

# REMOTE COLLABORATION VIA DISTRIBUTED HDF FILE ACCESS

Enabling collaborative, small accesses to large remote files in Python

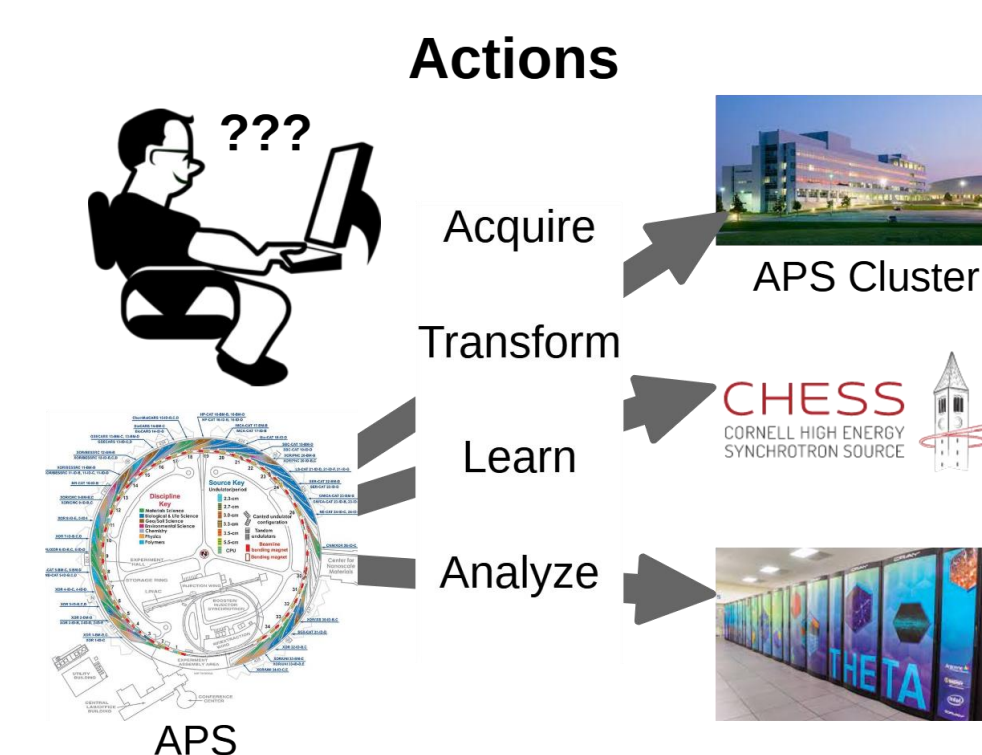
Justin M. Wozniak, Data Science and Learning Division, Argonne National Laboratory  
Raymond Osborn, Materials Science Division, Argonne National Laboratory

## ABSTRACT

Cross-institutional data sharing is still a challenging problem for the large datasets collected at the Advanced Photon Source (APS). Sector 6 at the APS routinely collects single-crystal x-ray diffraction data at a rate of several terabytes per day, which is streamed for automated data reduction in local file stores. Such large data volumes make it challenging to collaborate on data analysis with remote collaborators, without the inefficiencies of transferring full datasets.

In our approach, users are presented with a unified view accessible to all collaborators, which is populated on-demand with experimental data and/or theoretical simulations from multiple sources, and synchronized with distributed remote locations (cloud, exascale computing facilities) as needed. From the user perspective, this will have three advantages: (1) accelerated analysis and learning pipelines due to reduced data transfer overheads; (2) improved integration of data analysis and advanced modeling, and (3) increased productivity due to less management overhead and development overhead.

## MOTIVATION



Multiple stages in data analysis workflow distributed across multiple computing sites, each with differing processing and I/O capabilities. Typically, whole data sets are copied in bulk, incurring complexity and inefficiencies.

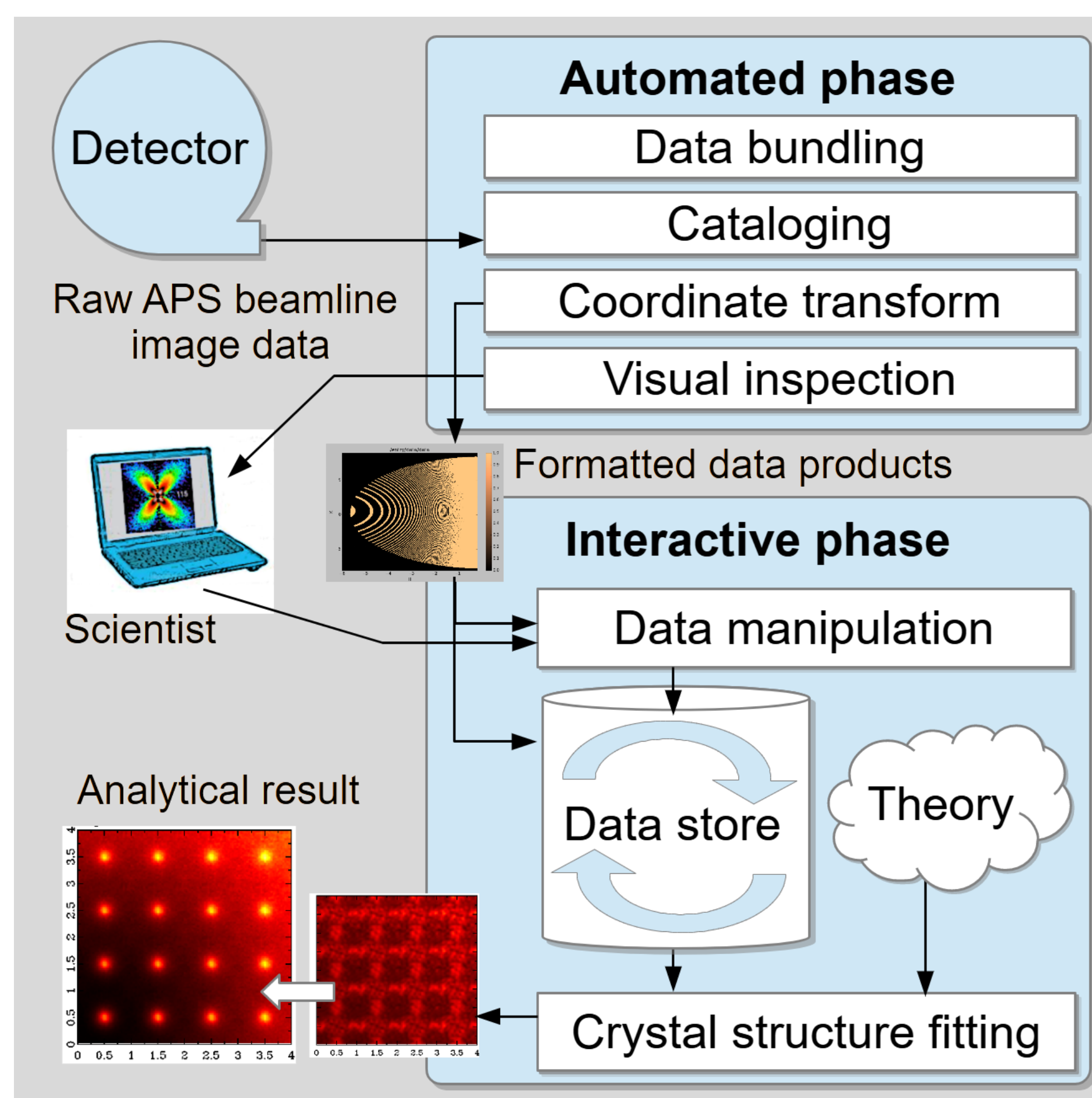
## METHODS

Client	Actions	Container
Detectors	Upload	User-Added Tools
Filters	Query	Reusable Data Tools
Replay	Slice	Remote Object Service
Assimilate	Analyze	
Learning	Learn	

Multiple stages in data analysis workflow supported by remote data tools, including serving Pythonic objects for learning from small slices of large data.

