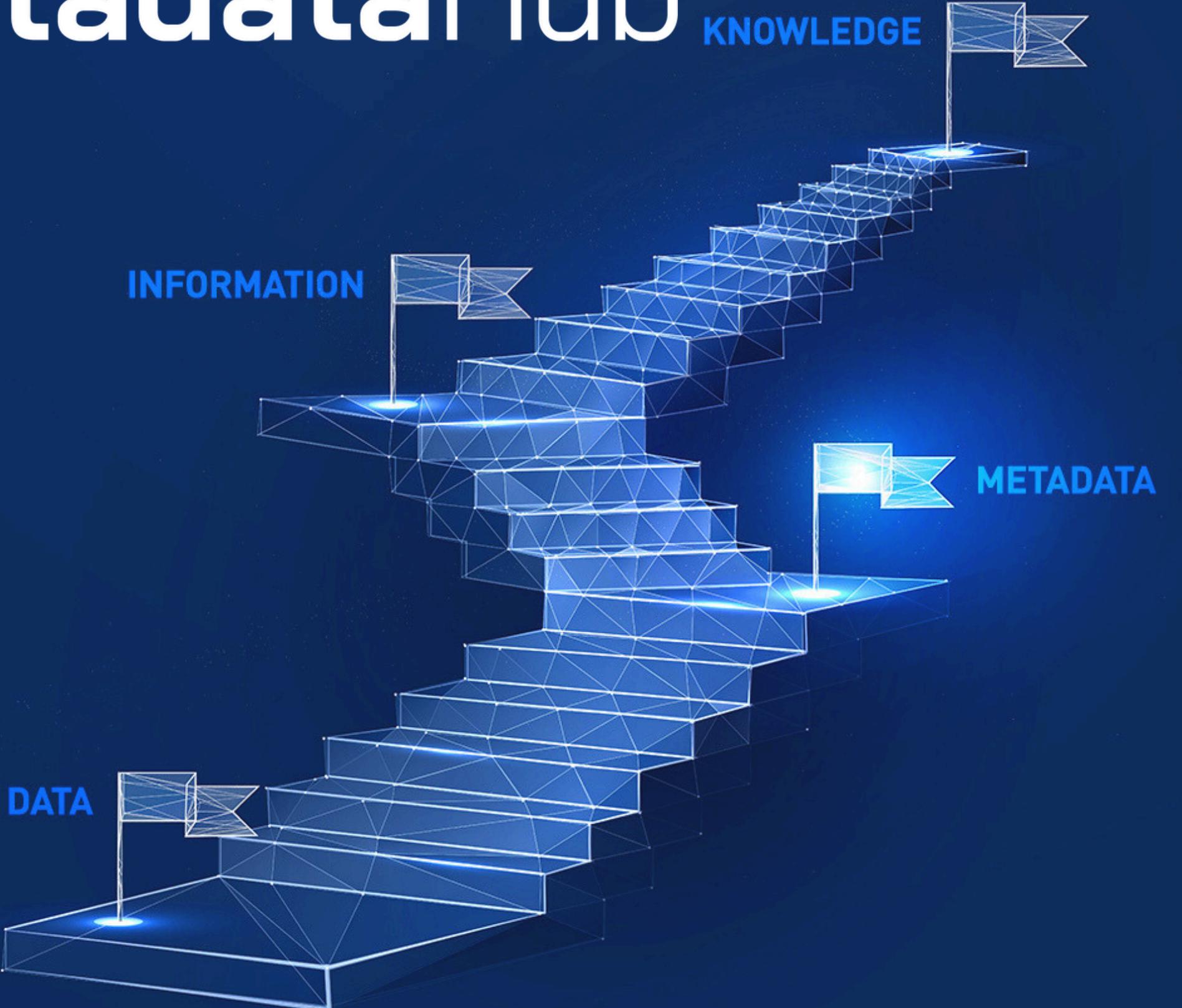




Unleash the Power
of your
Unstructured Data



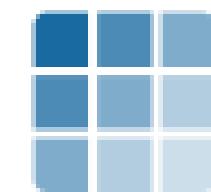
David Cerf
Chief Data Evangelist
david@graudata.us



David Cerf

Chief Data Evangelist
GRAU DATA

DAVID@GRAUDATA.US



GRAU DATA

Your data \ Your control _



Session Goals

- 1. Why is embedded metadata important**
- 2. Metadata's role in a data-driven world**
- 3. How embedded metadata enhances iRODS**

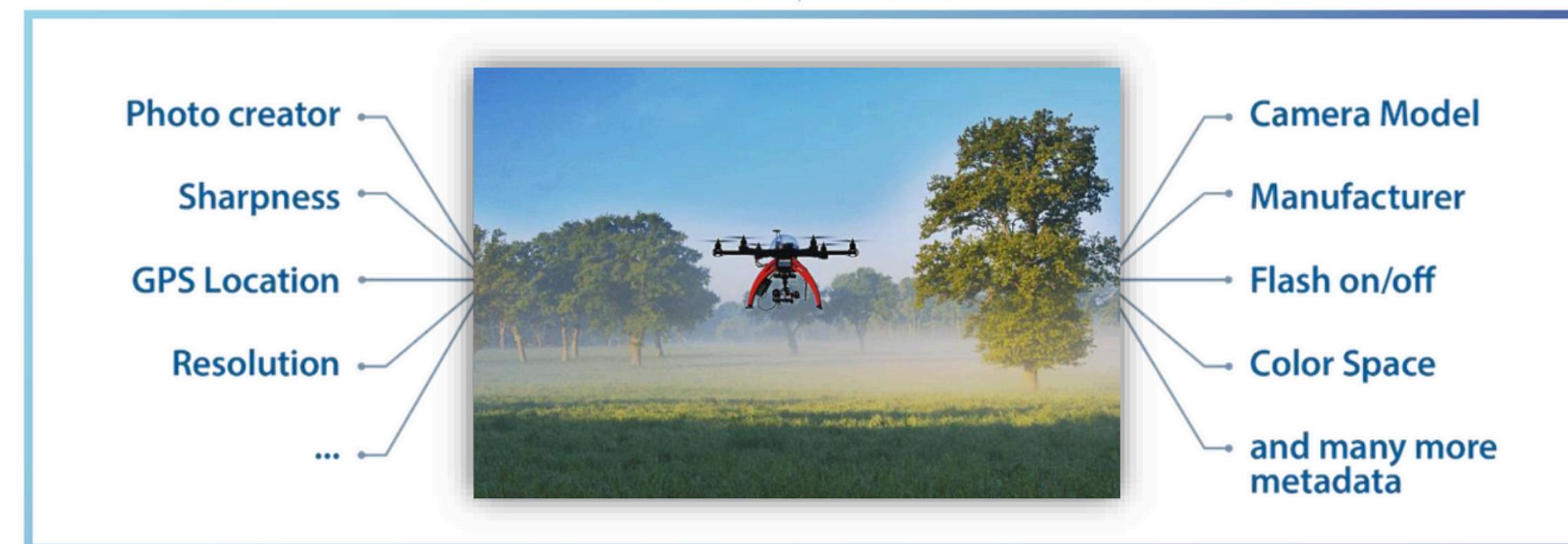
What is Embedded Metadata?

POSIX Metadata

File name, file size, creation time, last access time, modification time, etc.

Embedded File Metadata

- Created by the application, defines the content and substance of the file.
- File can contain hundreds to tens of thousand metadata tags.



Machine-Generated File Content is Critical Embedded Metadata

From DNA, RNA, and Proteins to useful metadata

FASTA, FASTQ, FB Files

The screenshot shows a window titled "ecoli_k12.fasta" containing a FASTA file. The file contains approximately 66,000 lines of sequence data, starting with line 66276 and ending with line 66311. A large blue arrow points from the FASTA file towards the metadata summary on the right.

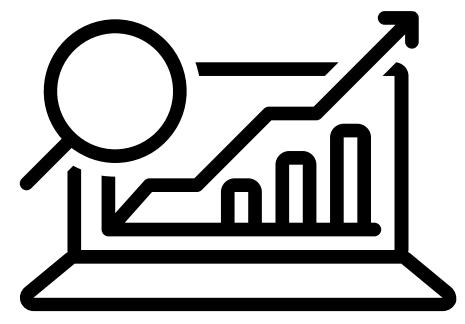
> 66.000 lines

```
66276 TGTGTTATTCGCTTGGCGTTCTGAACATACCGGTTTACCGTGGCATGGATAATGCCAGTCT
66277 GATTGGGCGATAATGAACCGTATTGAAAGGTTGGTCAACAGCAGTGTGTT
66278 ACCCTCGCTGTTGCTGGATGGTGTGATGGGACTGCTCAACTCTGCCATAGCGCGCTGTTGT
66279 TGGGAACAGTGTGCTGGTGGCGCATGATGTTGTGACCCGTAATATCGACTGGTATGC
66280 GTTTCACTGCCGAAATGAAAGCCAGTAAAGAACGTTACAACCGACGATGAGTTACGTATCTGAAATAA
66281 GGTTGAAAAATAAAACGGCTAAAAAGCGCCGTTTTTGACGGTGGTAAAGCGATTAACTTCCA
66282 GATCACCGCAGAACGATAACCTCACCGTGAATGGTGGCGATGATTCCGGTAGATT
66283 GAAATGTTACGAATACGGCGATCGTCAGTCTACAGTACGGTGTGCGGTTACGTCACGGCGGT
66284 ATTTCTTCAGCAGTCAGCACGGGACTGAATTTGCCCTGGGTTTACAGAAGTGAAGCATGGCGGGA
66285 ACTCGCTGCCGGCAGCTGTACTGCTGCCATCAGGGCGATCAACGAACGGCTGTTGATGCCAGTC
66286 CCAACCATTGAACTGTAGCTTCAACGCTACGACGTTCTCGCTGACAGTACCCAGATTGATGGTACGG
66287 GACAGTAGGTTGGCGTCAGAACGTCAGTTCACGGGGTTGAACGGTTGGTGTAGTGTAGTCATGAC
66288 CGATTTCGAGGCCGAGAATTATCGACTTCGTTGTACGGCCAGTCAGGAACATCAACGCAACATTGCG
66289 CTGCTCGCAGTCACCGCTAACAGAACGGCTTACCCGGCAGATTGATATCCATGATCACCAGG
66290 TTGATGTCATATTCAAGAGGAGCTGATGCAATTCCGCCATCTGCGCTTCAAACATCATAGCCTT
66291 CCGCTCGAAAATACTTCAACGTTGCTGTTACCAACTCGTCTTCAACGATAAGAATGTCGGGGGT
66292 CTGCATGTTGCTACTAAATGCCAACTAAATCGAAACAGGAAGTACAAAAGTCCCTGACCTGCTGAT
66293 GCATGCTGCAAATTAAACATGATCGCGTAACATGACTAAAGTACGTAATTGCGTTCTGATGACCTTCC
66294 ATCAACGTCACACATCATTAGCTTGGTGTGGTACTTCCCTCAGGACCCGACAGTGTCAAAACGG
66295 CTGTCATCTAACCAACATAACAGGCTAACAGGCTAACAGGGGGGGACACCCAAATAAACTACGCT
66296 TCGTTGACATATATCAAGTTCAATTGTTAGCAGCTAACAGTTGATGAAATCATCGTATCTAAATGCTAG
66297 CTTTCGTCACATTATTTAAATCCAACACTAGTTGTCATCATACAACATAAAACGTTGAATCCAATTG
66298 TCGAGATTATTTTATAAAATTATCTAAGTAAACAGAAGGGATATGTTGACATTAAACACTCAAC
66299 CGTTAGTACAGTCAGGAAATAGTTAGCCTTTTAAGCTAACGTAAGTAAAGGGTTTCTGCGACTTACGTT
66300 AAGAATTGTAATTGCAACCGCGTAATAAGTTGACAGTGTGATCACCCTGGTTCGCGGTTATTGATCAAGA
66301 AGAGTGGCAATATGCGTATAACGATTATTCTGGTGCACCCGCCAGAGCAGAAAATATTGGGCAGCGGC
66302 GCGGGCAATGAAAACGATGGGTTAGCGATCTGCGGATTGTCGATAGTCAGGCACACCTGGAGGCCAGCC
66303 ACCCGCTGGGCGCATGGATCTGGTGTATTGATAATTTAAAGTTCCGACATTGGCTGAAT
66304 CGTTACAGATGTCGATTTCAGTGCGCCACACTGCGCAGTCGGCGAAATATTCAACTACGCCAC
66305 GCCAGTTGAACTGGTGCCTTAGAGGAAAATCTTCATGGATGAGCCATGCCGCGCTGGTGTGGT
66306 CGCGAAGATTCCGGGTTGACTAACGAAGAGTTAGCGTTGGCTGACGTTCTACTGGTGTGCCAGTGGTGG
66307 CGGATTATCCTTCGCTCAATCTGGGGCAGGCCGTATGGTCTATTGCTATCAATTAGCAACATTG
66308 ACAACCGCGAAAAGTGTGCAACGGCAGACCAACATCAACTGCAACGAGGACATGACA
66309 TTGCTGACGACTCTGGCAGTGGCAGATGCAACGAGGACATGACA
66310 TAGAGCAACGAGCAGACCAACATCAACTGCAACGAGGACATGACA
66311 ACGCTTGTAGTAAGTATTTTC
```

