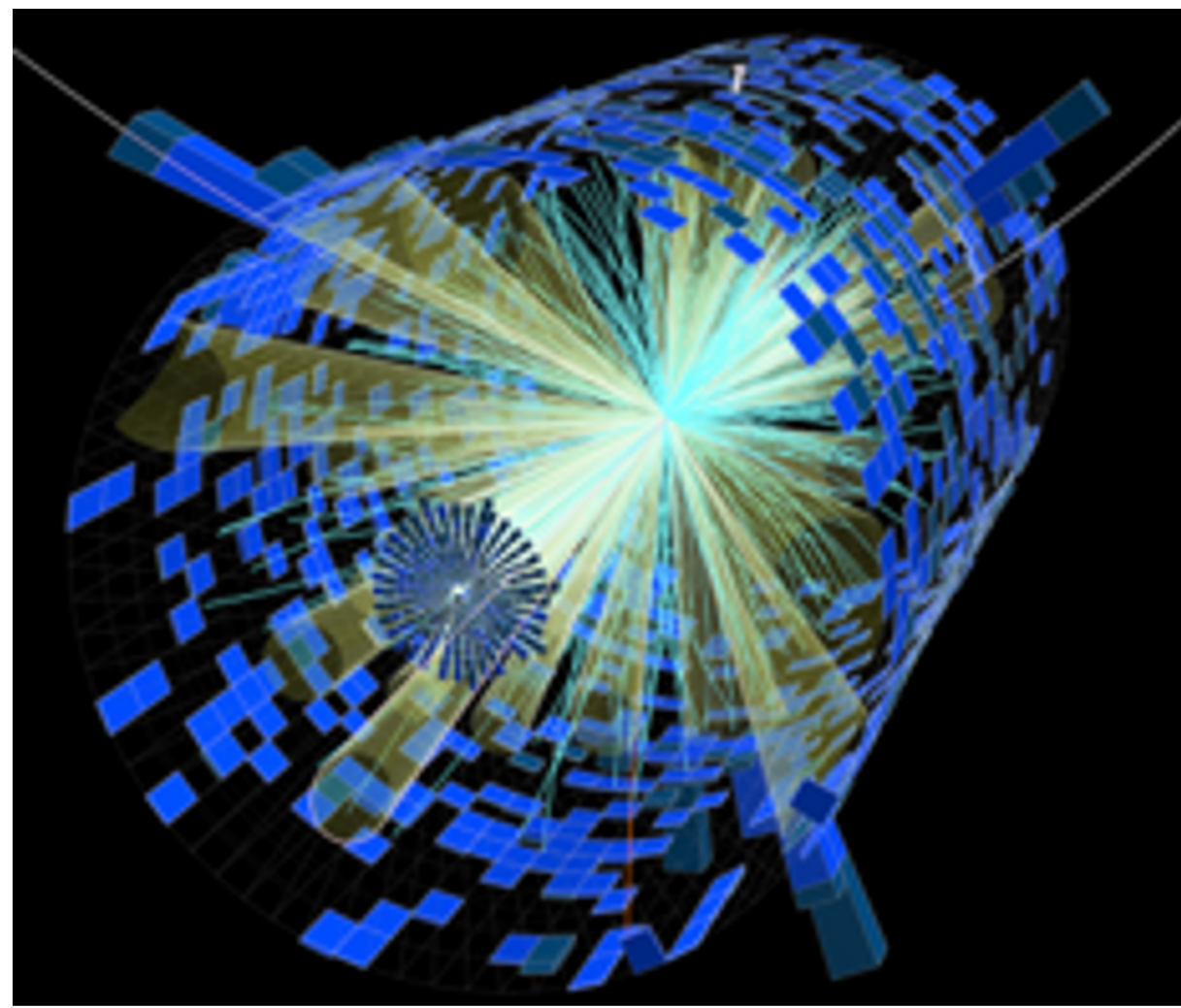


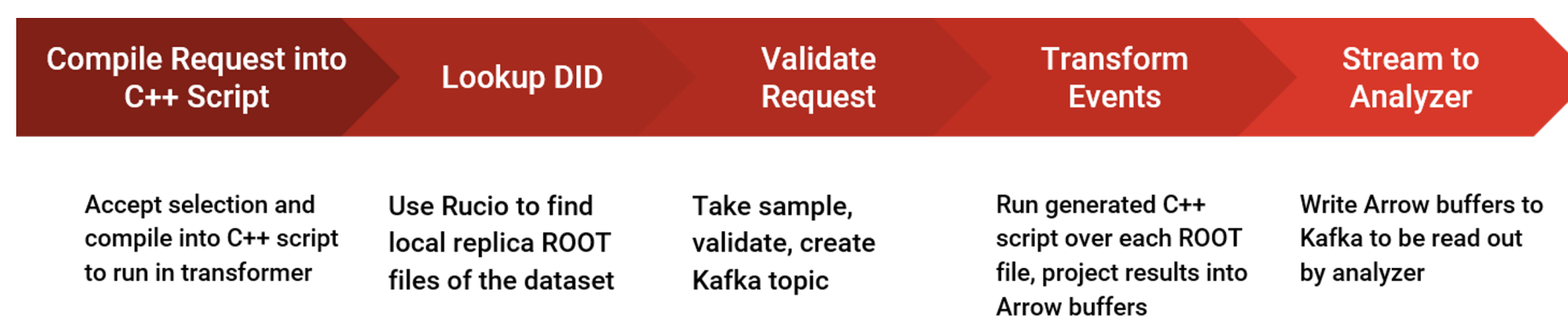
## Common data challenges from ATLAS and CMS

The High Luminosity Large Hadron Collider (HL-LHC) faces enormous computational challenges in the 2020s. The HL-LHC will produce exabytes of data