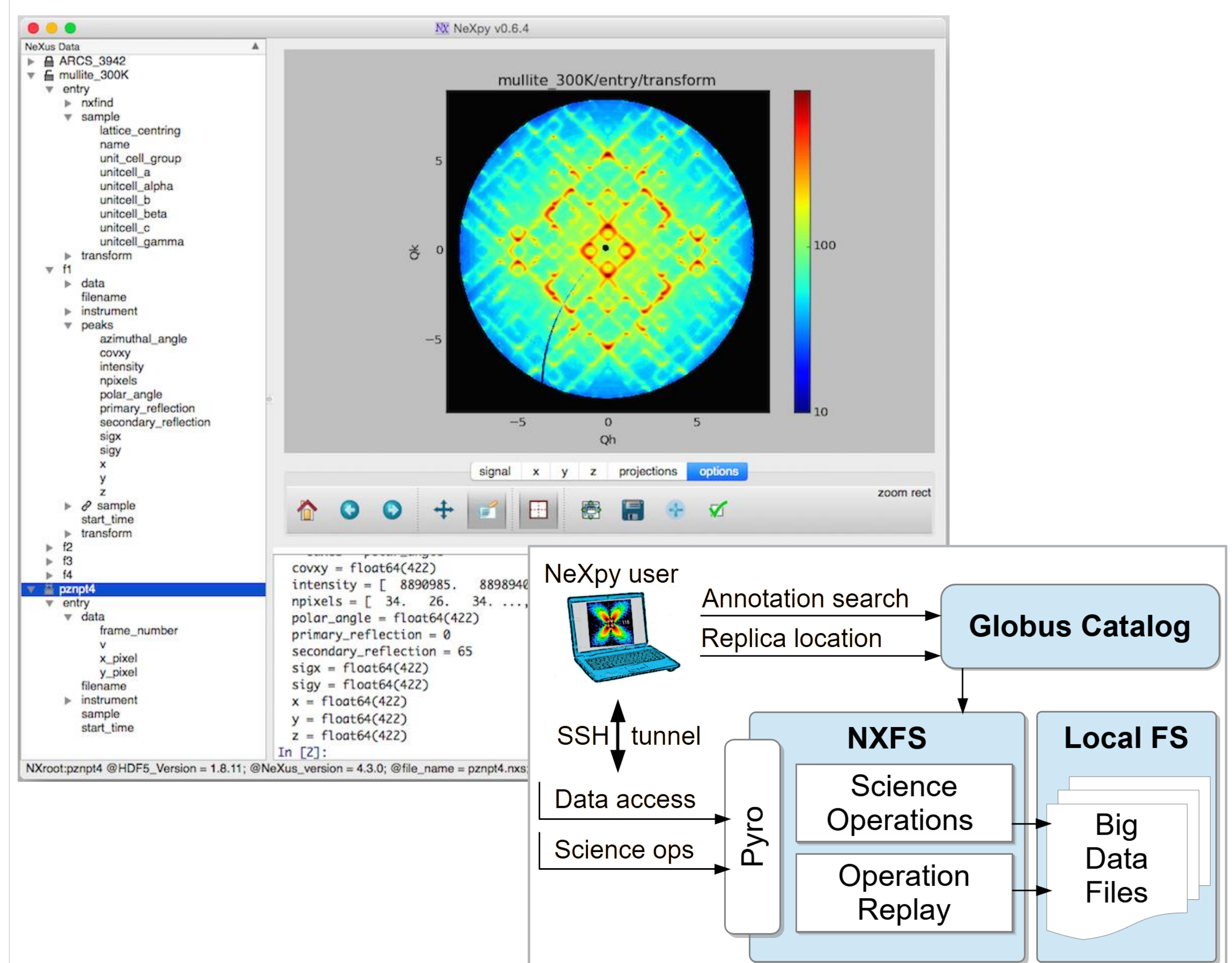
## SMALL ACCESSES TO LARGE DATA

- There is a need to be able to collaborate with remote users performing visualization, lightweight analysis, and small data modification on these datasets, tasks for which the full dataset is not needed. Thus, we desire to create a system in which it is possible to enable remote data slicing and selection.
- Want to be able to perform slicing, visualization, analysis, and learning on remote data.



## USING THE HDF OBJECT INTERFACE

- The HDF interface allows for the possibility of lightweight data access, at least at the interface level. We are prototyping this approach using the a Pyro wrapper around HDF, a distributed service that can be hosted on an institutional cluster or commercial cloud and accessed by a remote client. This proposed solution integrates with existing analysis routines and visualization tools such as NeXpy.



## CONCLUSIONS

- Our approach is to develop a toolkit of learning and analysis-ready Python libraries that can be integrated into workflows by users from a broad range of disciplines.
- Enhanced data analysis and learning will be enabled through extensions to a remote object toolkit which provides a familiar numerical interface, and optimizable aggregated data pipelines for learning frameworks.
- We want to enable the orchestration of remote data sources, selecting from available datasets, providing high level indexing and query support to extract, group and present the data from multiple data services.