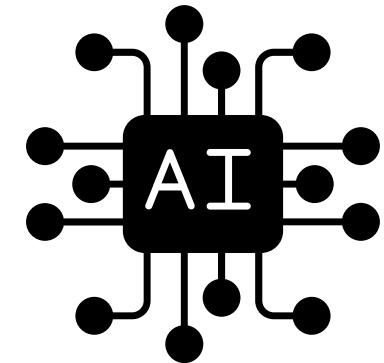
Metadata

GC(%) = **50.79**
N50 = **4.641.652**
Q1 = **2.320.826**
Q2 = **4.641.652**
Q20(%) = **0**
Q3 = **2.320.826**
Q30(%) = **0**
U00096_3.Note = **Escherichia coli str. K-12 substr. MG1655**
avg_len = **4.641.652**
max_len = **4.641.652**
min_len = **4.641.652**
num_seqs = **1**
sum_gap = **0**
sum_len = **4.641.652**
type = **DNA**

Embedded metadata is critical to data-intensive industries and for AI



Analytics



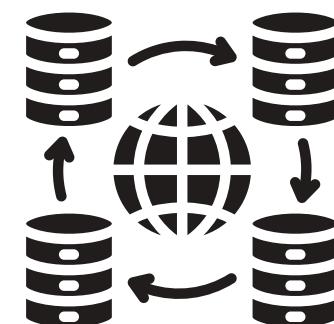
AI



Discovery



**Workflow
Automation**



**Data
Orchestration**

**AI and ML have changed the
data landscape**

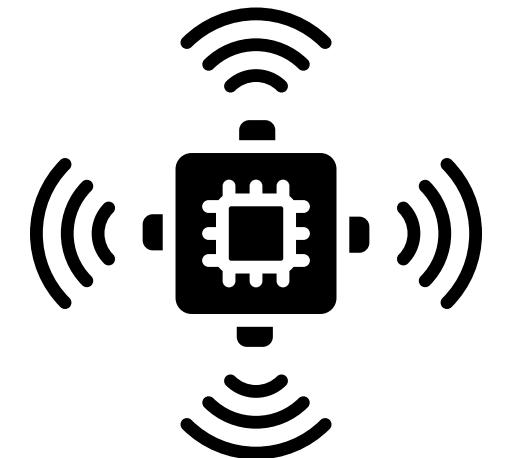
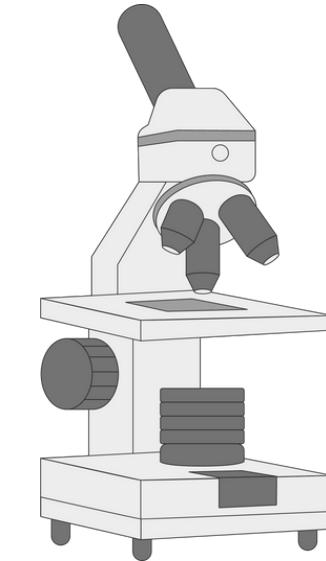
All data is needed all the time



The Machine-Generated Data Problem

- 98% of all new data is machine-generated (by volume i.e., TB)
- Not typically read by humans but for computation

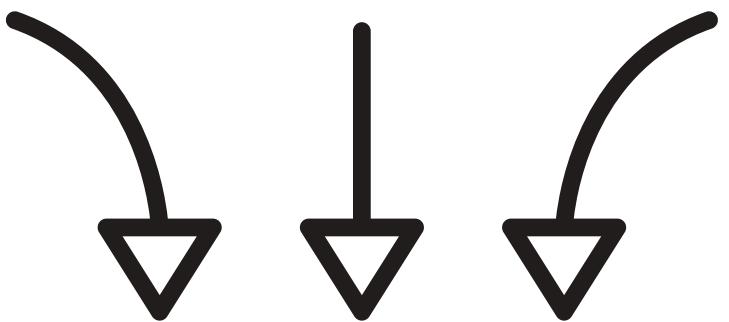
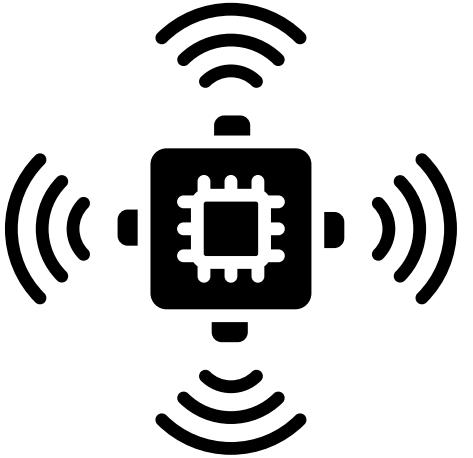
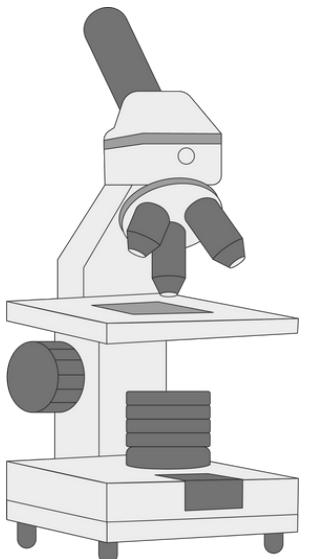
Data Generators



The Machine-Generated Data Problem

Machine-generated data resides on expensive storage

Data Generators



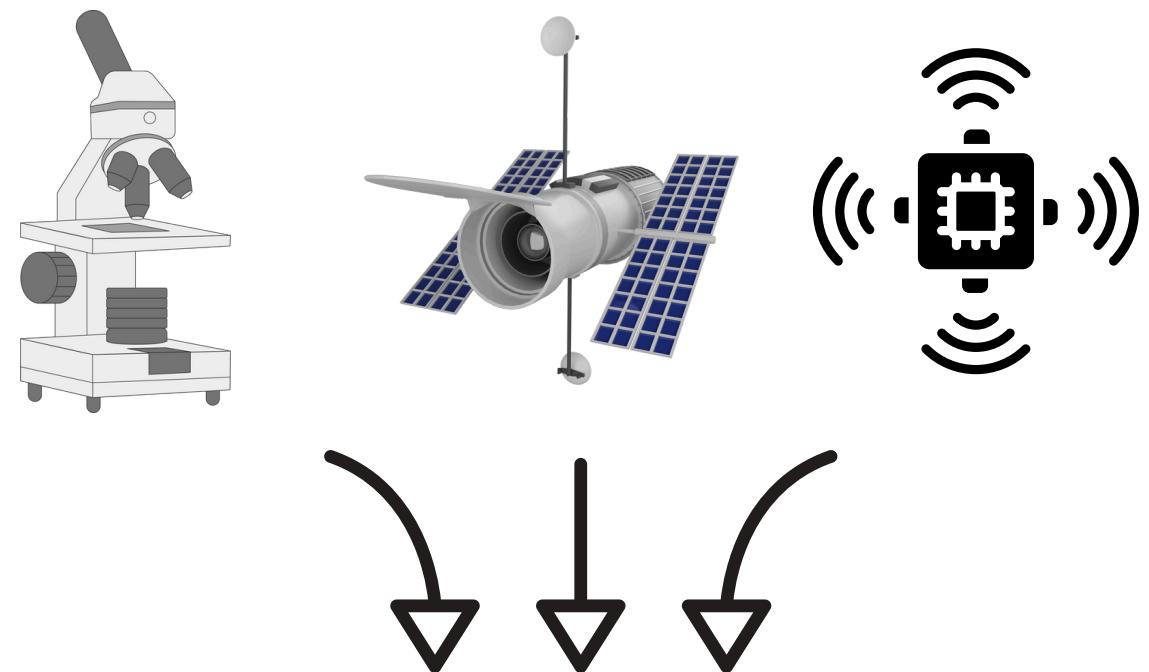
Primary Storage

The Machine-Generated Data Problem

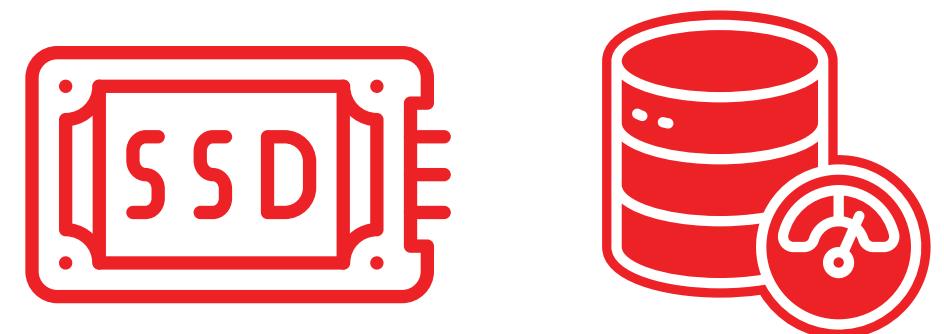
**Copies can go to archive
(cheap storage)**
**but machine-generated data
is still on primary storage**

