5 is a file format used to organize and store large amounts of numerical scientific data. A python program was designed that parses the CSV data, and writes it into an HDF5 format using the h5py library. In this implementation, each item of floating point data is reduced from double to single precision, which decreases the output file size by half. Header metadata and information about the HDF5 file are stored as attributes.

```
def write_tbl(outputFile, headerList, dataArray, stepCount, stepList):
    dset = outputFile.create_dataset(("Timestep%05d" % (stepCount-len(stepList)) + ".tbl" % (stepCount - 1)), data=dataArray, dtype='<f4')
```

Function in the parser that writes HDF5 Datasets

```
flax:lager sonne1$ ls -l ExtractedZoneParameters_0_1.tbl_out.hdf5 ExtractedZoneParam
-rw-r--r-- 1 sonne1 55445 1256648544 Jul 28 15:29 ExtractedZoneParameters_0_1.DOUBL
-rw-r--r-- 1 sonne1 55445 628330464 Jul 28 15:13 ExtractedZoneParameters_0_1.tbl_0
```

File size comparison of single vs. double float precision



Generated HDF5 File shown in HDFView, with attributes

Conduit:

- A flexible way to describe complex in-core data.
- A C++ API for accessing that data.
- Designed to be adaptable to other languages, such as Python.
- Designed for in-core data exchange and sharing, thus has the potential to be very useful in the future as a component of the data processing pipeline between simulations and data analytics.

Discussion:

It has been demonstrated that HDF5 is a very useful format for storing very large amounts of data produced by simulation codes. More interestingly, it is very useful to have data in this format when a user needs to be able to extract specific sections of the data in order to work with it. Additionally, once the data is in this format, accessing it can be done much more efficiently than doing so when the data is stored as a CSV.

Continuing Work:

In the time remaining, the plan is to take an existing output from Conduit used with the Livermore Unstructured Lagrangian Explicit Shock Hydrodynamics mini-app (LULESH) and output the data in HDF5 instead of using the SILO library.

Both this research group and the developers of Conduit are interested to know whether Conduit data can be directly and easily exported into HDF5 and back again with little to no additional manipulation of the data. This avenue will be explored very soon.

Another future goal is to insert Conduit into the simulation codes for our data analytics. Another is to output directly to HDF5 and skip the costly and inefficient step of writing the data out to cumbersome CSV files.

Acknowledgements:

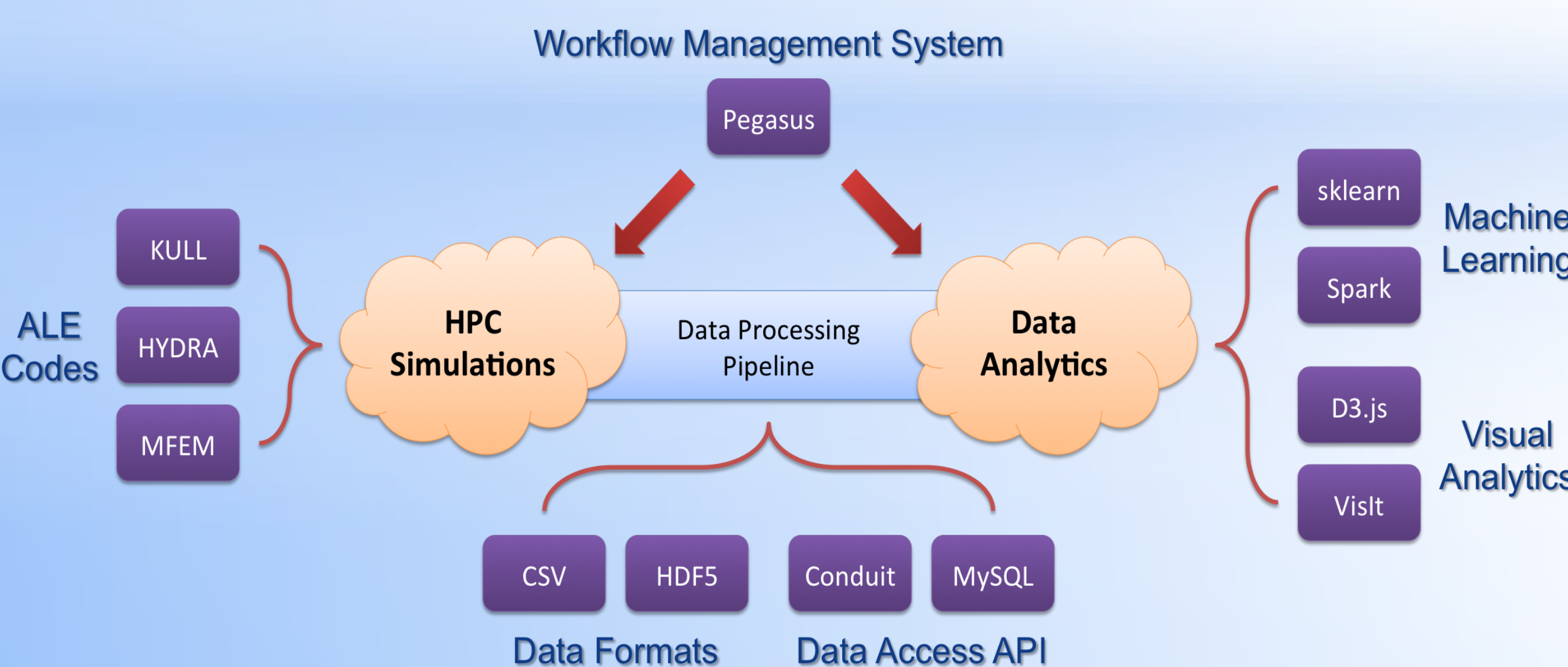
Many thanks to Joshua Kallman, Jay Salmonson, Joe Koning, Luc Peterson, Judy Thomas, and Amy Huang. Most of all, thank you to LLNL for providing the opportunity to do this research.

```
Converting CSV file ExtractedZoneParameters_0_1.txt
Writing steps 0-2999 to file.
Writing steps 3000-3999 to file.
Writing steps 4000-4999 to file.
Writing steps 5000-5999 to file.
Writing steps 6000-6999 to file.
Writing steps 7000-7999 to file.
Writing steps 8000-8999 to file.
Writing steps 9000-9999 to file.
Writing steps 10000-10999 to file.
Writing steps 11000-11999 to file.
Writing steps 12000-12999 to file.
Writing steps 13000-13999 to file.
Writing steps 14000-14999 to file.
Writing steps 15000-15999 to file.
Writing steps 16000-16999 to file.
Writing steps 17000-17999 to file.
Writing steps 18000-18999 to file.
Writing steps 19000-19999 to file.
Writing steps 20000-20999 to file.
Writing steps 21000-21999 to file.
Writing steps 22000-22999 to file.
Writing steps 23000-23999 to file.
Writing steps 24000-24999 to file.
Writing steps 25000-25999 to file.
Writing steps 26000-26999 to file.
Writing steps 27000-27999 to file.
Writing steps 28000-28999 to file.
Writing steps 29000-29999 to file.
Total time: 113
flax:lager sonne1$
```

Parser takes 113 seconds to convert a 2.6 GB CSV File

```
flax:lager sonne1$ python test_read_0_1.txt
1.txt -t 58345
[ [ 1.55186898e-03 3.82403946e+01 ...
  3.58725369e-01 1.33439549e-03 ...
  2.0837762e+01 ...
  1.41145588e-03 1.57215710e+01 ...
  3.45559150e-01 ...
  9.72413830e-03 1.78375447e+00 ...
  4.71846968e-01 1.18654179e-02 ...
  1.53458953e+00 7.16429855e-01 ...
  1.28019294e-02 1.87283787e+00 ...
  5.76188696e-01 ...
  195, 16)
Total time: 33
```

Returning array from final time step:
 - From CSV takes 33 seconds
 - From HDF5 takes 0.1 seconds



High Level Overview of the Workflow Management System