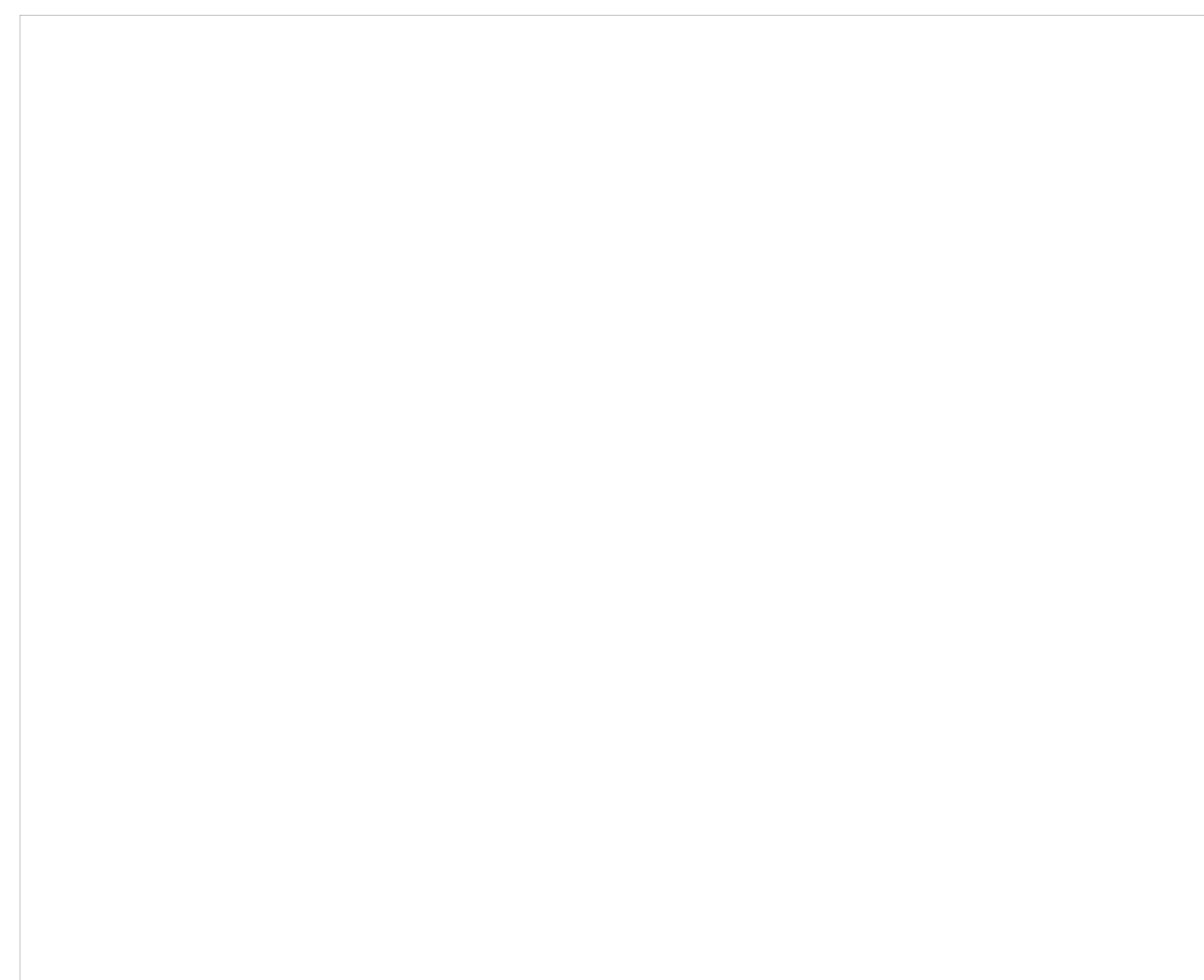
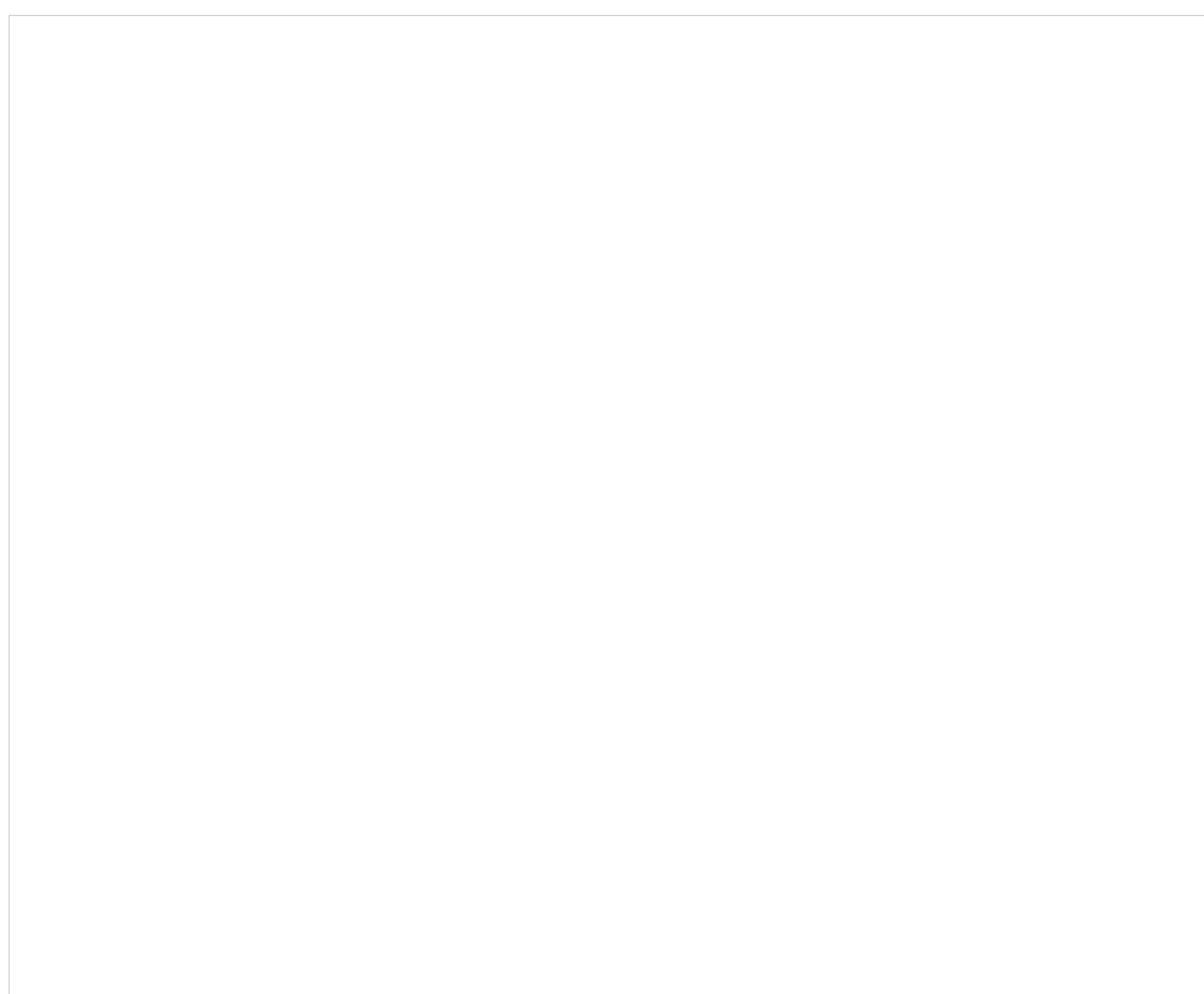
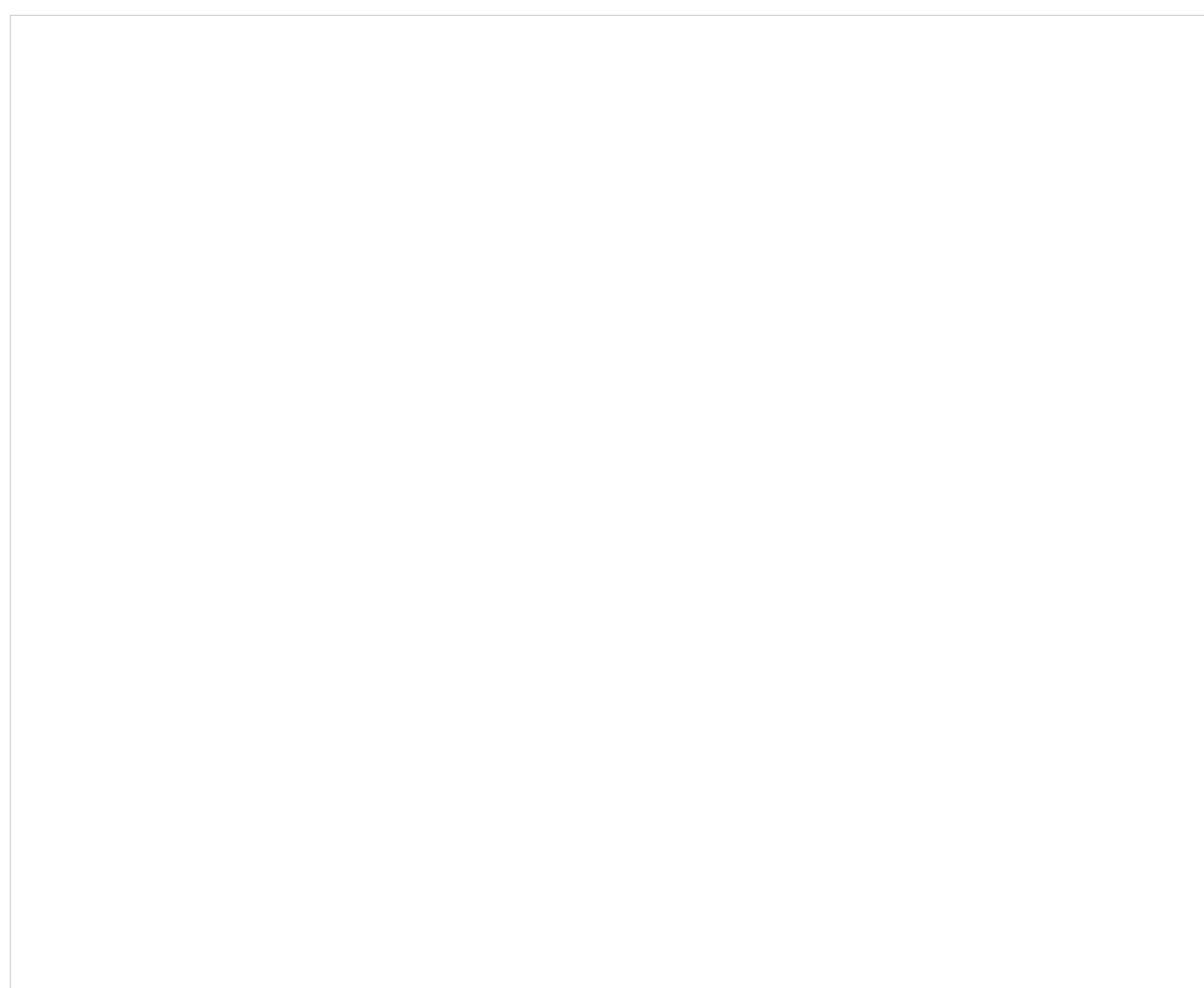