- Critical for AI / ML,
analytics
- Requires fast access

Data Generators



Primary Storage



Archive / Copy

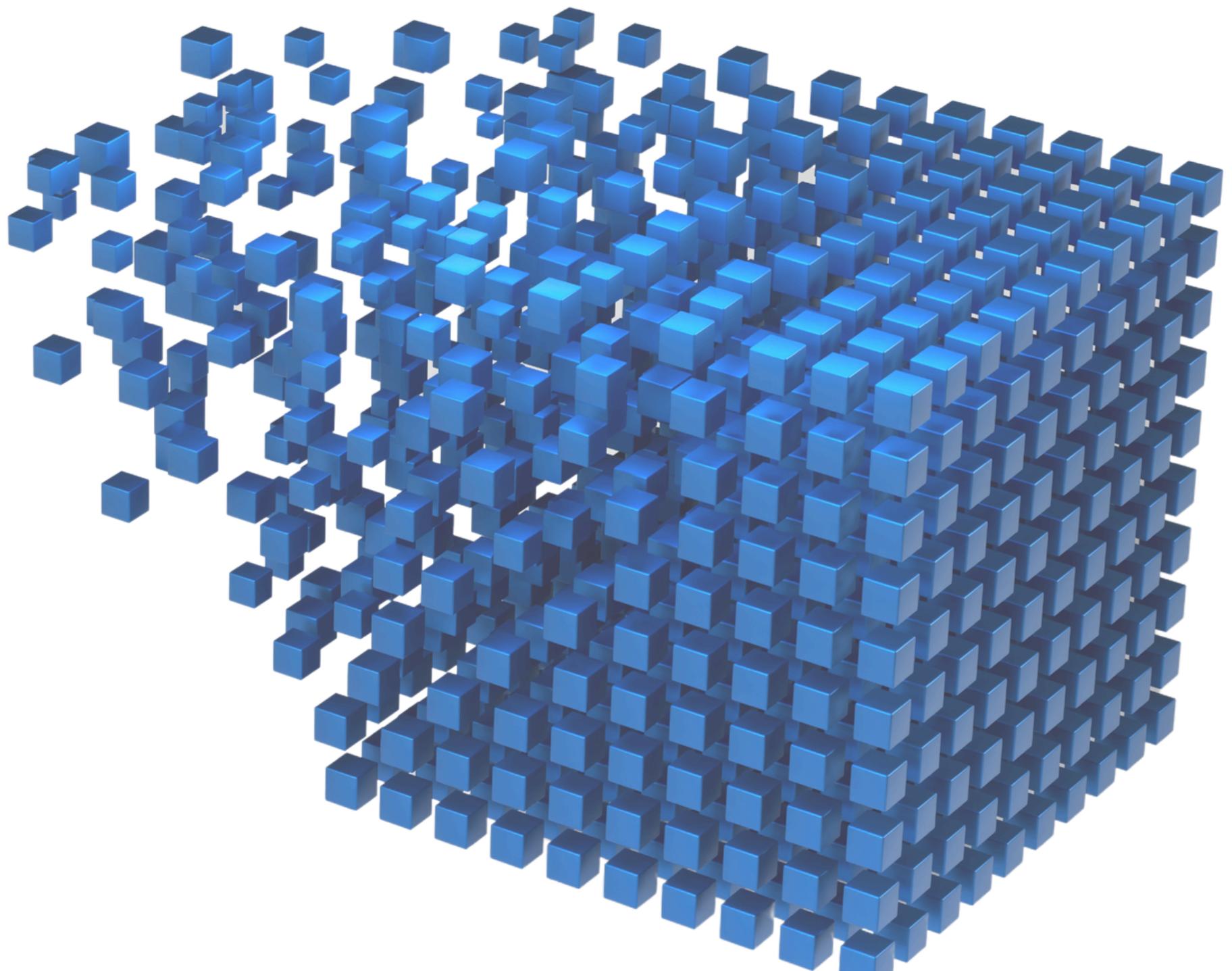


Unleash the Embedded Metadata



MetadataHub

Transforms unstructured data
into structured, easily searchable
data sets by harvesting the
embedded metadata



The Challenge of Diverse Unstructured Data Formats

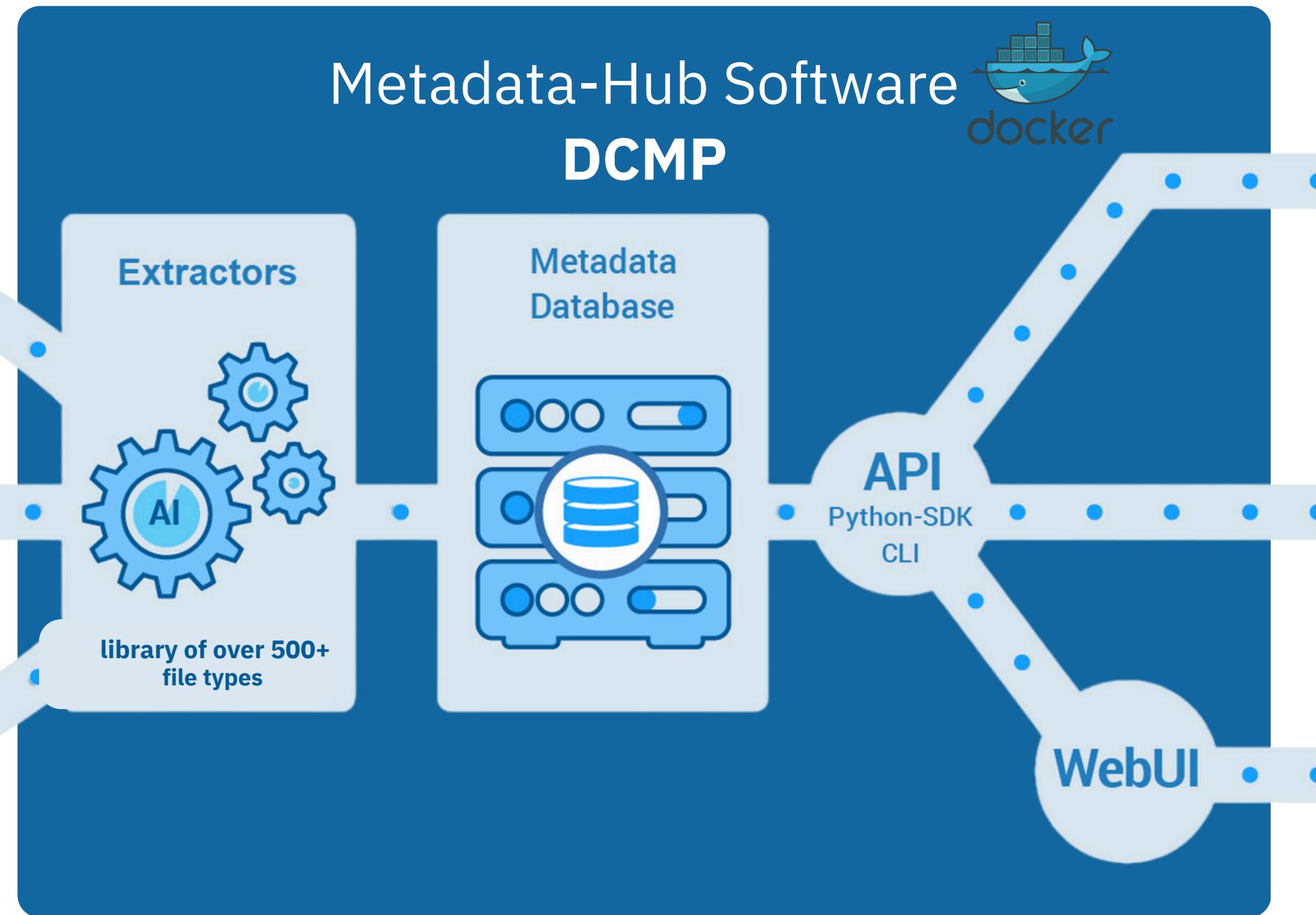
Metadata-Hub support 100's of unique file formats

360 3FR 3G2 3GP 3GP2 3GPP 1SCA AA AAE AAX ACFF ACFM ACR AFM AI AIF AIFC AIFF AIT AIM AL3D ALI AMFM
APE APNG ARF ARQ ARW ASF AVI AVIF AZW AZW3 BMP BPG BTF CDF CH5 CHM CIF CIFF COS CR2 CR3 CRM CRW
CS1 CSV CXD CZI DC3 DCM DCP DCR DFONT DIB DIC DICM DIVX DJV DJVU DLL DM2 DM3 DM4 DNG DOC DOCM
DOCX DOT DOTM DOTX DPX DR4 DS2 DSS DV DVB DVR-MS DYLIB EIP EMI EPS EPS2 EPS3 EPSF EPUB ERF ETS EXE
EXIF EXR EXV F4A F4B F4P F4V FFF FIT FITS FLA FLAC FLEX FLI FLIF FLIR FLV FRM FPFFPX FTS GEL GIF GPR GZ GZIP
H5 HDP HDR HEIC HEIF HIF HIS HTM HTML I2I ICAL ICC ICM ICS IDML IIQ IMG IMS IND INDD INDT INR INSP INSV
INX ISO ITC J2C J2K JNG JP2 JPC JPE JPEG JPF JPG JPM JPS JPX JSON JXLJXR K25 KDC KEY KTH LA LFP LFR LIF LNK
LRV M2T M2TS M2V M4A M4B M4P M4V MACOS MAX MDF MEF MIE MIF MIFF MKA MKS MKV MNG MOBI MODD
MOI MOS MOV MP3 MP4 MPC MPEG MPG MPO MQV MRC MRW MSR MTS MXF NAF ND2 NEF NEWER NII NKSC
NMBTEMPLATE NRW NUMBERS O ODB ODC ODF ODG ODI ODP ODS ODT OFR OGG OGV ONPOPUS ORF ORI OTF PAC
PAGES PBM PCD PCT PCX PDB PDF PEF PFA PFB PFM PGF PGM PICT PLIST PMP PNG PNL POT POTM POTX PPAM
PPAX PPM PPS PPSM PPSX PPT PPTM PPTX PRC PS PS2 PS3 PSB PSD PSDT PSP PSPFRAME PSPIMAGE PSPSHAPE
PSPTUBE QIF QT QTI QTIF R3D RA RAF RAM RAR RAW REC RIF RIFF RM RMVB RPM RSRC RTF RV RW2 RWL RWZ
SCN SEQ SER SIF SKETCH SLD SO SR2 SRF SRW SVG SVS SWF THM THMX TIF TIFF TORRENT TS TTC TTF TUB TXT
VCARD VCF VOB VRD VSD VSI WAV WDP WEBM WEBP WMA WMV WOFF WOFF2 WTV WV X3F XCF XDCE XHTML
XLA XLAM XLEF XLS XLSB XLSM XLSX XLT XLTM XLTX XMP XQD XQF XV ZIP ANY MANY MORE ...