each year, with increasingly complex event structure due to high pileup conditions.

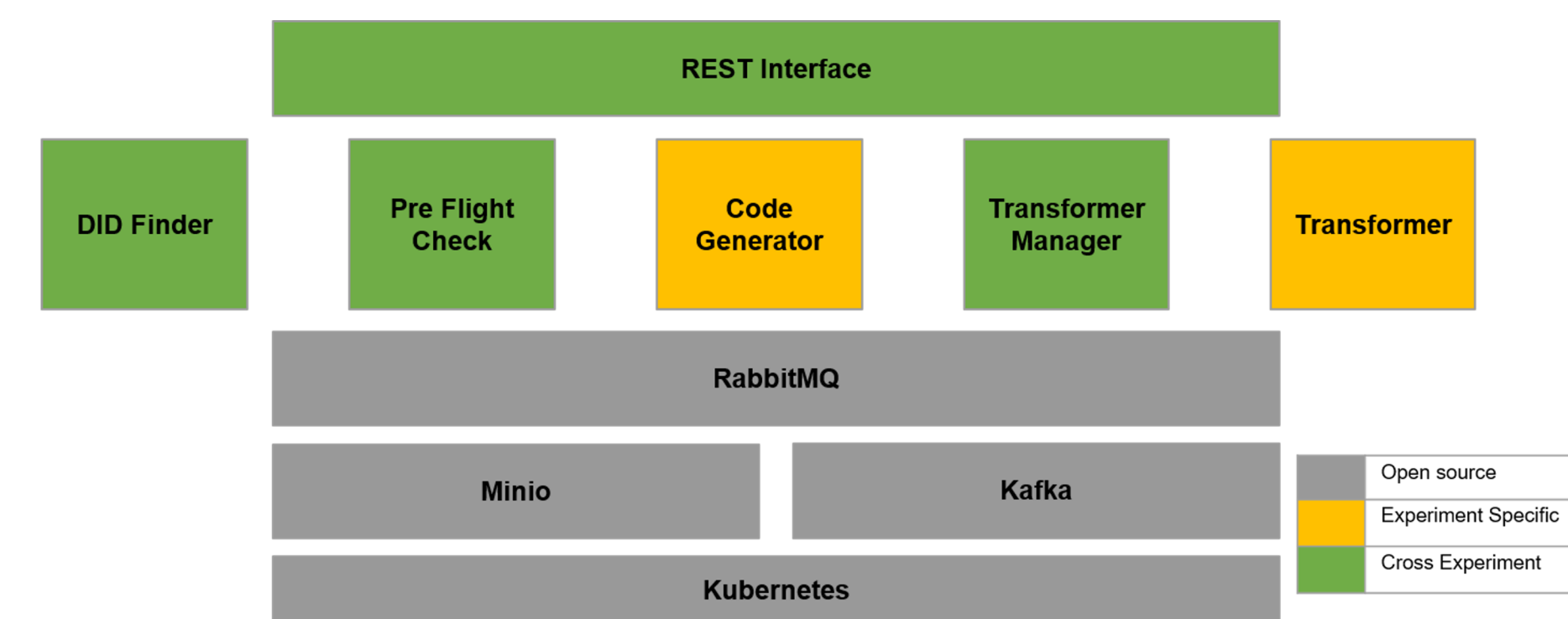


The ATLAS and CMS experiments will record ~ 10 times as much data from ~ 100 times as many collisions as were used to discover the Higgs boson.