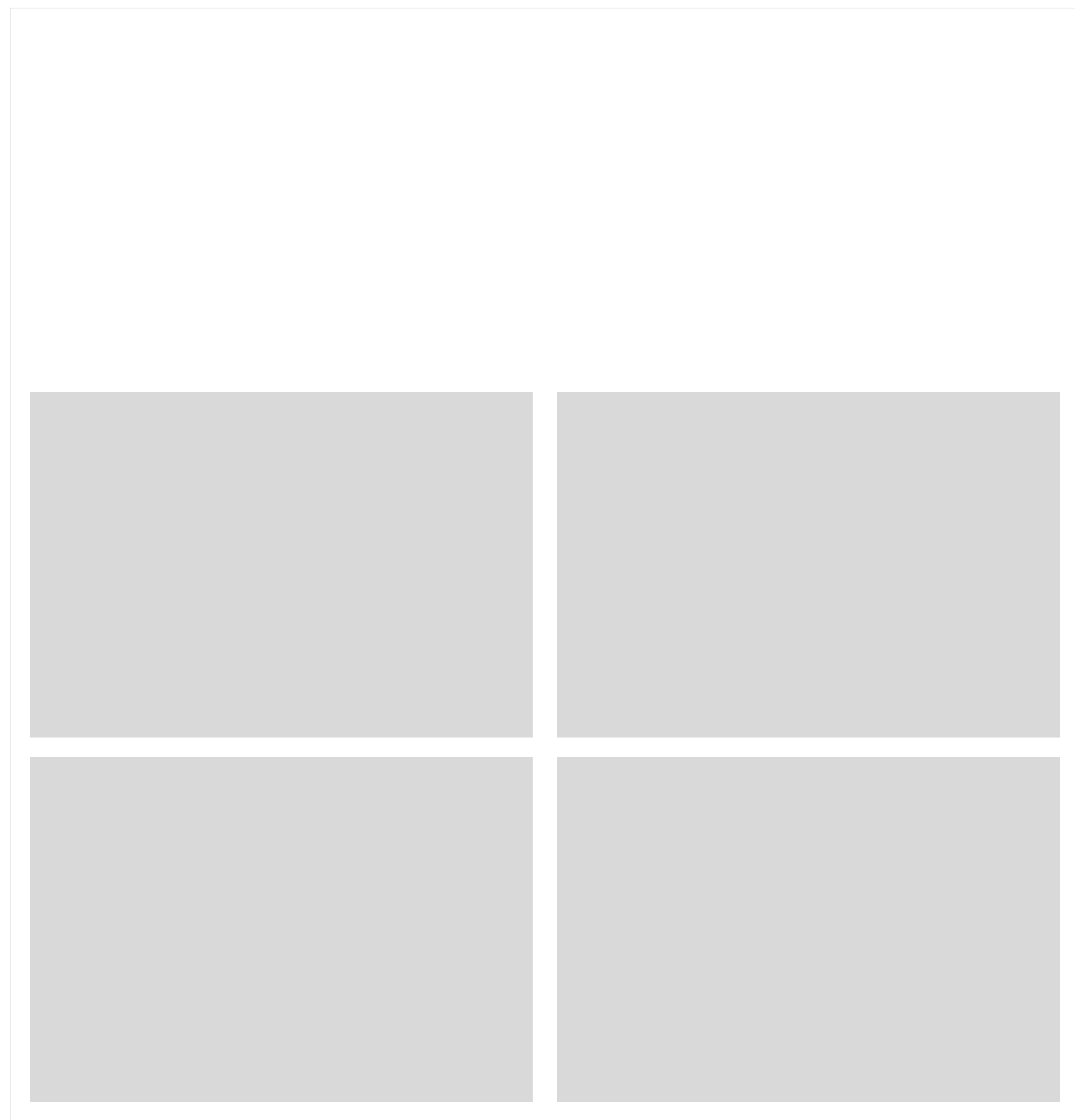
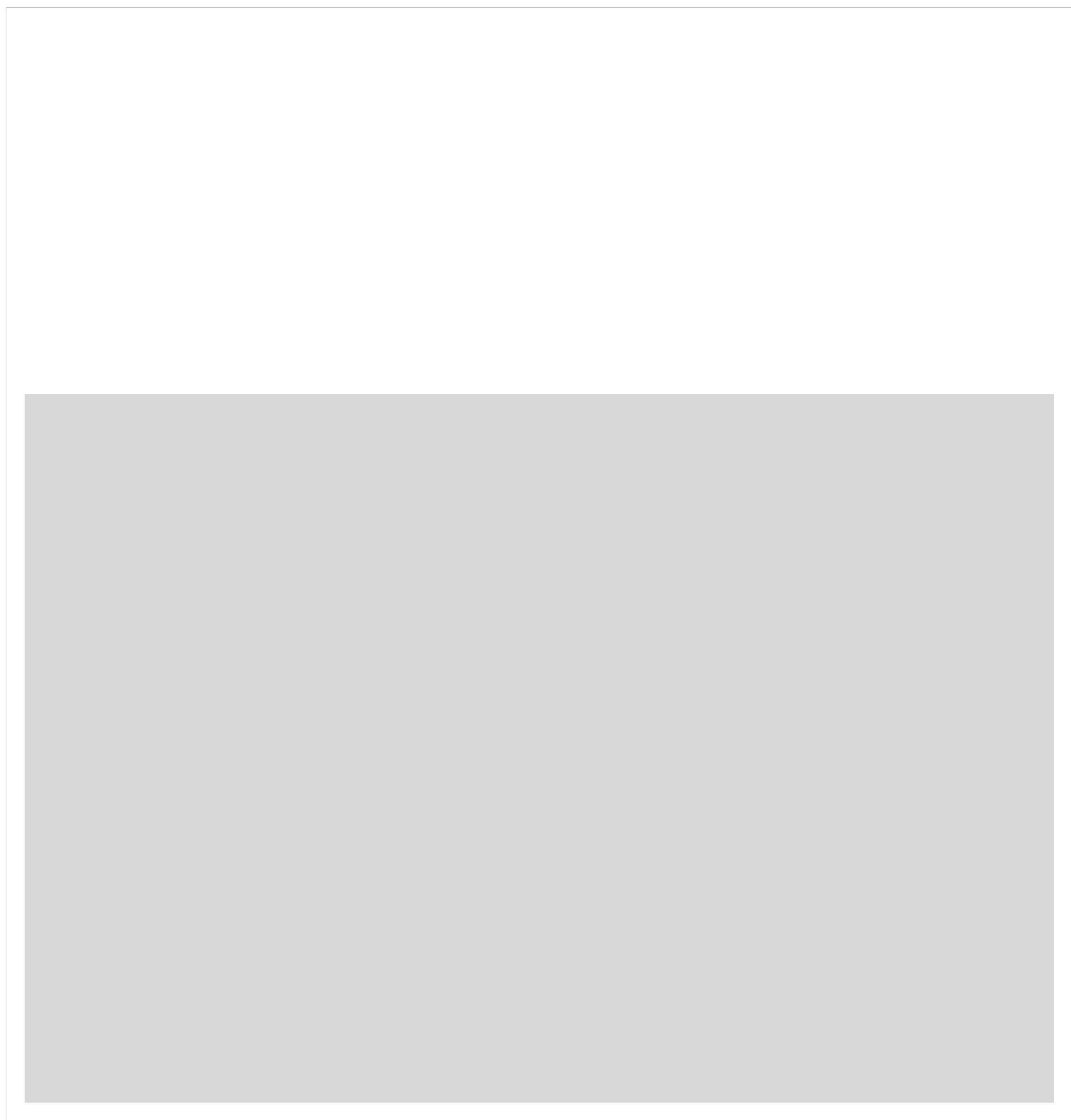
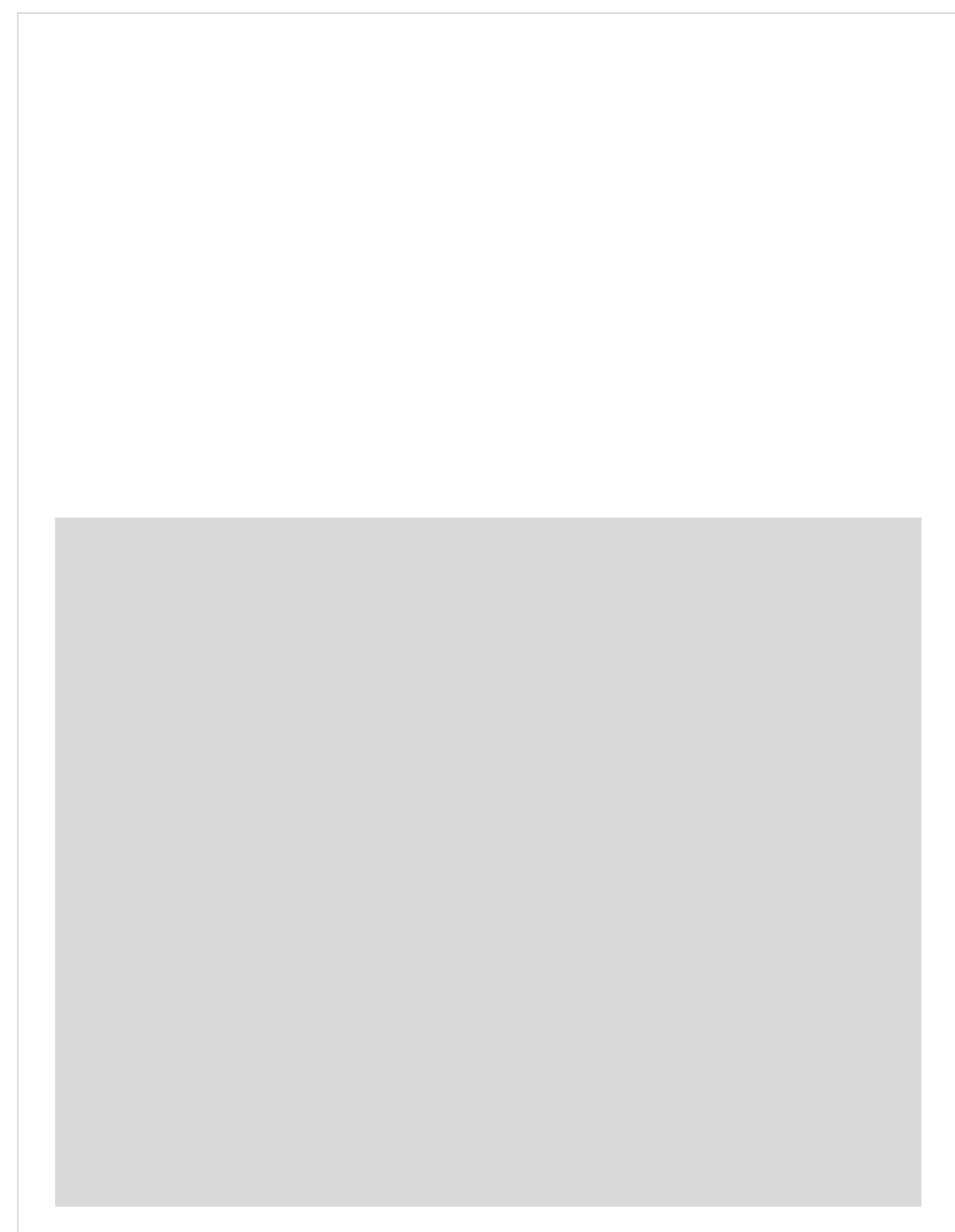