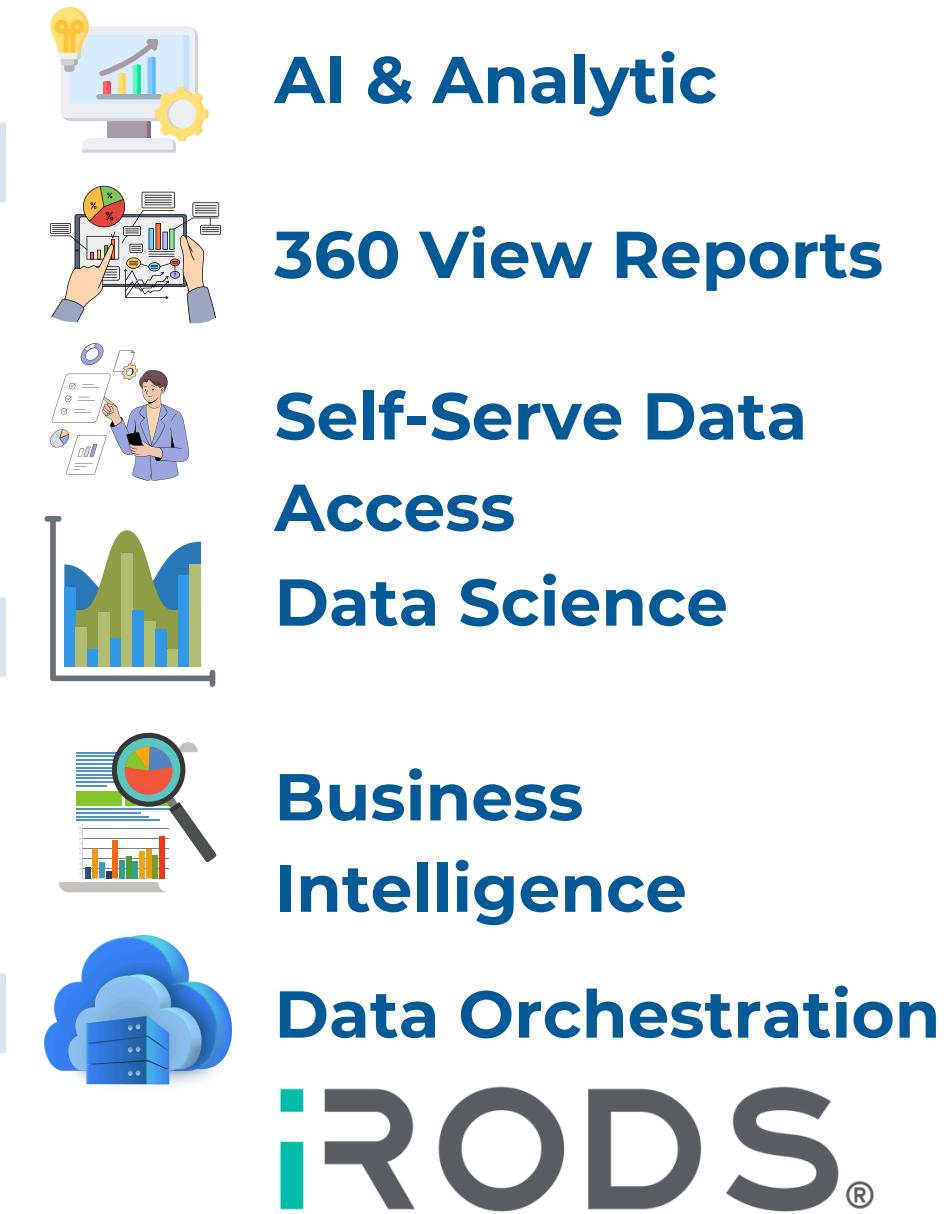
NOTE: ALL SUB-FORMATS OF THE LISTED FILES ARE SUPPORTED TOO → FOR EXAMPLE: JPG ALONE HAS [41](#) SUB-FORMATS