ServiceX [1] is an experiment-agnostic service to enable on-demand data delivery, tailored for nearly-interactive high-performance array-based analyses. It provides a uniform backend to data storage services, ensuring the user doesn't have to know how or where the data is stored, and is capable of on-the-fly data transformations into a variety of formats (ROOT files, Arrow arrays, Parquet files, ...) The service offers preprocessing functionality via a simple selection language that allows users to filter events in place, unpack compressed formats, request columns, and specify computations to be applied to the results. This enables the user to start from any format and extract only the data needed for an analysis.



## ServiceX: Columnar data delivery

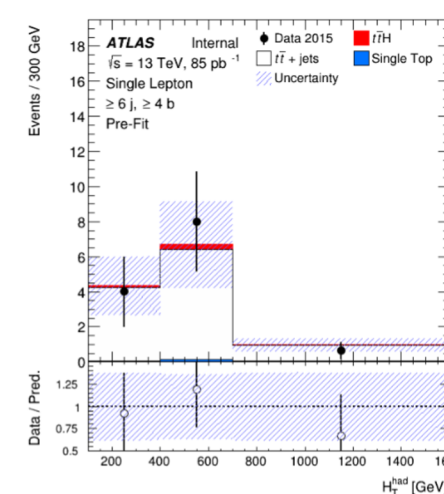


First users specify the needed events and columns, the desired output format, and any required preselection. ServiceX queries the Rucio backend for the data and provides a unique token to the user to identify the request. It accesses the data and performs the required transformation, and tracks the data that's processed and delivered to ensure the transformation completes successfully.

Output is managed by the Kafka message broker, which caches the results to make the available for instant replay.

## Connections to other parts of the IRIS-HEP ecosystem

**ServiceX as an analysis driver:** The service is a highly performant analysis system designed to reduce time-to-insight. Analyzers at the University of Texas are early adopters of (and contributors to) the project, with the goal of using the service to select branches from very large, flat ntuples, applying preselections, and delivering columns to be fitted via TRExFitter.



**ServiceX as a component of the iDDS ecosystem:** The techniques supported by ServiceX are also relevant to iDDS [2], enabling data transformations developed in individual environments to be scaled up to production-based operations.

**ServiceX as a prototyping test case:** The service includes a reproducible pattern for deployment, and the entire project is currently implemented as a central service in a Kubernetes cluster running on the Scalable Systems Lab (SSL) [3], taking advantage of its infrastructure support to develop new features quickly.

Further, **ServiceX interfaces naturally with Skyhook [4]**, with output columns being sent to a programmable object storage system. Development is currently underway to connect these projects, reading Kafka output directly into object storage for analyses.

**ServiceX makes use of Xcache [5]** to accelerate data delivery by caching popular input datasets for efficient running of new transform requests.

## Collaboration across institutions, experiments

The ServiceX project draws expertise from many collaborators across multiple institutions. Contributions to the project code base come from the University of Chicago, University of Illinois at Urbana-Champaign, University of Washington, and Princeton University. The service employs Awkward tools and accommodates use cases that span the ATLAS and CMS collaborations.

## Taking advantage of industry standard tools

ServiceX utilizes a number of industry-standard tools for tasks ranging from creating the REST API to coordinating the transformers to caching the transformation output.



## References

- [1] <https://github.com/ssl-hep/ServiceX>
- [2] <https://iris-hep.org/projects/idds.html>
- [3] <https://iris-hep.org/ssl.html>
- [4] "Skyhook Data Management: Scaling Databases and Applications with Open Source Extensible Storage", Jeff LeFevre, CROSS Research Symposium 2019
- [5] "Creating a content delivery network for general science on the backbone of the Internet using xcaches.", Igor Sfiligoi, CHEP 2019