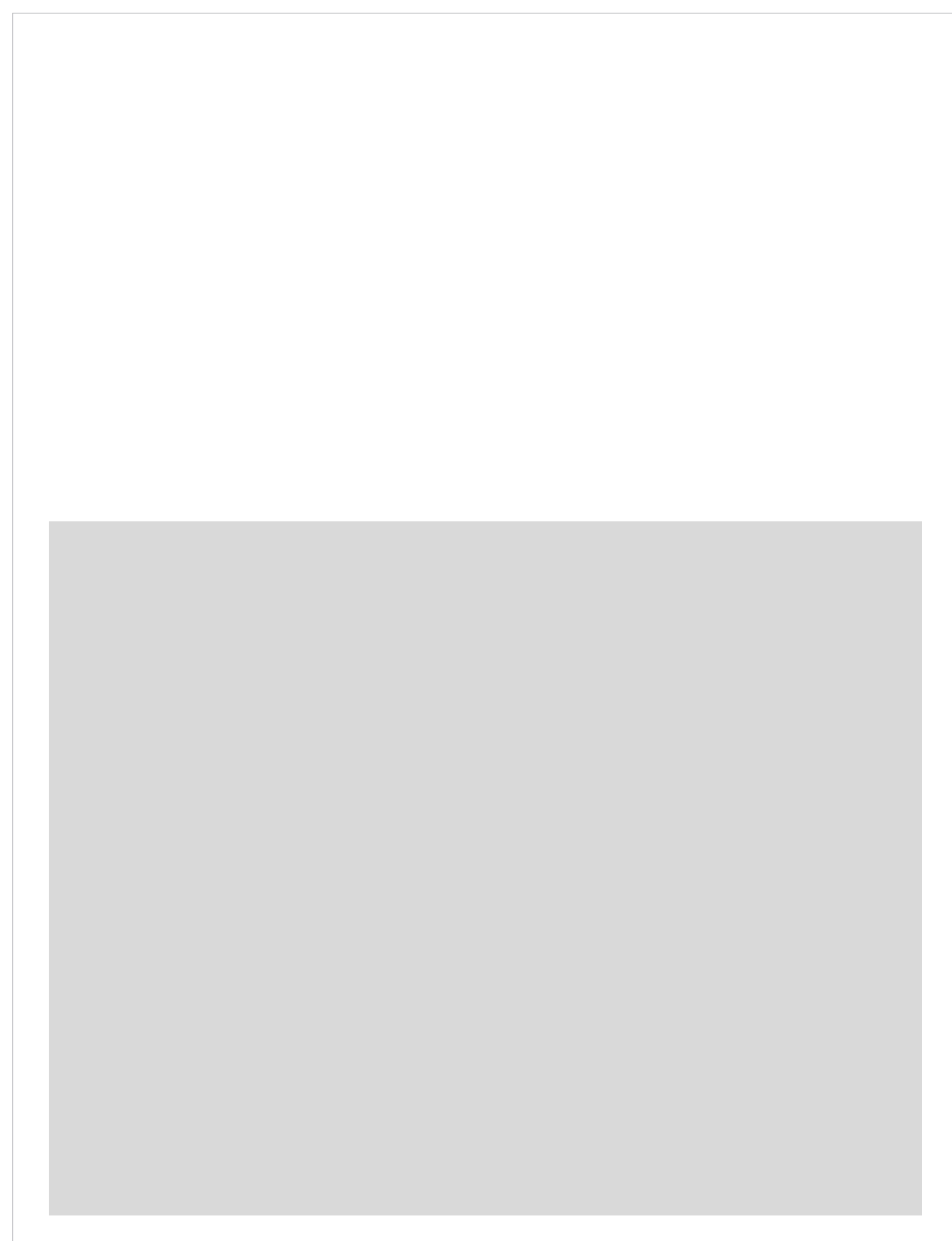
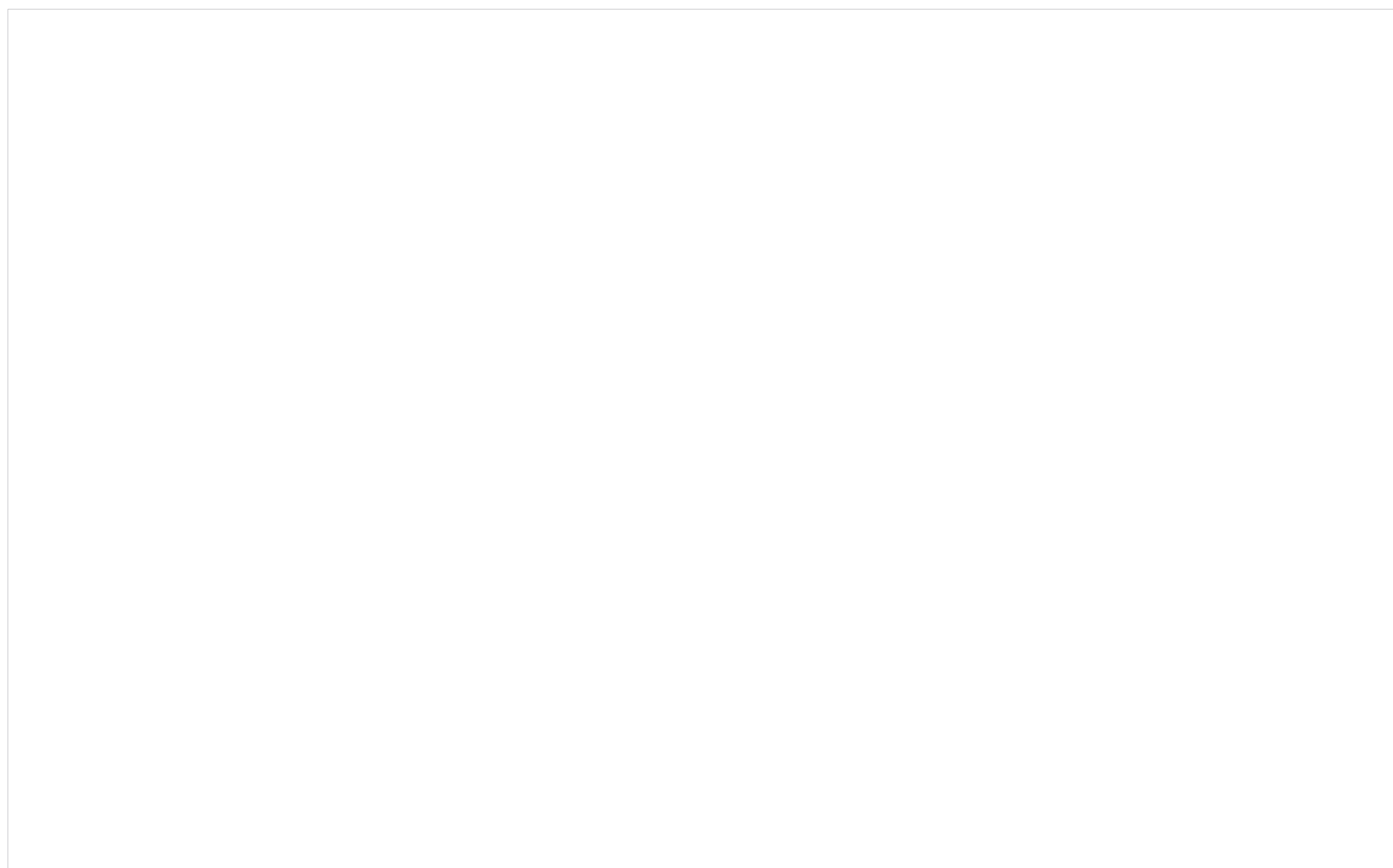
## SOFTWARE CAPABILITIES