Metadata-Hub: Creating a Metadata Fabric

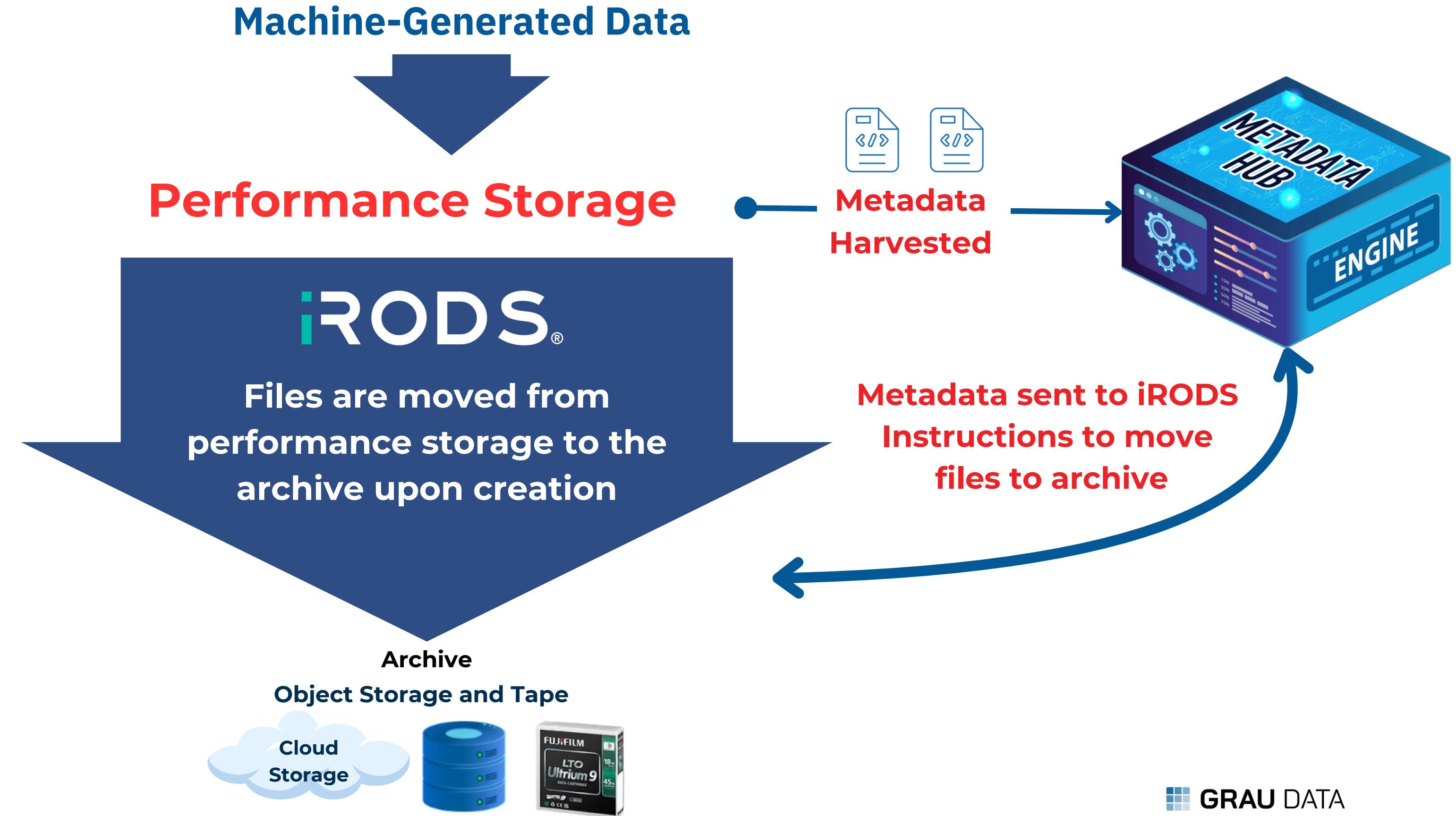
Storage



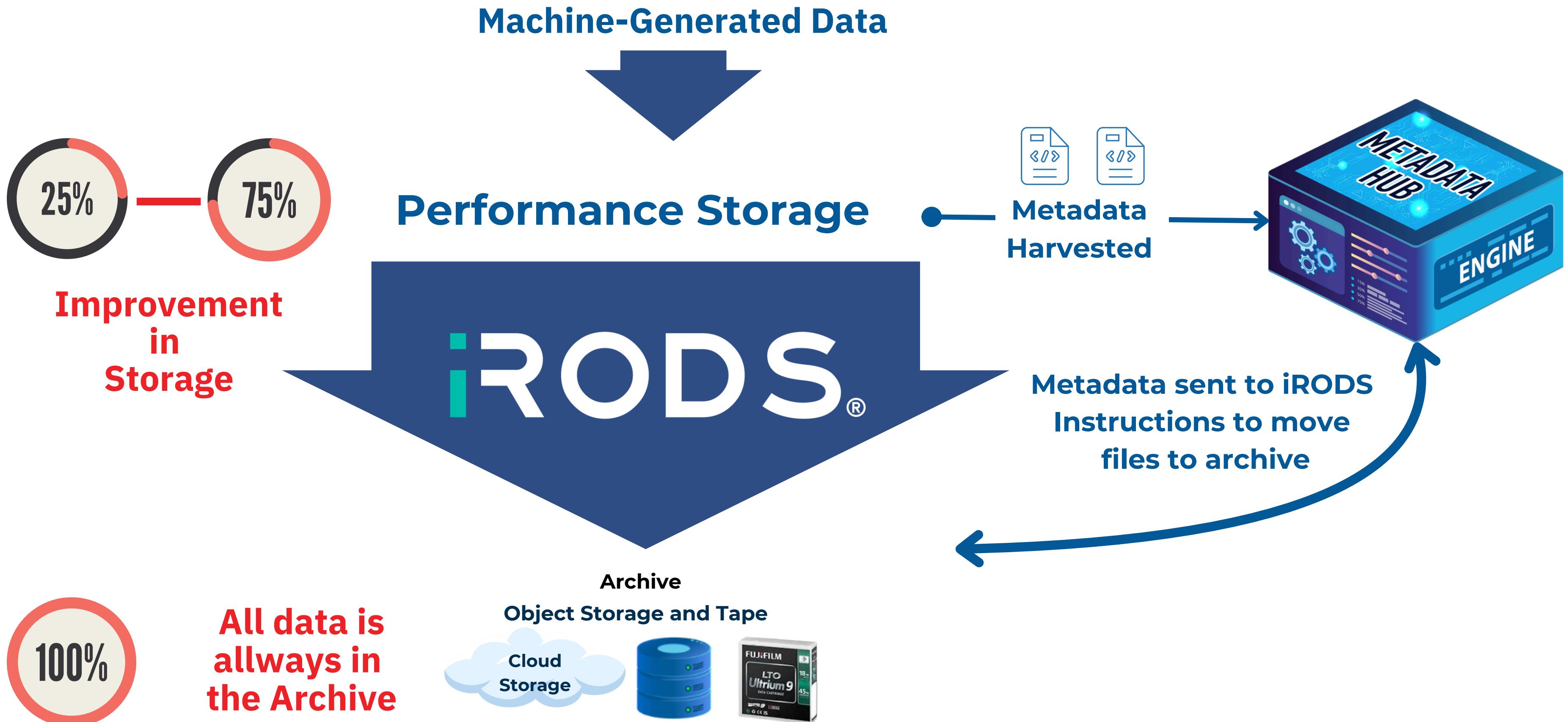
Data Users



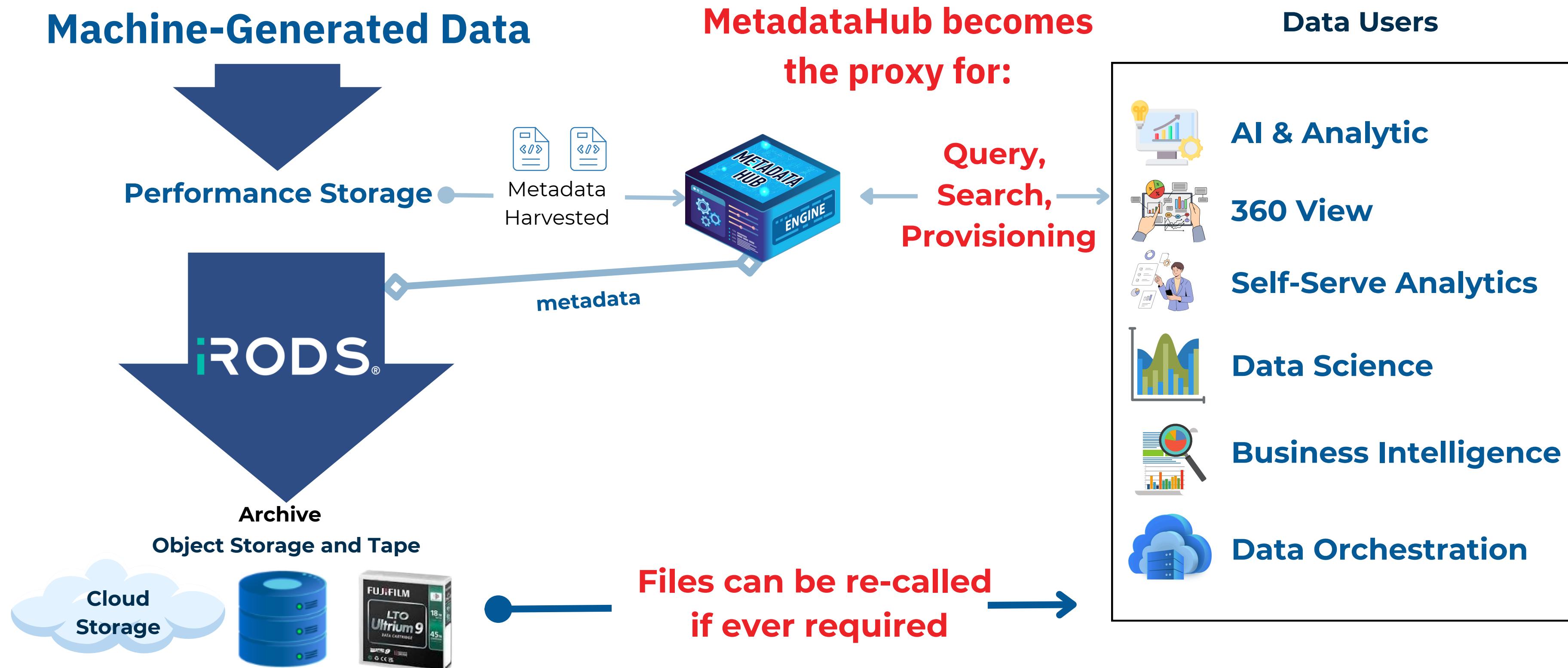
iRODS Metadata Powered Workflow



MetadataHub Helps Free Disk Capacity



Feeding Data to Data Users



S3 Glacier Tape Archive



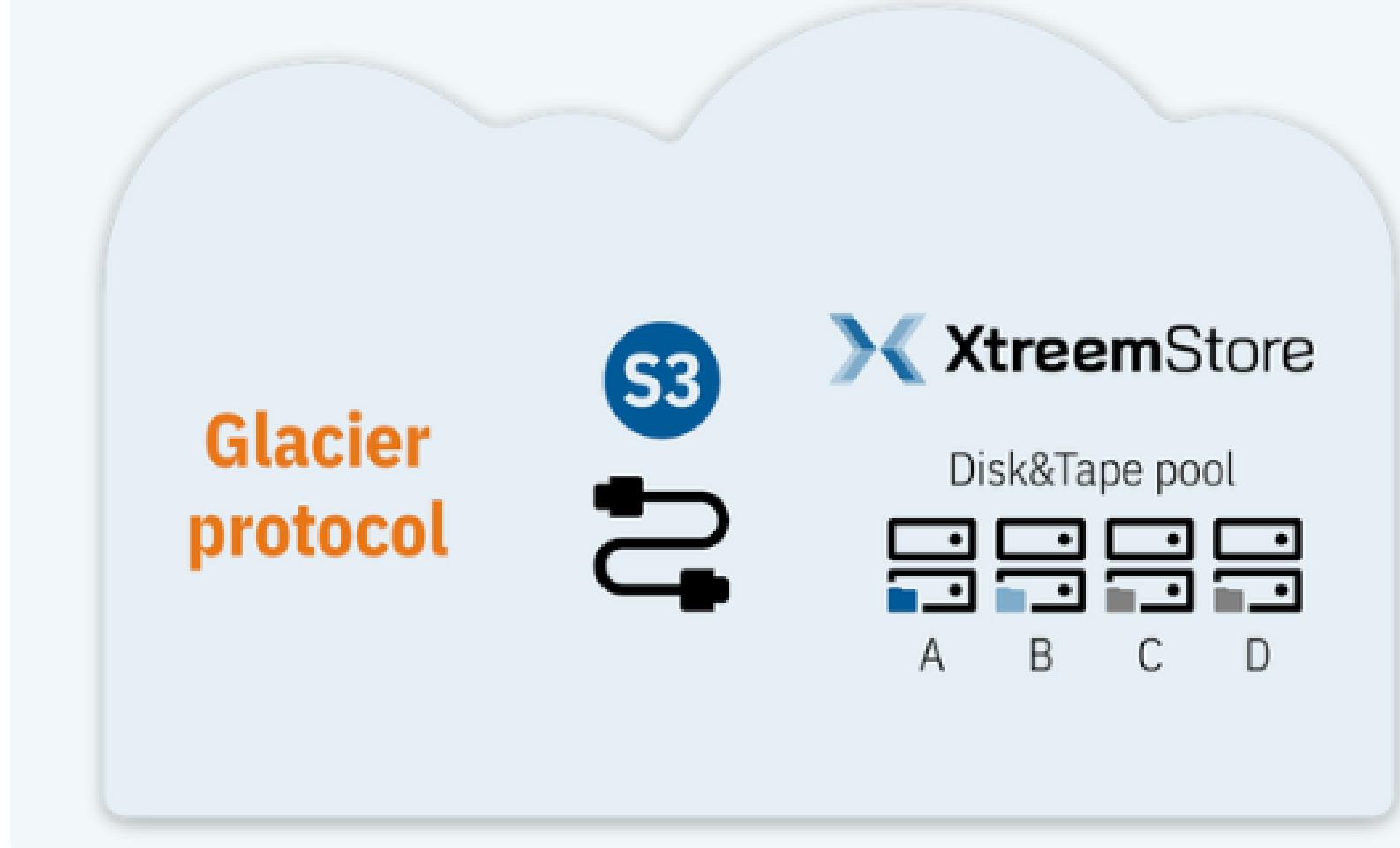
S3 / Glacier



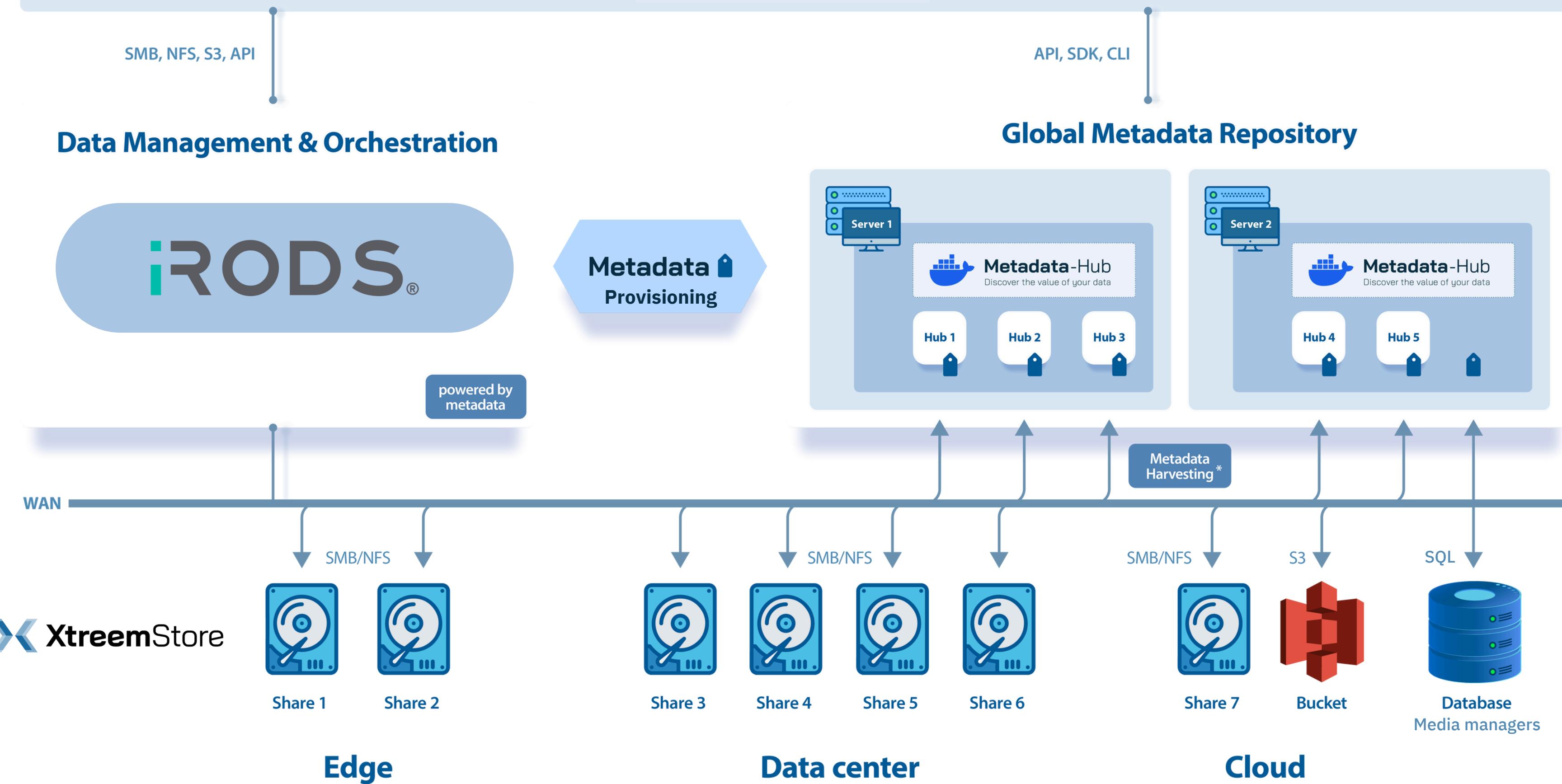


TURN ANY TAPE LIBRARY INTO A GLACIER CLOUD ARCHIVE

- Simplicity - S3 & Glacier compatible
- Open standards and vendor-neutral hardware
- No hidden costs



Applications, Data Lakes, Reporting, etc



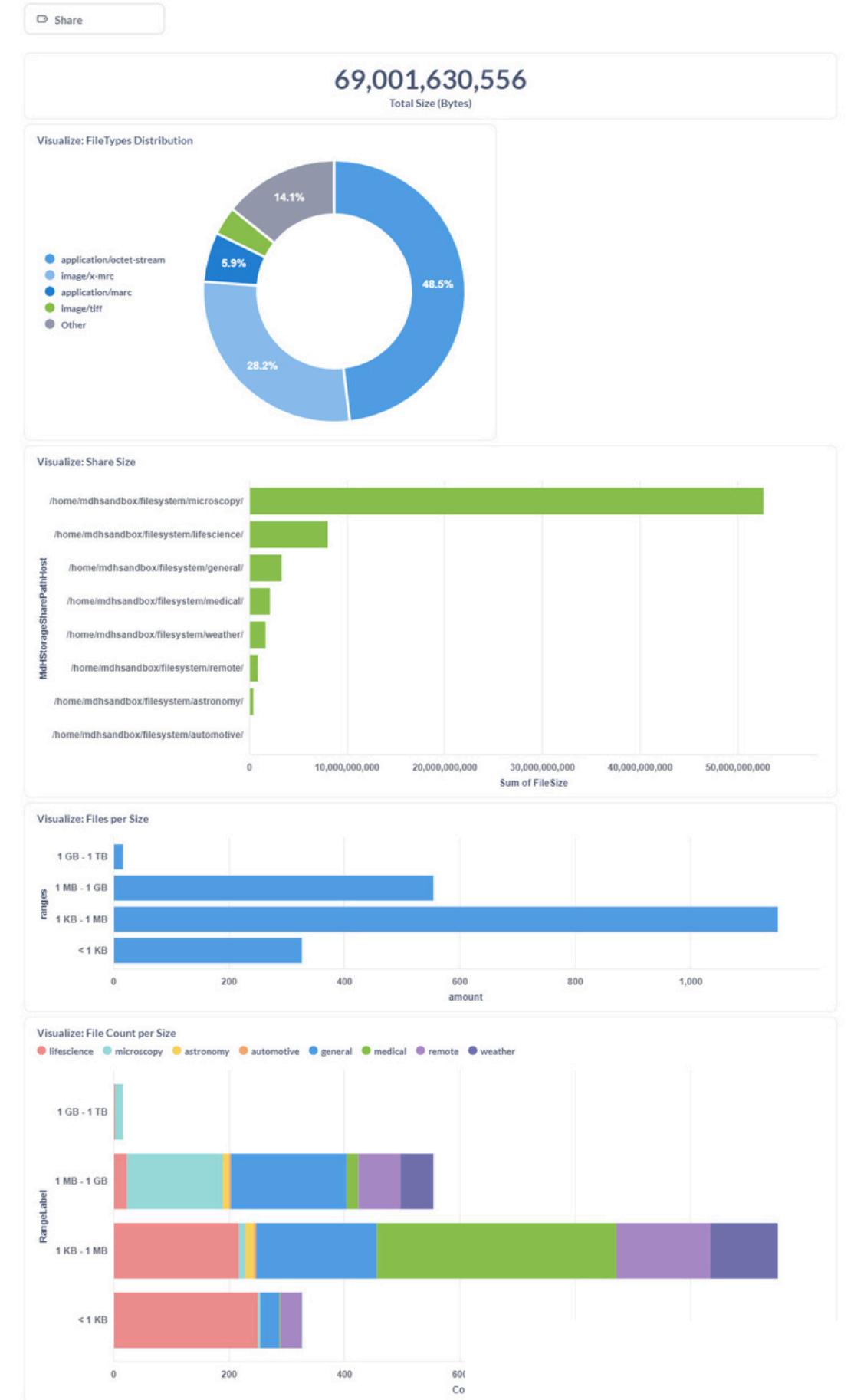
DATA LANDSCAPE REPORTS

Storage Report

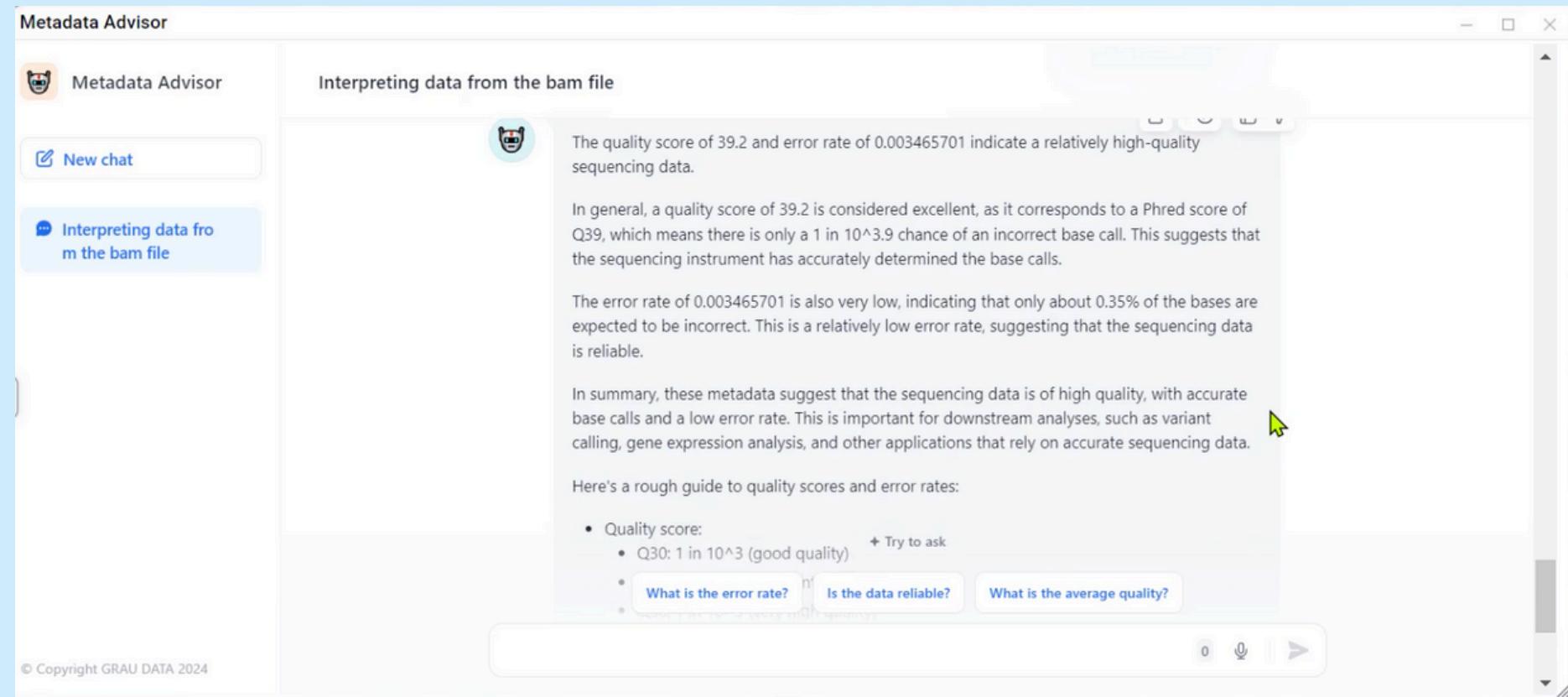
Gain Unparalleled Visibility into Your Data Assets

Unified view of data across cloud, data centers, and edge devices

Insights into data distribution, storage utilization, and growth trends



METADATA ADVISOR



AI-powered features to provide users with advanced insights into their unstructured data and metadata.

Provides expert guidance to:

- **Interpretation** metadata extracted by Metadata-Hub.
- **Confidence in Data Usage:** Make informed use of your data by ensuring a clear grasp of its context and significance.