- We have made widespread use of HDF5 external links in the prototype as a way of separating the raw data, that should be immutable and protected from corruption, and the higher level metadata both from the experiment and added through the pipeline.
- This prototype has been successfully demonstrated at both the APS and CHES, and will soon be tested for joint analyses of diffuse scattering at APS and the Spallation Neutron Source. This project should be of use in encouraging interfacility cooperation through the standardization of remote data access procedures, and, potentially, in the design of remote data facilities.

## REFERENCES

- Big data remote access interfaces for light source science.** Justin M. Wozniak, Kyle Chard, Ben Blaiszik, Ray Osborn, Michael Wilde, and Ian Foster. Proc. Big Data Computing 2015.
- Data object distribution for experimental science pipelines.** Justin M. Wozniak, Ray Osborn, and Jacob Ruff. ASCR Workshop on the Management and Storage of Scientific Data 2022.





# 30" X 40" HOW CAN WE USE ENERGY FROM THE SUN TO MAKE FUELS?

## Nature driven photochemical approaches to solar fuels

Sarah Soltau, Sunshine Silver, and Lisa Utschig, Chemical Sciences and Engineering Division, Sarah Soltau, Sunshine Silver, and Lisa Utschig, Chemical Sciences and Engineering Division, Sarah Soltau, Sunshine Silver, and Lisa Utschig, Chemical Sciences and Engineering Division

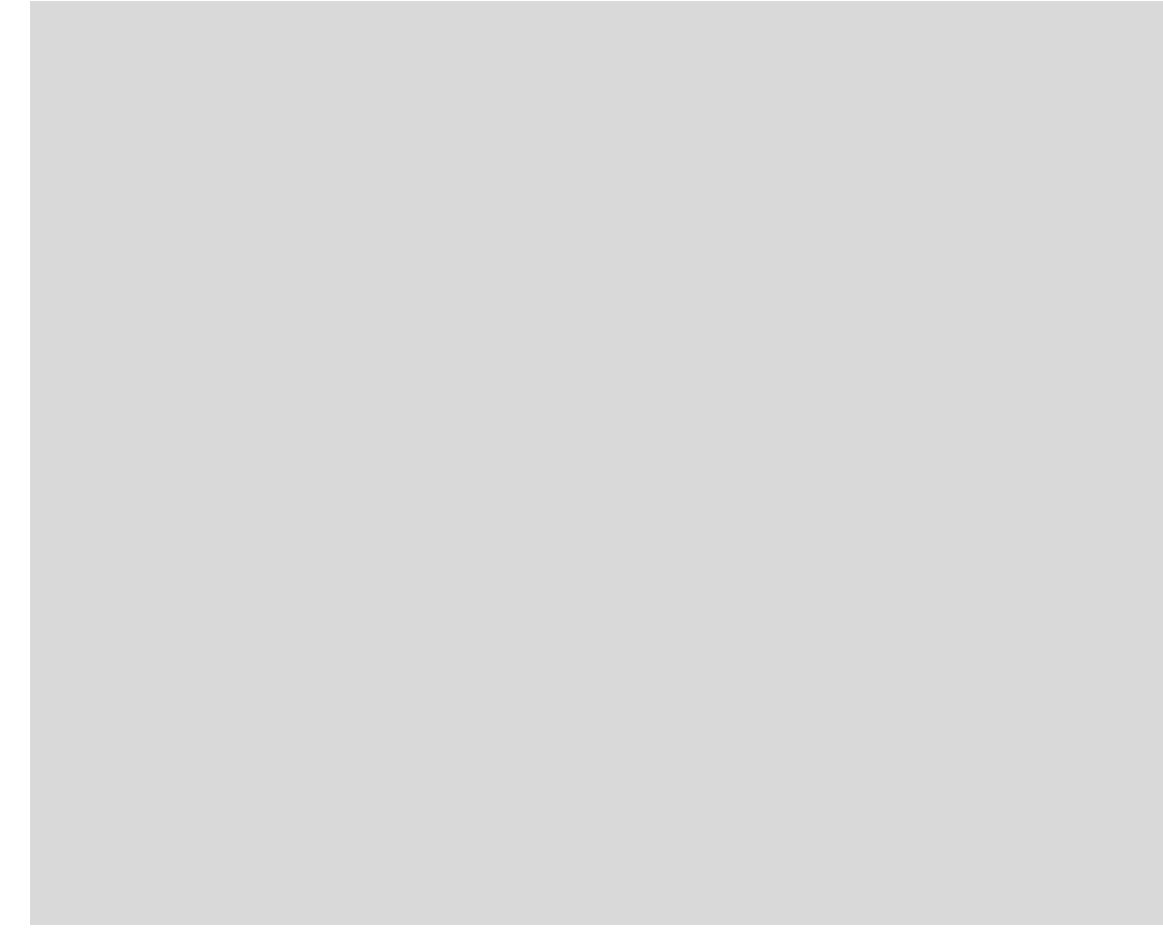
### ABSTRACT

Et volorepudis et lis essit am eum et quas cone sinciet laboresed maiost et modisit voleni cone volor ratemporum voluptassi dolupta dis explistrum doluptaqui sunt quisita comnis que porrupt atione cuptat est, con nis exceat aut debit eossimaio tecercidebis molorporepe volupta tatur, od quam id mollabo.

- Nem sequidendi ate volorum facculpa veni teni rerit andunt vendige ntorest, ipsunt omnimil ium doluptatini vendus consequue mos vendebis.
- Modigna tiurehe nimus, cum aut as rerit ma vollaut remporum volupta de dis ut ipsaepe ditaquam, quis ex endia que con ni doluptia comnist iostione repudam cust ad est hillici isitiun ditatur, quidebis mi, ad maximet des exerum volore sequid maximodis quiatemolum eius, tem
- Fugitatat. Pel inis debis est ape vellere rition niminve rectoribus sunt quatem di reium ratentionet
  - Igendi officit aerferu mquate nat offic to eossiminctis magni omnimod igendem alique modi iliquia pa voloreptas et et mod et, simped eos rem

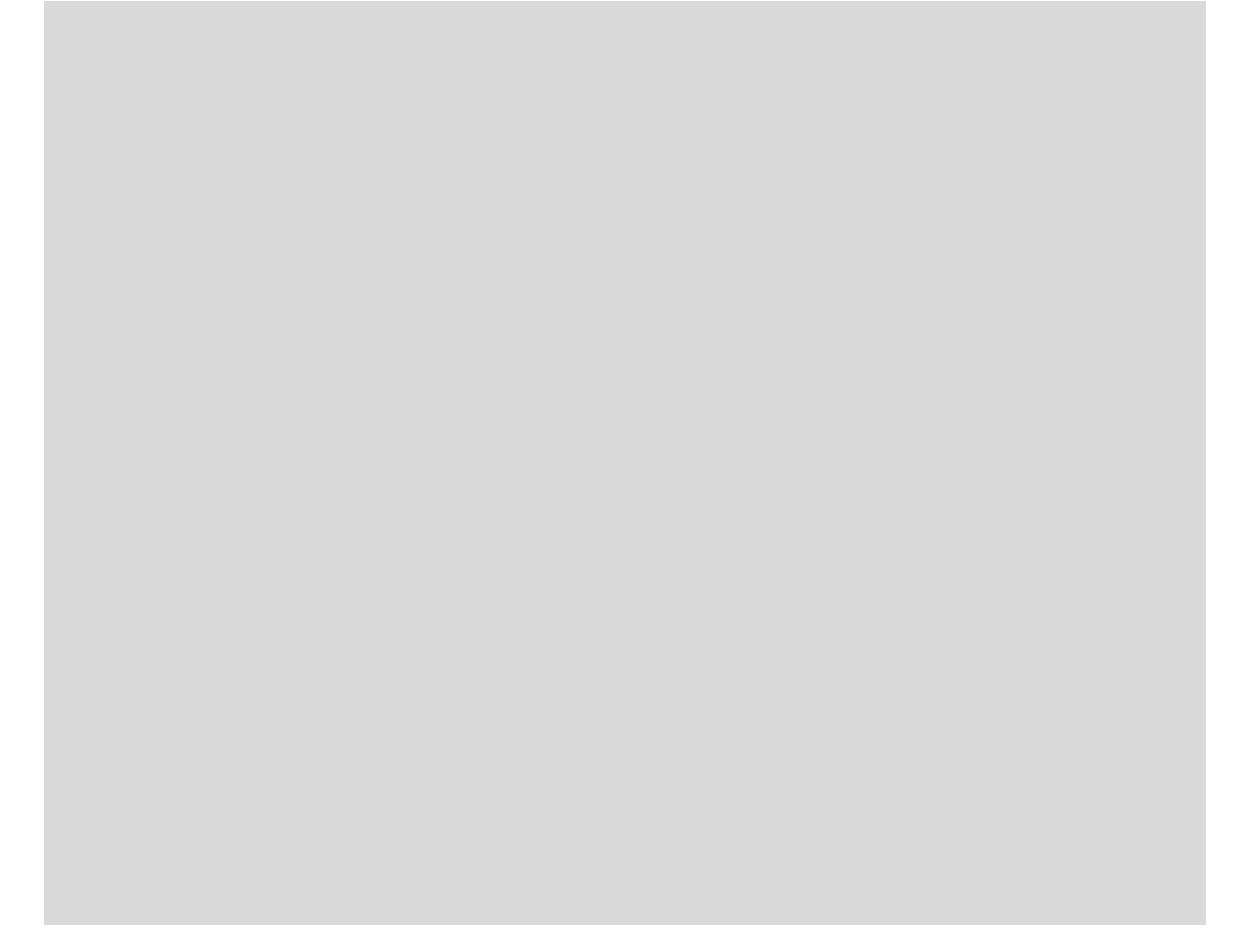
### MOTIVATION

- Et volorepudis et lis essit am eum et quas cone sinciet laboresed maiost et modisit voleni cone volor ratemporum voluptassi dolupta dis