Core-Instance: Essen

Dashboard Harvest Metadata Search Files via Metadata System

Simple Search Advanced Search Load query Explorer Clear Search

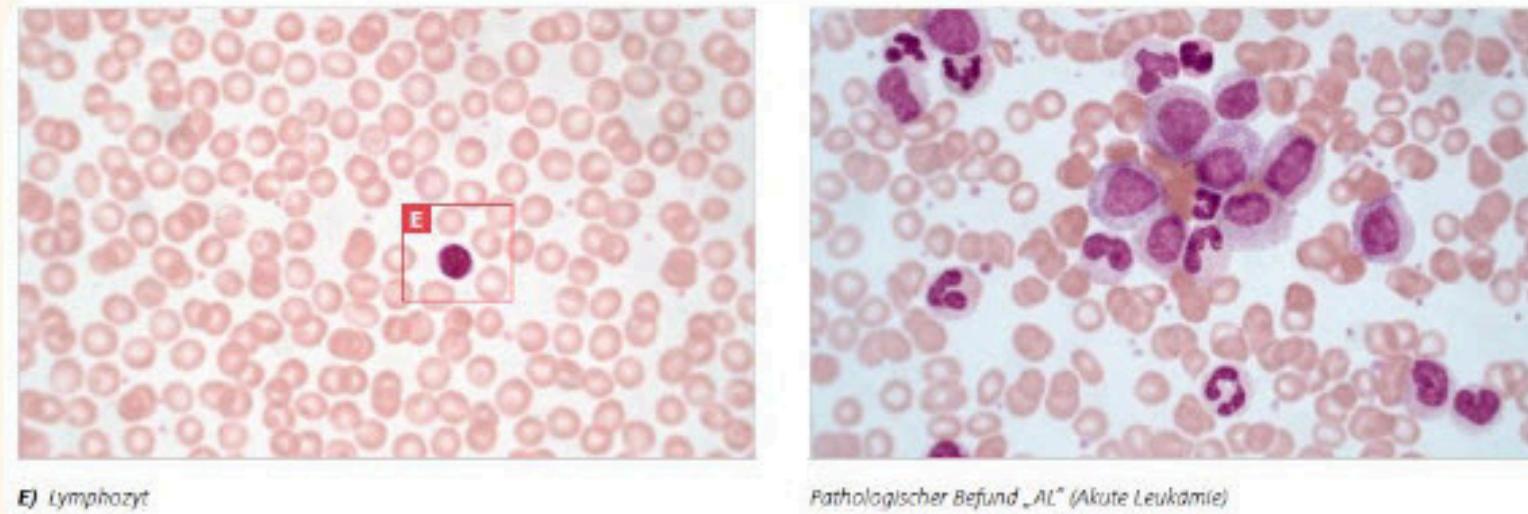
Metadata tag contains Value

Metadata tag contains Value

Search logic All AND Help

First Prev 1 Next Last Records per page: 25





Prompt:

Describe the images and compare them.

The image shows two microscopic views of blood smears, which are used for various medical tests. The left side is labeled "E) Lymphozyt" (German for "lymphocyte"), indicating that the image is focused on lymphocytes, a type of white blood cell. The right side is labeled "Pathologischer Befund, AL (Aktue Leukamie)" (German for "Pathological finding, AL (Acute Leukemia)"), suggesting that the image is showing a pathological finding related to acute leukemia.

In the left image, you can see a cluster of lymphocytes, which are typically round or oval-shaped cells with a distinct nucleus. In the right image, there is a more disorganized pattern of cells, with some cells appearing abnormal or atypical, which could be indicative of a leukemia diagnosis. **The cells in the right image are more irregular in shape and size, which is characteristic of leukemia cells.**

These images are likely used for diagnostic purposes in a clinical setting to identify and classify different types of blood cells and to detect abnormalities that could be indicative of a disease like leukemia.

MetaSens:

Creates text descriptions and keywords from files content.

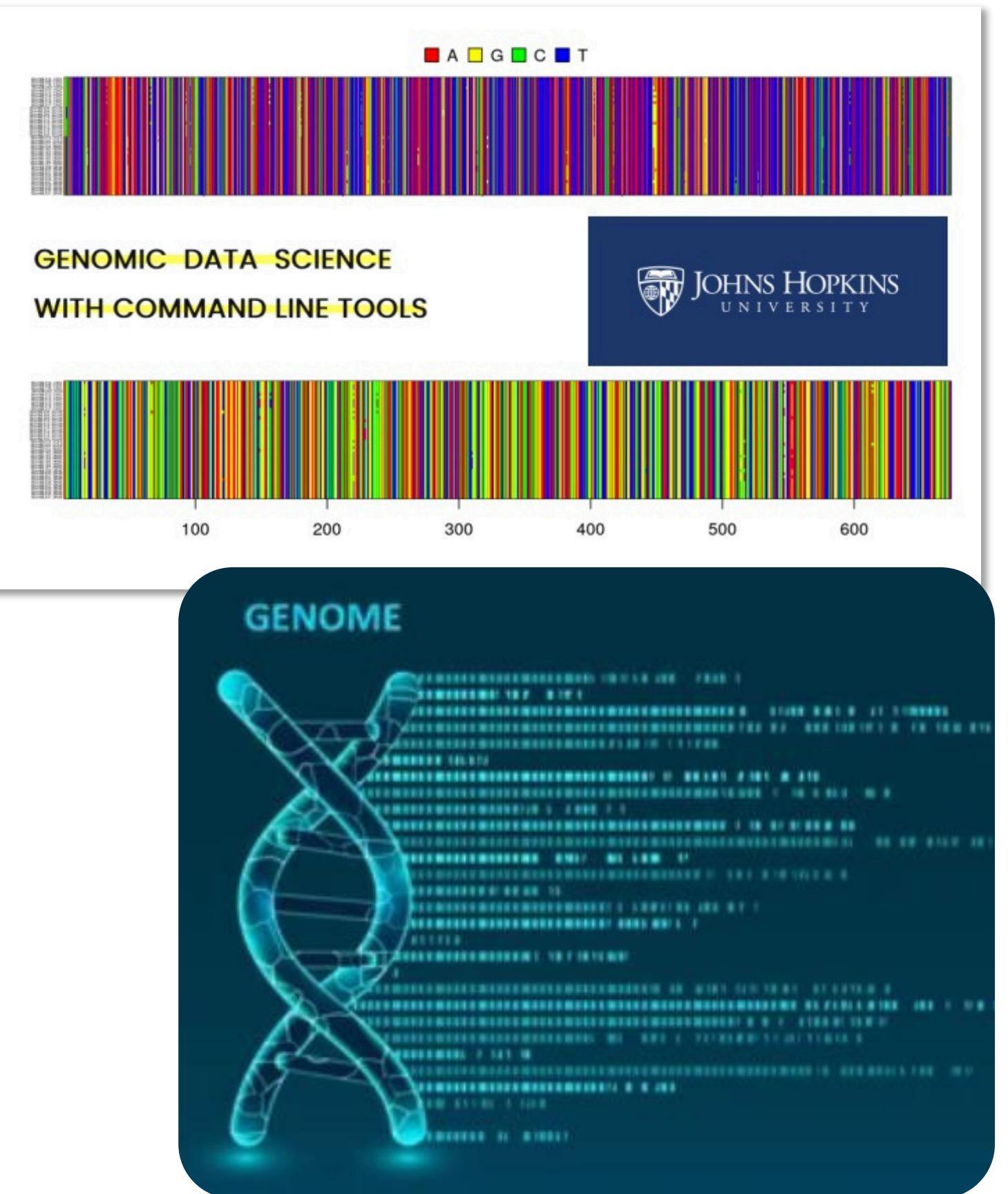
Use Cases



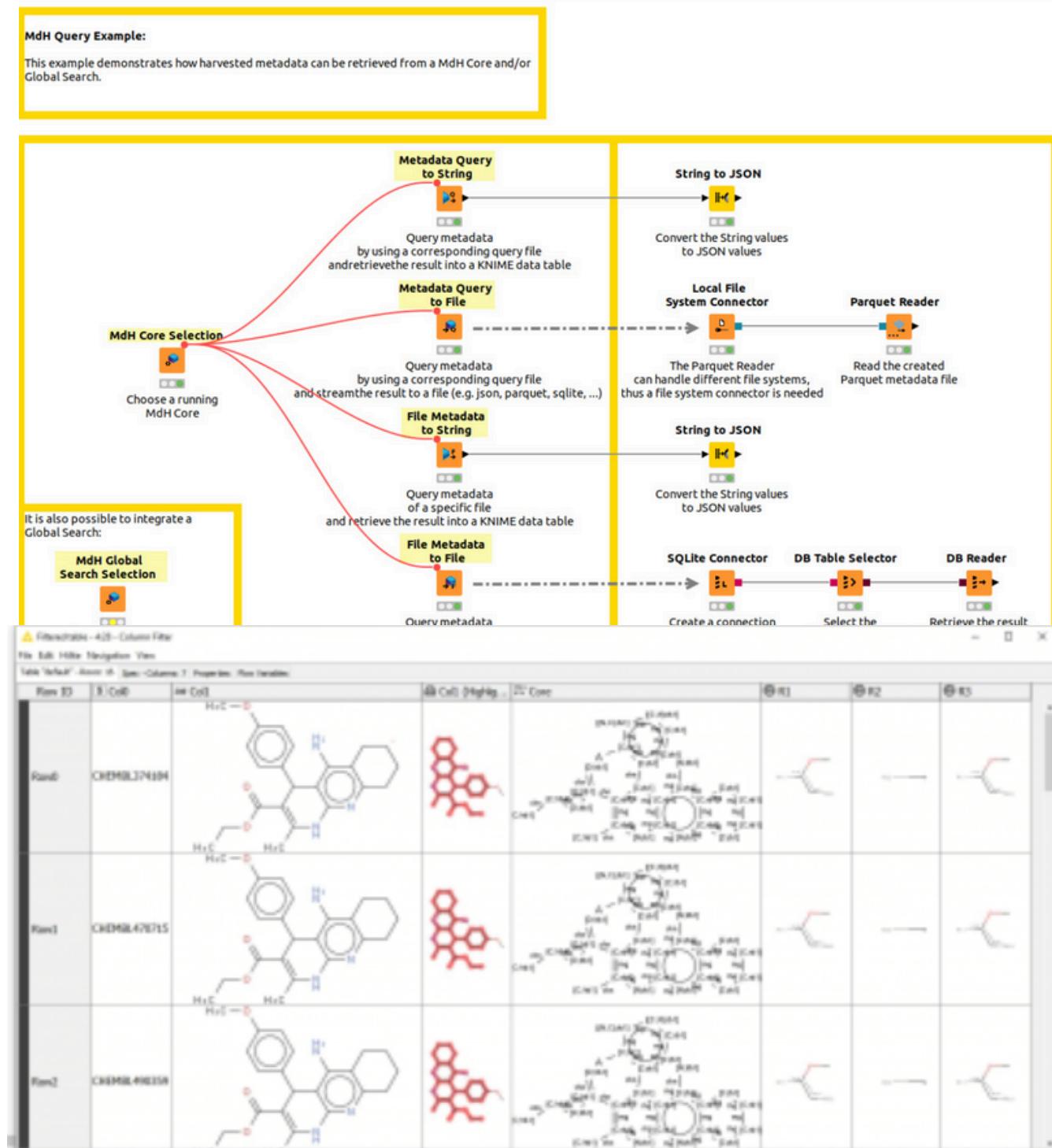
Use Case: ZUSE Institute Clinical Center

Challenges

- 1. Improve data accessibility, search and analysis**
- 2. Support Diverse Format Extraction:**
 - **Medical Imaging:** NIFTI (brain scans), DICOM (CT & MRI images)
 - **Genomic Formats** (FASTA, FASTQ, SAM, VCF): enable precise genetic data analysis
 - **Proteomics Format** (PDB): in-depth protein structure studies
- 3. Enable collaboration in accordance with FAIR practices**
- 4. Reduce Data Preparation**
- 5. Security & Governance**



ZUSE Institute - Management of Research Data / Life Science Data



Results

Enhancing Workflow Automation:

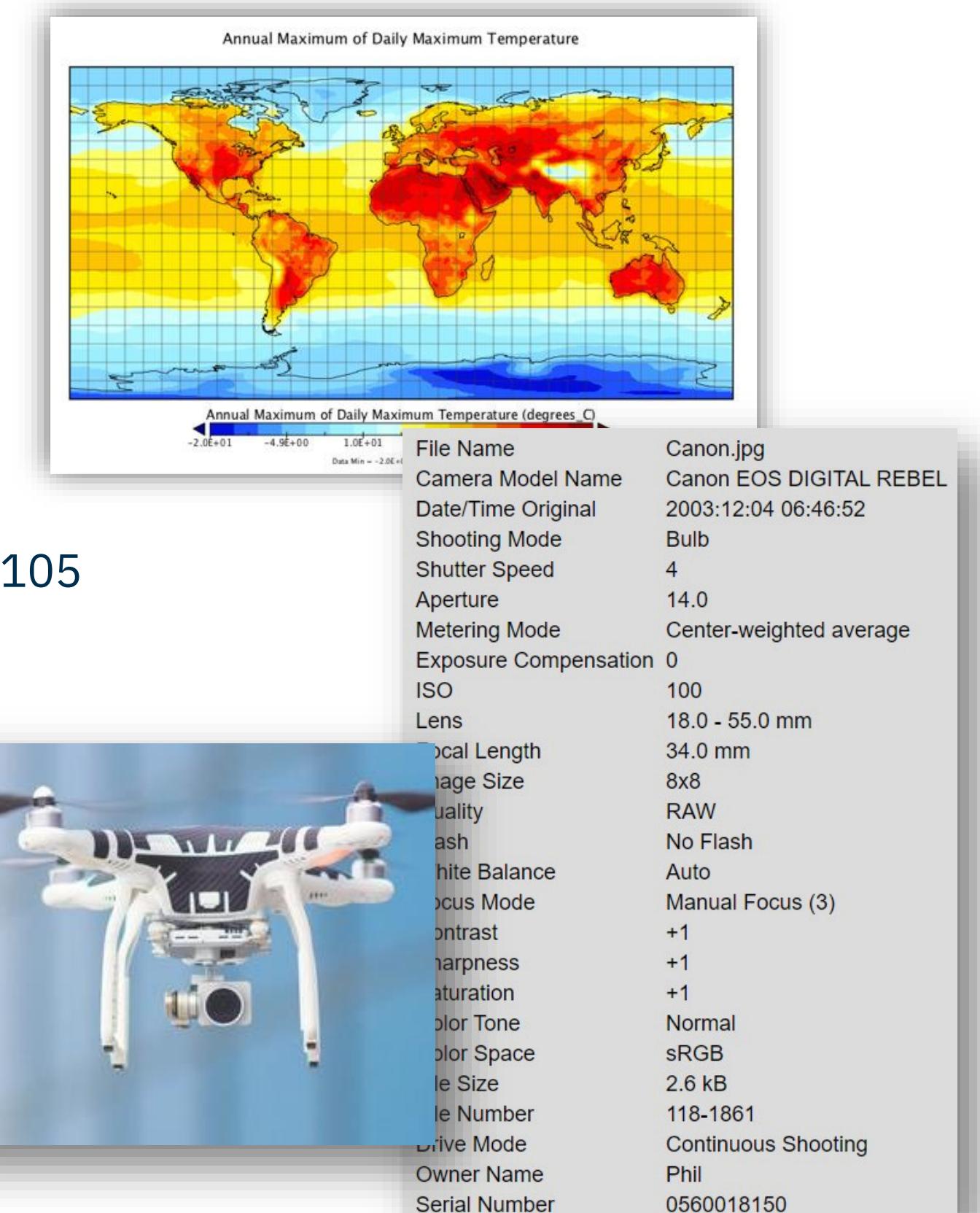
- **Streamlined data preparation** saving countless hours
- **Global data accessibility** to all metadata enforced by access controls
- **Accelerated insight**

Optimizing Research Data Management:

- Holistic Data Landscape View
- Reduced Data Transfers
- Improved storage utilization

Use Case: Weather & Climate

- Climate research -NetCDF , GRIB, BUFR, ZARR, GRID
 - File format which saves multidimensional scientific data such as temperature, humidity, pressure, wind speed and wind direction
- Drone shots –RAW Image Format
 - Master metainformation from around 7000 different camera models of 105 camera manufacturers
 - Used in various industries such as forensics, geo-information services, logistic, agronomics usw.



Challenges

Initial 45 PB of storage (telescope images, blueprints, etc.).

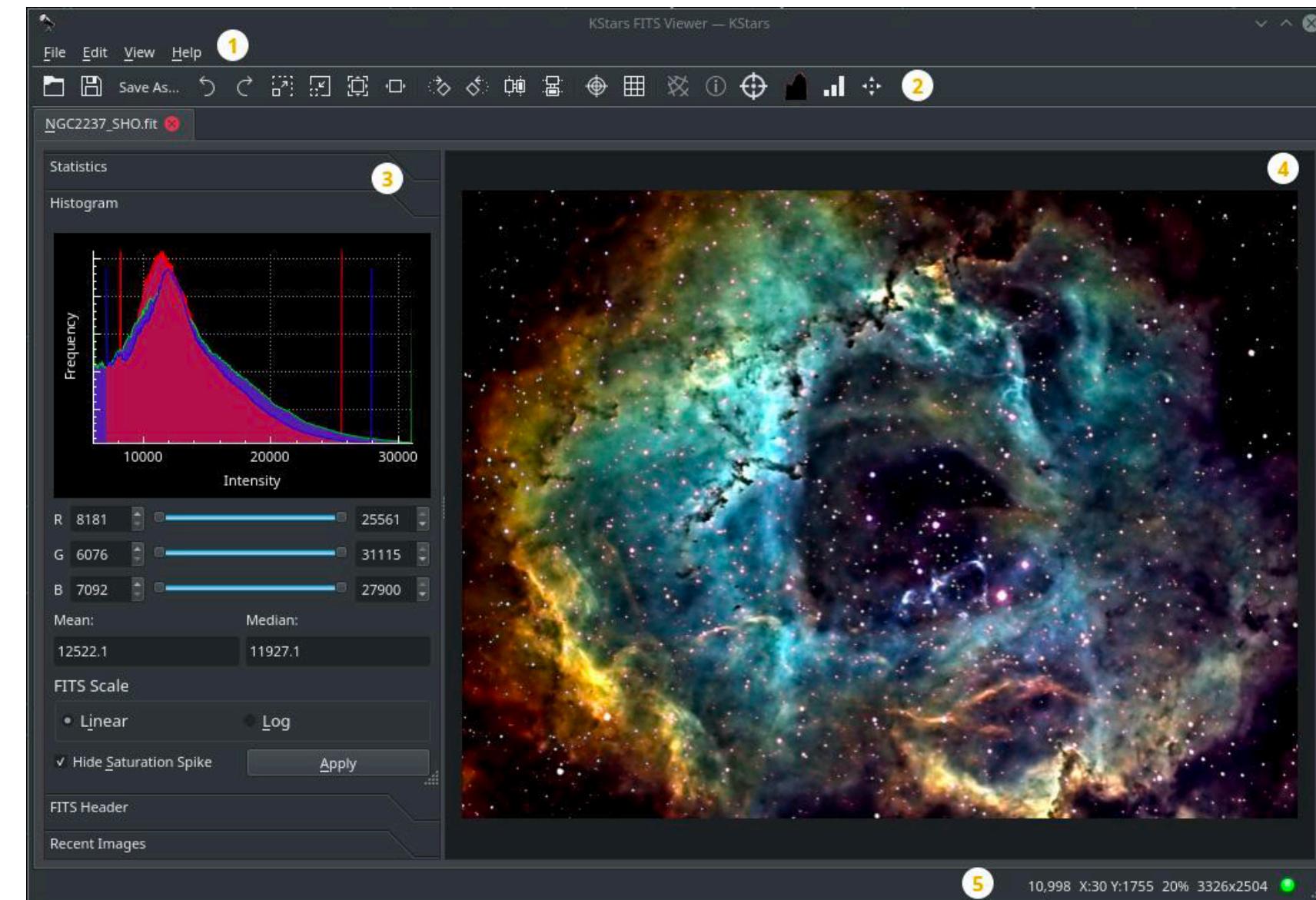
"No idea what data we have or where it is!"

No transparency to billions of files → "a black hole"

No possibility of filtering out relevant data from complex, heterogeneous datasets

Requirement:

1. Organize, catalog, and share research data effectively.
2. Improve collaboration, reproducibility, and data integrity.



FITS

(Flexible Image Transport System) was designed specifically for astronomical data, and includes provisions such as describing photometric and spatial calibration information, together with image origin metadata.

More info: fits.gsfc.nasa.gov

Results

Metadata-Hub created a centralized metadata warehouse.

- Enforcing data governance policies, and tracking data lineage.
- Streamlining data migration, integration, and compliance processes.
- Created custom extractors.

Enabled Metadata Mart:

- Extract embedded metadata from FITS file format (NASA.org).
- Reduced data prep in evaluation and search.
- Automated data provisioning to Jupyter Notebooks.
- Introduced “no-code/low-code” self-service data access.



CONTACT

David Cerf

Chief Data Evangelist

david@graudata.us

Phone: +1 972 896 3164

Nina Mangold

Channel Manager

nina.mangold@graudata.com

Phone: +49 7171 187-125

Learn more at:

www.Moremetadata.com

www.graudata.com



GRAU DATA

Your data \ Your control _