### METHODS

- Et volorepudis et lis essit am eum et quas cone sinciet laboresed maiost et modisit voleni cone volor ratemporum voluptassi dolupta dis



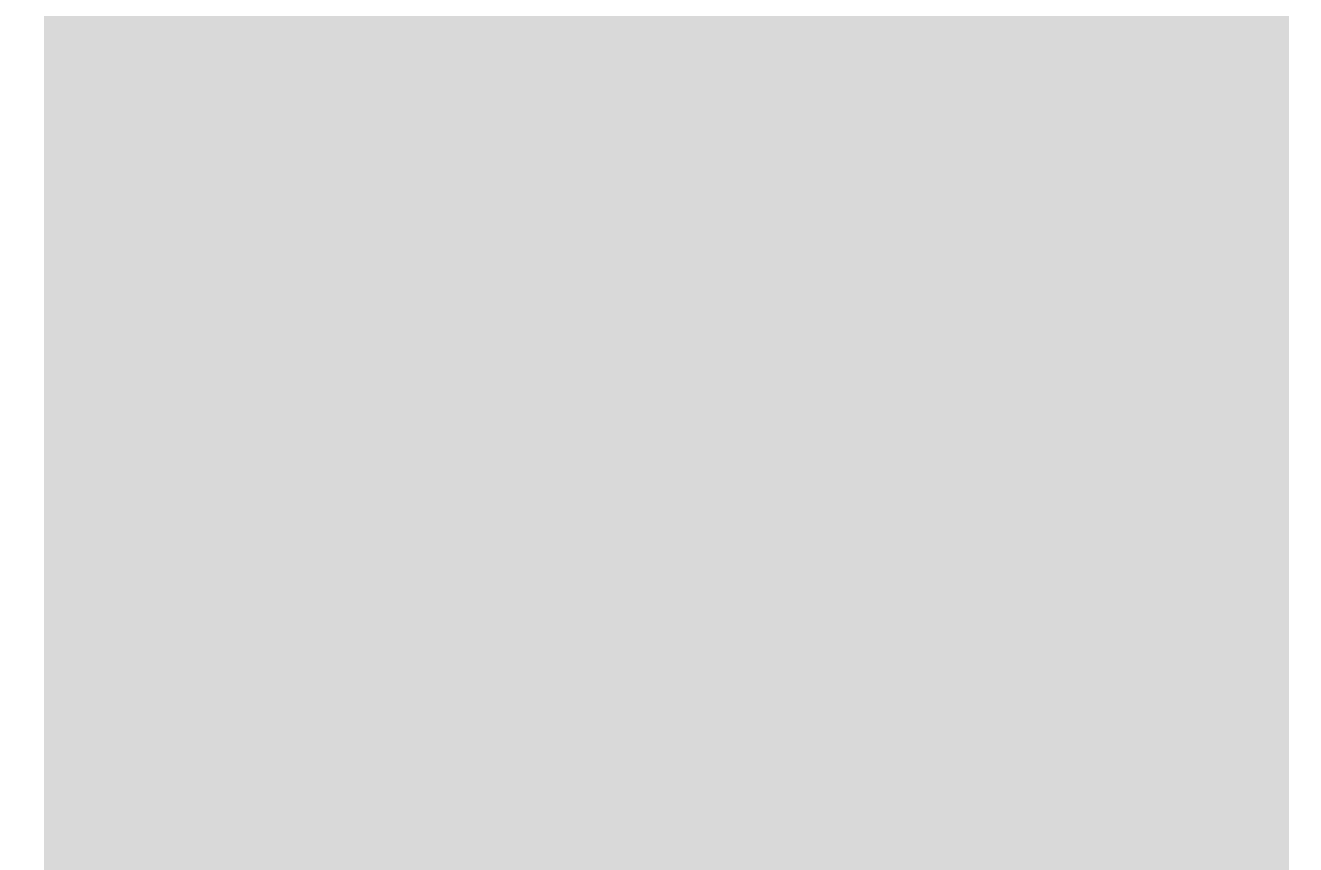
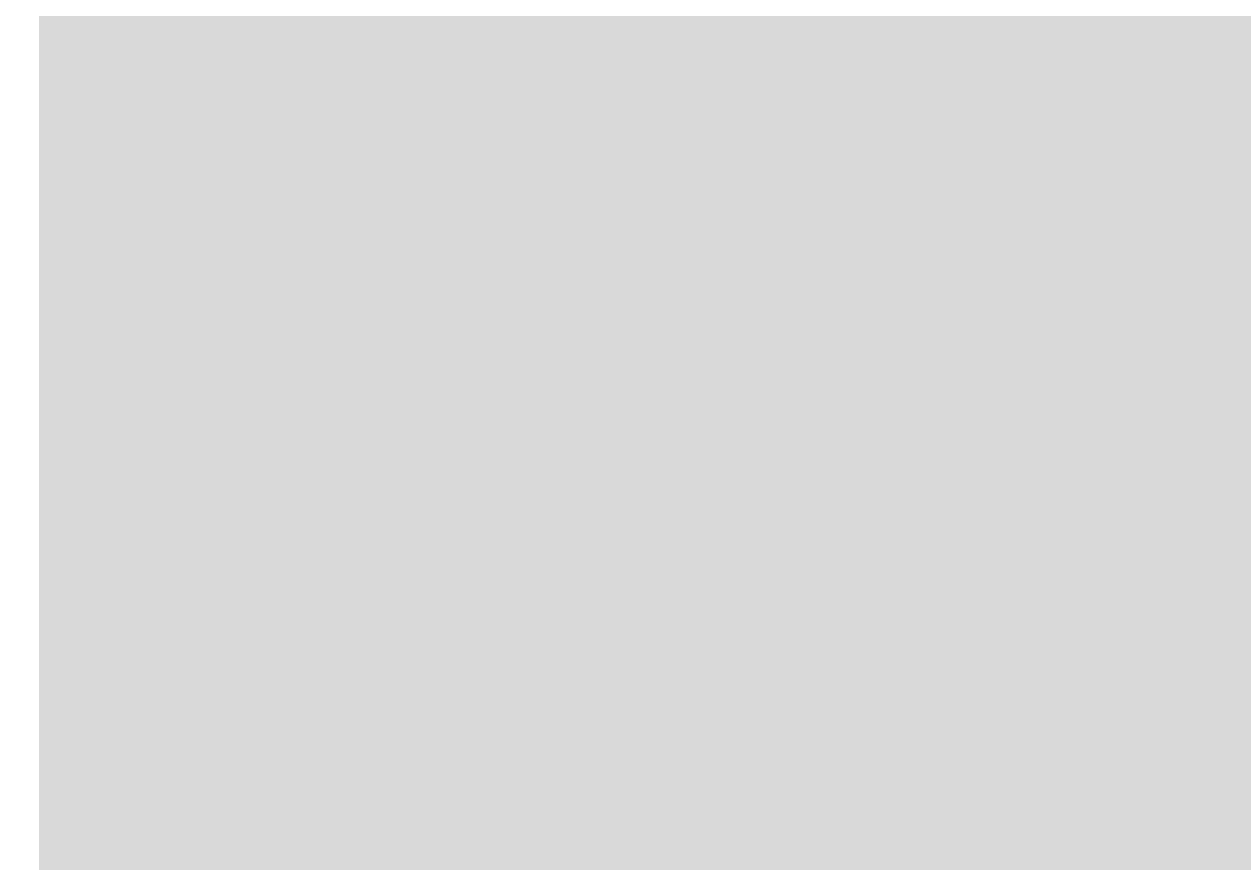
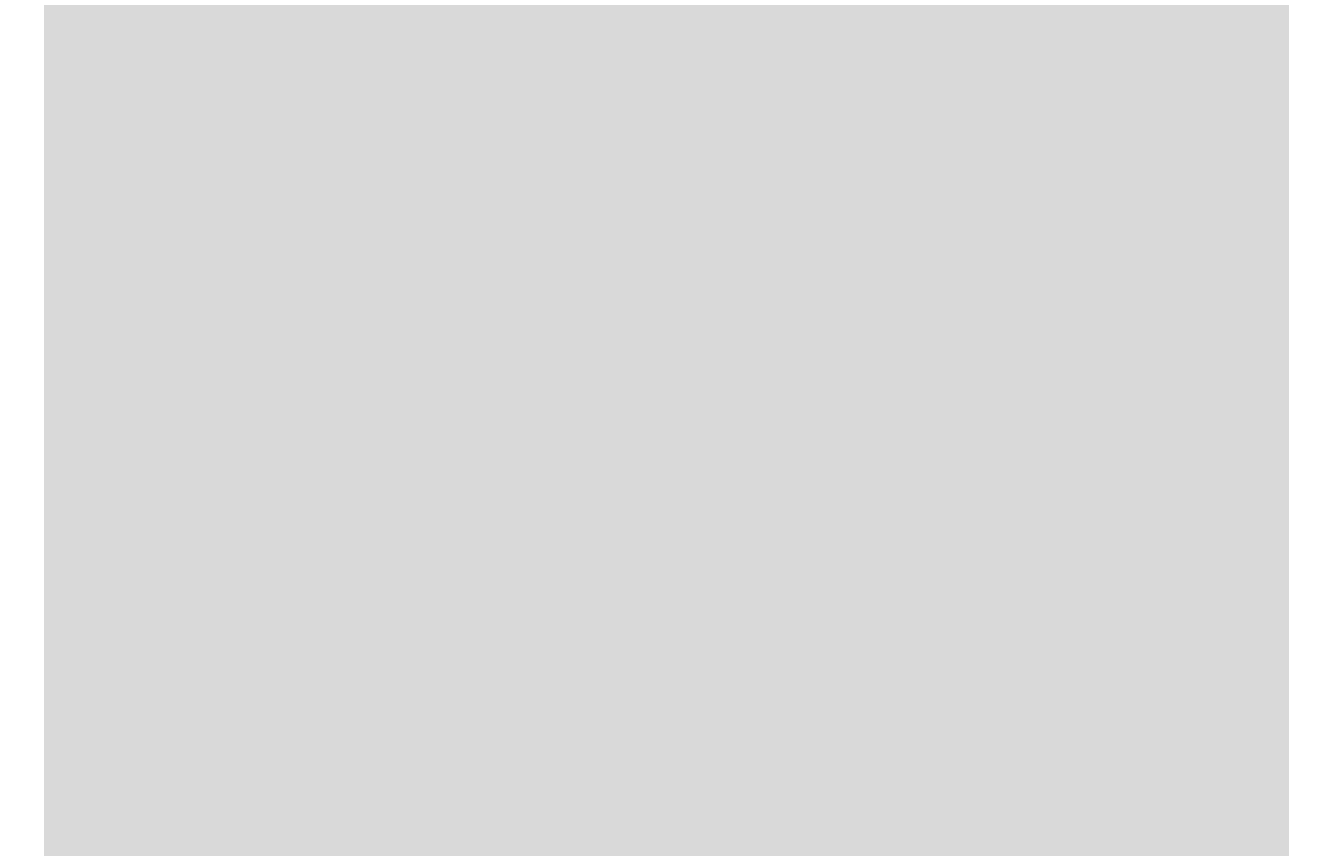
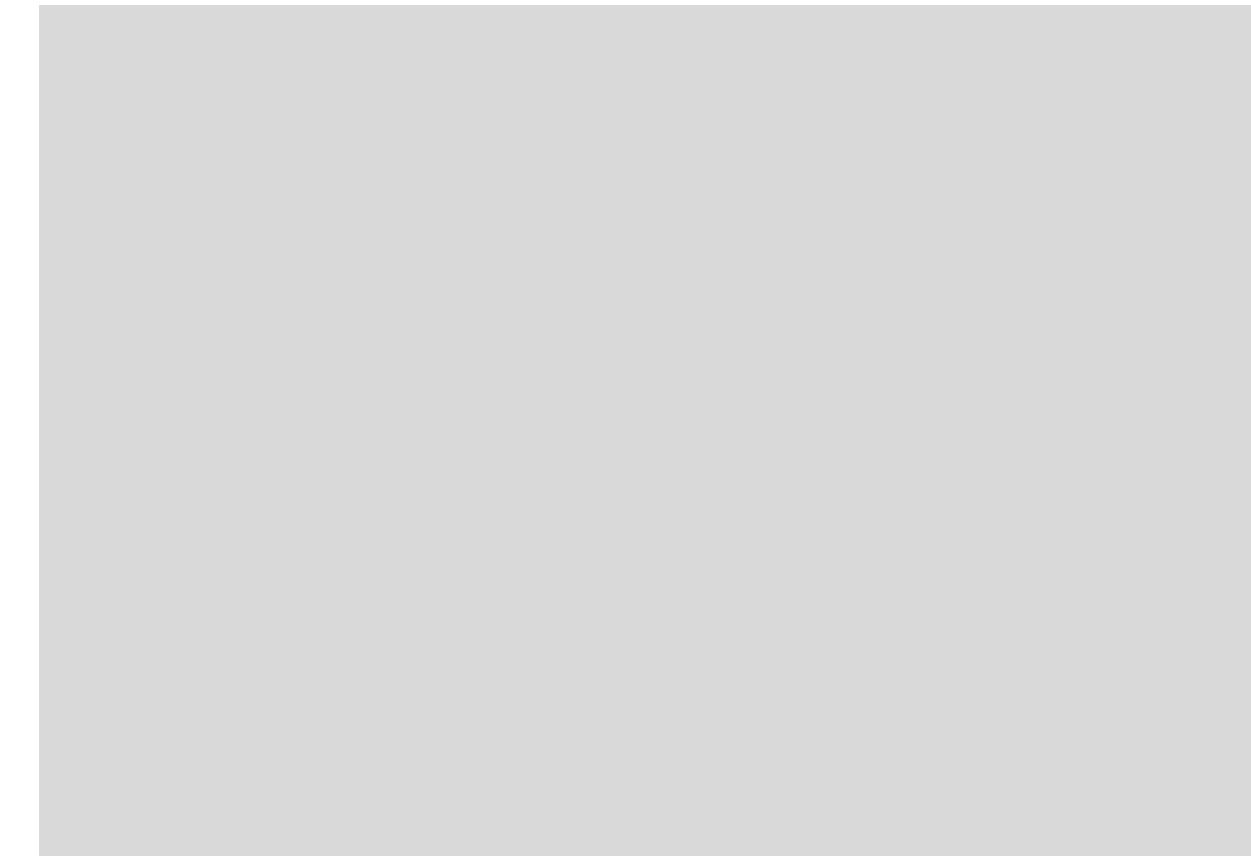
### SECTION 1

- Et volorepudis et lis essit am eum et quas cone sinciet laboresed maiost et modisit voleni cone volor ratemporum voluptassi dolupta dis explistrum doluptaqui sunt quisita comnis que porrupt atione cuptat est, con nis exceat aut debit eossimaio tecercidebis molorporepe volupta tatur, od quam id mollabo. Nem sequidendi ate volorum facculpa veni teni rerit andunt vendige ntorest, ipsunt omnimil ium doluptatini vendus consequue mos vendebis modigna tiurehe nimus, cum aut as rerit ma



### SECTION 2

- Et volorepudis et lis essit am eum et quas cone sinciet laboresed maiost et modisit voleni cone volor ratemporum voluptassi dolupta dis explistrum doluptaqui sunt quisita comnis que porrupt atione cuptat est, con nis exceat aut debit eossimaio tecercidebis molorporepe volupta tatur, od quam id mollabo. Nem sequidendi ate volorum facculpa veni teni rerit andunt vendige ntorest, ipsunt omnimil ium doluptatini vendus consequue mos vendebis modigna tiurehe nimus, cum aut as rerit ma



### CONCLUSIONS

- Et volorepudis et lis essit am eum et quas cone sinciet laboresed maiost et modisit voleni cone volor ratemporum voluptassi dolupta dis explistrum doluptaqui sunt quisita comnis que porrupt atione cuptat est, con nis exceat aut debit.
- Eossimaio tecercidebis molorporepe volupta tatur, od quam id mollabo. Nem sequidendi ate volorum facculpa veni teni rerit andunt vendige ntorest, ipsunt omnimil ium doluptatini vendus consequue mos vendebis modigna tiurehe

### NEXT STEPS

- Et volorepudis et lis essit am eum et quas cone sinciet laboresed maiost et modisit voleni cone volor ratemporum voluptassi dolupta dis explistrum doluptaqui sunt quisita comnis que porrupt atione cuptat est, con nis exceat aut debit.
- Eossimaio tecercidebis molorporepe volupta tatur, od quam id mollabo. Nem sequidendi ate volorum facculpa veni teni rerit andunt vendige ntorest, ipsunt omnimil ium doluptatini vendus consequue mos vendebis modigna tiurehe

### REFERENCES

- Et volorepudis et lis essit am eum et quas one sinciet laboresed maiost et modisit voleni.
- Cone volor ratemporum voluptassi dolupta dis explistrum doluptaqui sunt quisita comnis que porrupt atione cuptat est, con nis exceat aut debit eossimaio tecercidebis.
- Et volorepudis et lis essit am eum et quas one sinciet laboresed maiost et modisit voleni.
- Cone volor ratemporum voluptassi dolupta dis explistrum doluptaqui sunt quisita comnis que porrupt atione
- Cuptat est, con nis exceat aut debit eossimaio tecercidebis.

